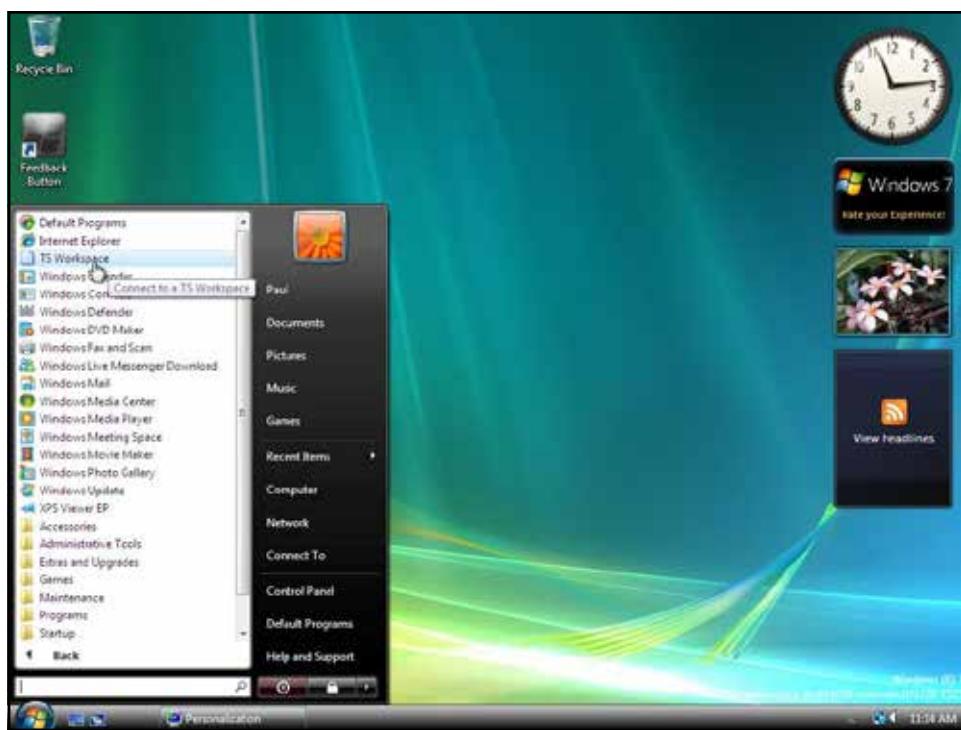


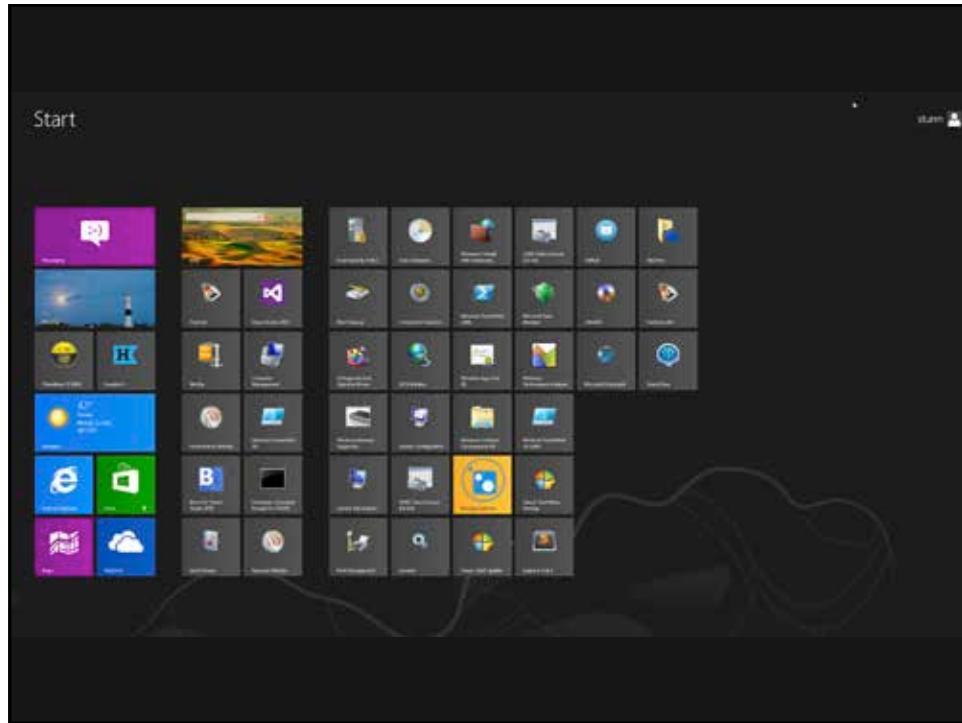
Organisatorisches

- Mittwochs, 10 – 11.30 Uhr
 - Grundlagenvorlesung
 - HS11
- Es gibt Übungen
- Klausur: Februar 2016
- Folienkopien und mehr
 - Asysob-Blog
 - studip



Bekannte Vertreter



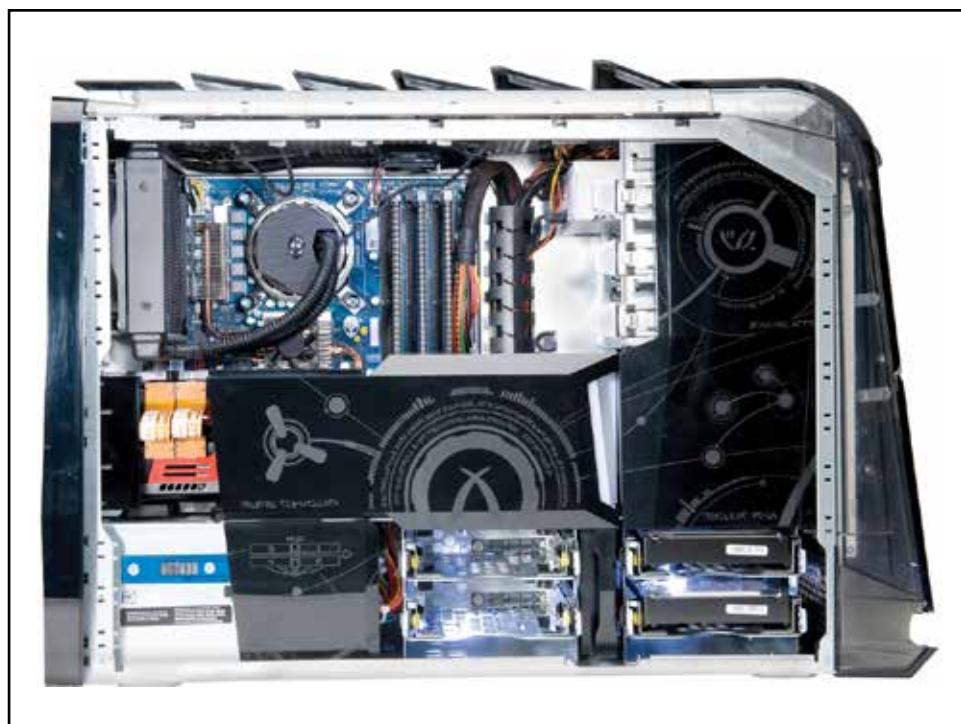






Hardware?

- Datenspeicherung (nicht-flüchtig, langsam)
 - Festplatten, SSD, CD/DVD/BD
- Datenspeicherung (flüchtig, schnell)
 - Hauptspeicher (RAM)
- Datenverarbeitung
 - Prozessor(en): CPU
 - Graphikkarte: GPU
 - Soundkarte usw.
- Datenkommunikation



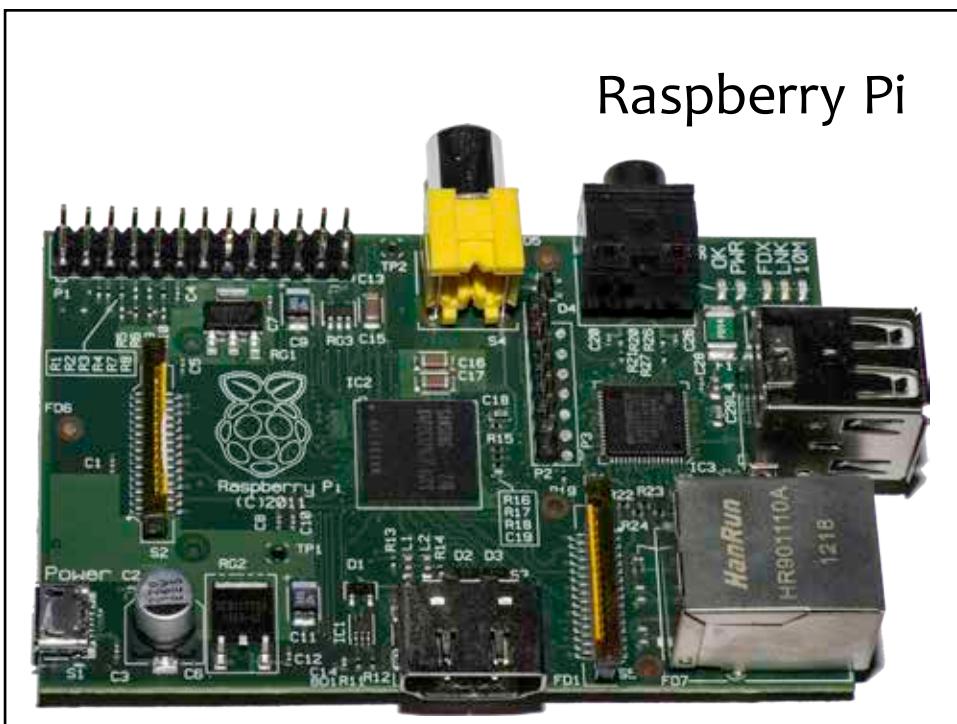


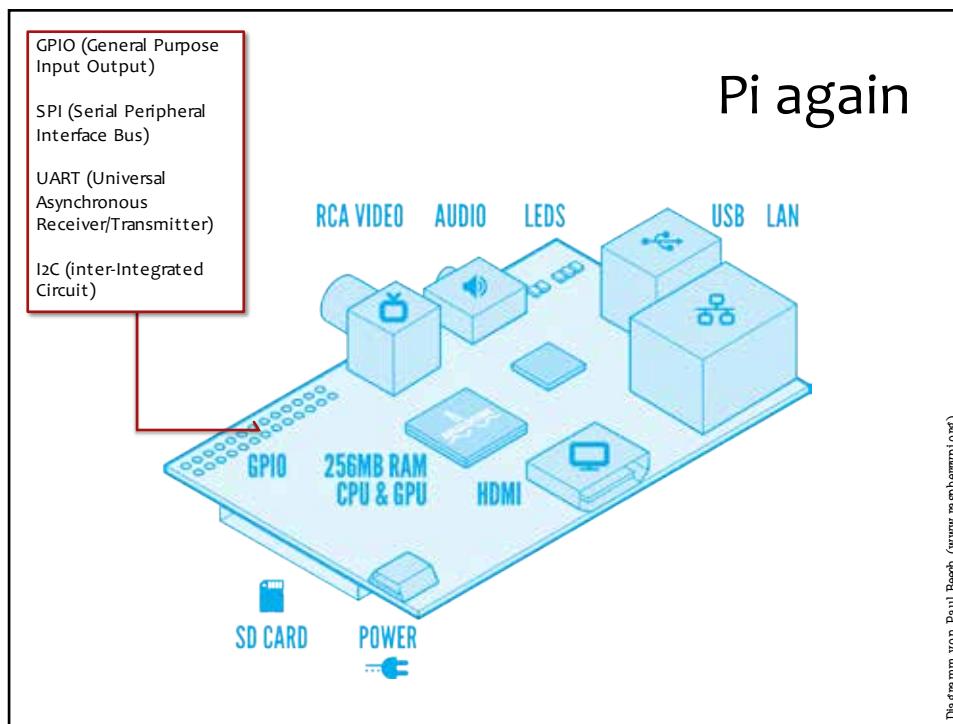
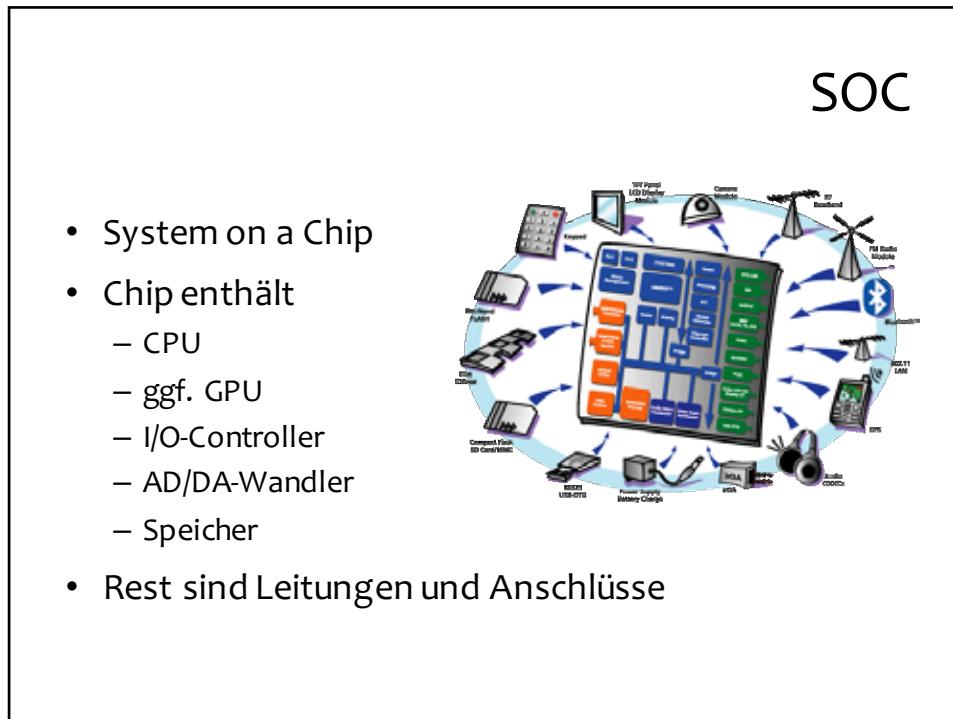


The screenshot shows a presentation slide titled 'INTEL® IA-32 ARCHITECTURES'. The slide contains three main sections:

- 2.1.17 2010 Intel® Core™ Processor Family (2010)**: Describes the Intel Core processor family based on the Westmere microarchitecture, using 32 nm process technology. It includes features like Intel® Scalable Memory Interface (Intel® SMIF) channels, Intel® 7500 Scalable Memory Buffer, Advanced RAS, and software recoverable machine check architecture.
- 2.1.18 The Intel® Xeon® Processor 5600 Series (2010)**: Details the Intel Xeon processor 5600 series, based on the Nehalem microarchitecture code name Westmere using 32 nm process technology. It supports up to six cores per physical processor package, up to 12 MB enhanced Intel® Smart Cache, and instruction sets AESNI, PCLMULQDQ, SSE4.2, and SSE4.1.
- 2.1.19 Second Generation Intel® Core™ Processor Family (2011)**: Describes the second-generation Intel Core processor family based on the Sandy Bridge microarchitecture code name. It includes Intel Turbo Boost Technology for Intel Core i5 and i7 processors, Intel Hyper-Threading Technology, Enhanced Intel Smart Cache and integrated memory controller, Processor graphics and built-in visual features like Intel® Quick Sync Video, Intel® Insider™, and instruction sets AVX, AESNI, PCLMULQDQ, SSE4.2, and SSE4.1.

At the bottom of the slide, there is a footer: '2.2 MORE ON SPECIFIC ADVANCES'.



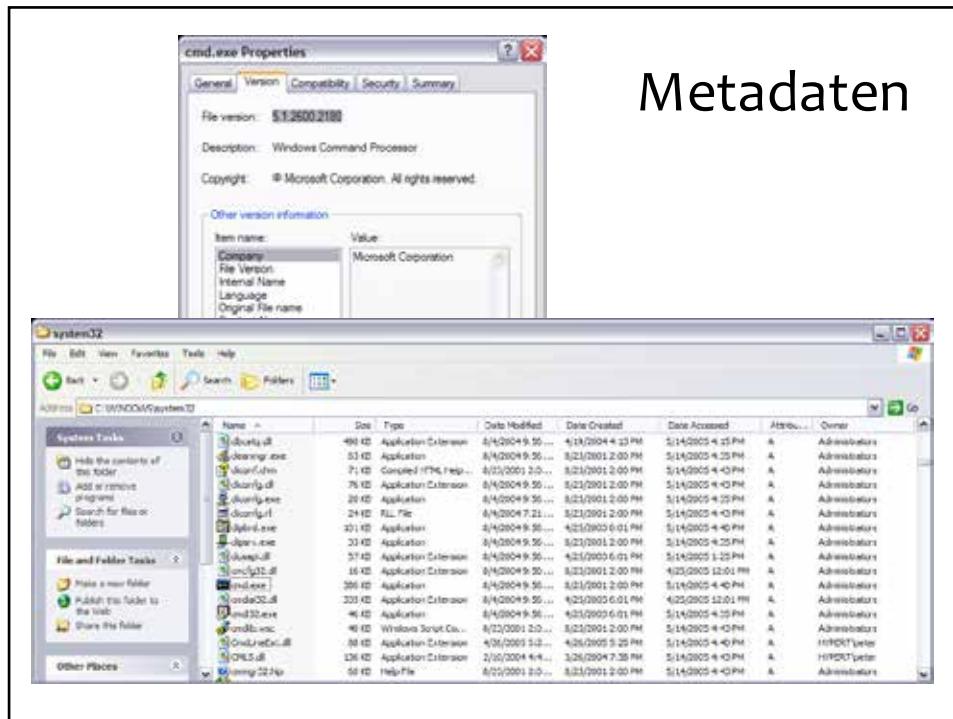




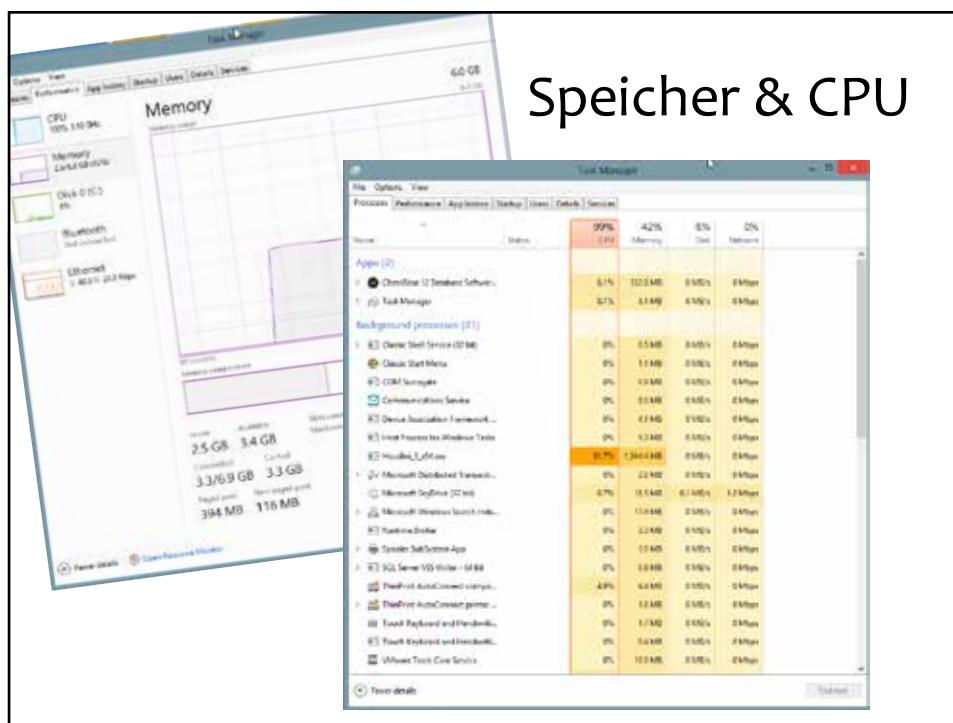
Dateien und Verzeichnisse



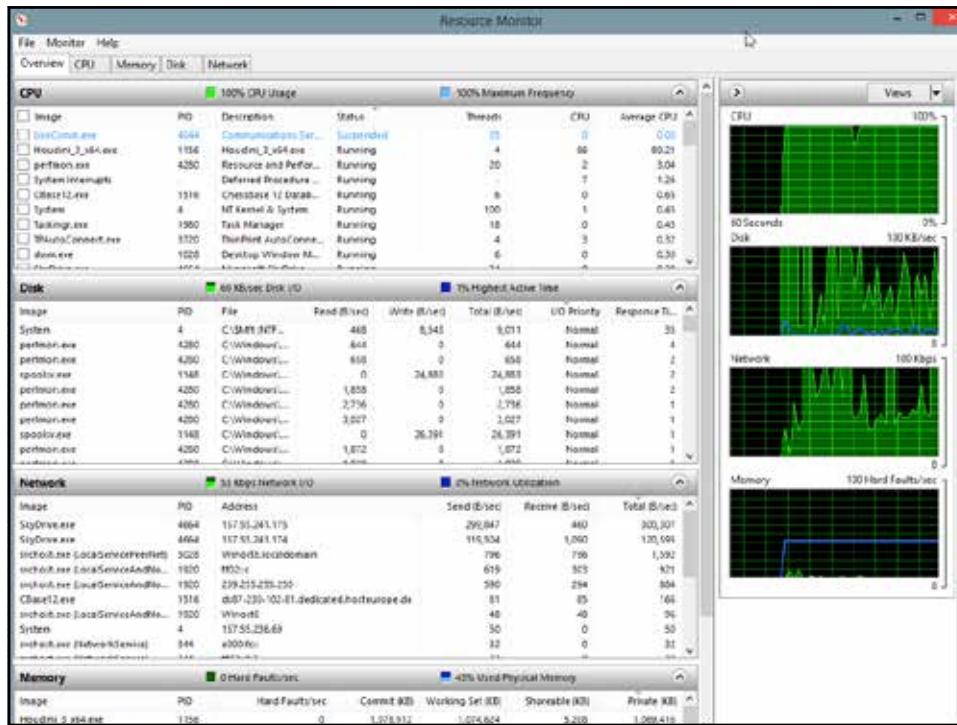
- Graphischer und textueller Zugang



Metadaten



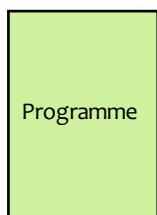
Speicher & CPU



Mediatorfunktion



User



Programme

- Vermitteln zwischen
 - Hardware
 - Programmen und Benutzern
- Ziele
 - Von technischen Details abstrahieren
 - Einfache Benutzbarkeit
 - Keine Monopolisierung
 - Vermeidung von Engpässen (Virtualisierung)
- Schnittstellen
 - User Interface (Shell)
 - Programmierung (API)

Hardware

CPUs
Bus
MMU
Speicher
Controller
Disks
Graphics
3D
Audio
Video
CD Writer
...
Netzwerk



Mission

- Develop a model that defines the functionality, semantics, and abstractions suitable for application developers
- Despite many physical limitations, the goal is to
 - provide unlimited memory (Virtual Memory)
 - provide an unlimited number of CPUs (Virtual processors)
 - provide means to cooperate and synchronize (Inter-Process Communication IPC)
 - provide persistent storage of data

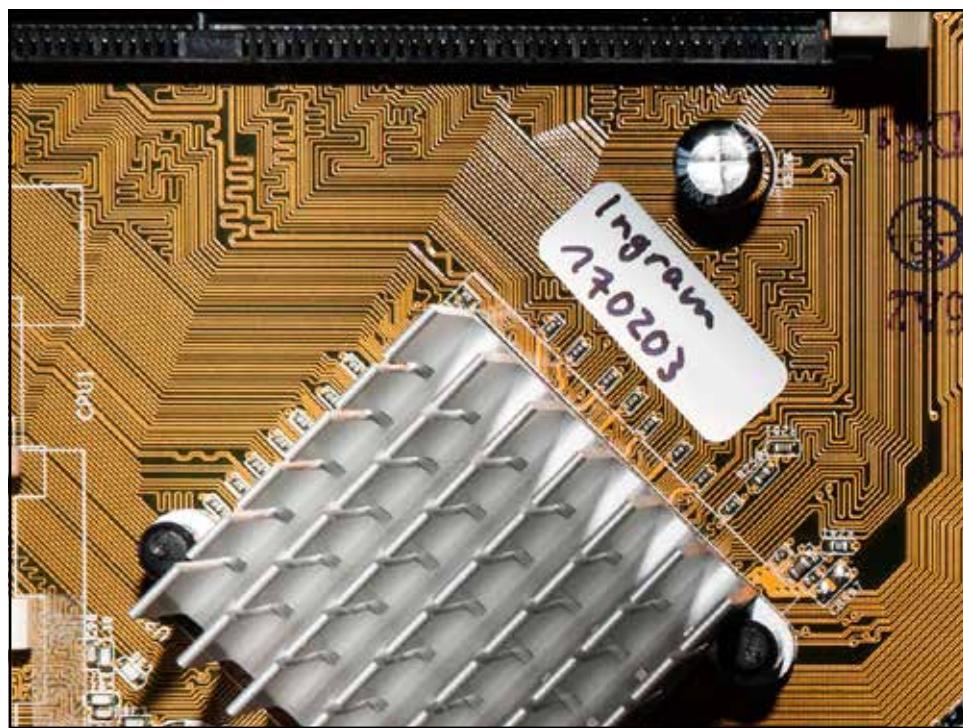
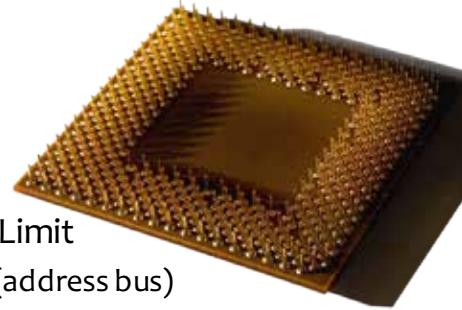
Das Kontinuum



- Aktivität (Zeit)
 - Abfolge ausgeführter Instruktionen
 - Kontrollfluss oder Thread
 - Unterstützung für mehrere Threads
 - Concurrency Control
- Raum
 - Menge der adressierbaren Speicherzellen
 - Adressraum oder "Address space"
 - Unterstützung für mehrere Adressräume

Adreßraum

- CPU definiert ein theoretisches, oberes Limit
 - Anzahl der Adreß-Pins (address bus)
- Beispiele
 - 32 bit Architektur
(Intel Pentium, AMD Athlon XP, ...)
 - 32 Adressleitungen = 4 GByte Adreßraum
 - 64 bit Architektur
(AMD Athlon 64, Intel Itanium, ...)
 - 64 Leitungen = 4194304 TByte theoretischer Adreßraum
 - Real limitiert auf ~40 Pins (Server 48 Pins)



Vorstellbar?

- 1 Adresse = 1 Smartie = 1 Gramm
- Wie schwer sind 2^{32} Smarties?
- Wie schwer sind 2^{64} Smarties?



Lösung

- Gewicht von 2^{32} Smarties
 - 4294 Tonnen
- Gewicht von 2^{64} Smarties
 - 18,446,744,073,709 Tonnen

13 Meter

22.1 Kilometer



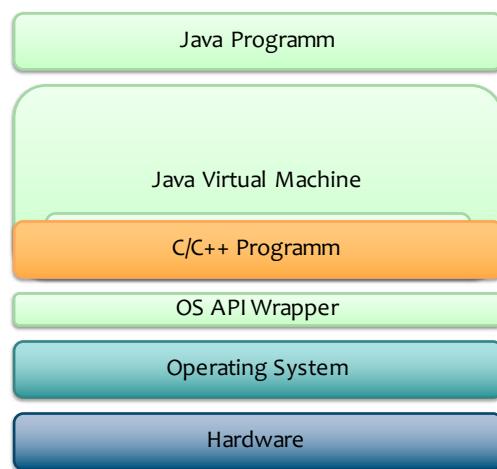
Was ist im Adressraum?



- Instruktionen (Maschinenbefehle)
- Programmdaten
 - Verschiedene Bereiche



Back to the roots

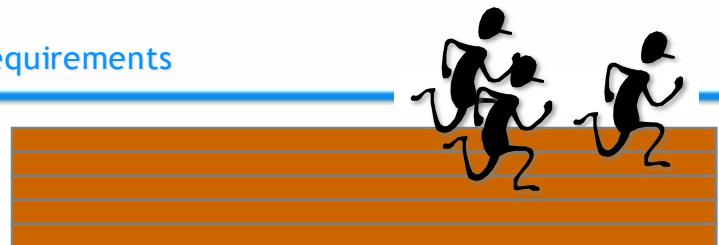


Operating Systems

Introduction
Space-Time Continuum

System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

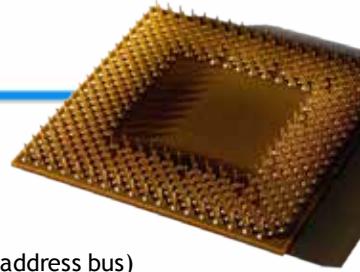
Basic Requirements



- Activity
 - Sequence of instructions executed in time = **Thread of control**
 - Support for multiple threads
 - Cooperation or concurrency among threads
- Space
 - Set of accessible memory cells = **Address space**
 - Support for multiple address spaces
 - Communication among shared memory

Operating Systems

Adress Space



- Theoretical limit defined by the CPU architecture
 - Specific number of address pins (address bus)
- Examples
 - 32 bit architecture (Intel Pentium, AMD Athlon XP, ...)
 - 32 address lines = 4 GByte address space
 - 64 bit architecture (AMD Athlon 64, Intel Itanium, ...)
 - 64 addresses = 4194304 TByte theoretical address space
 - Address bus of any 64 bit CPU limited to ~40 pins

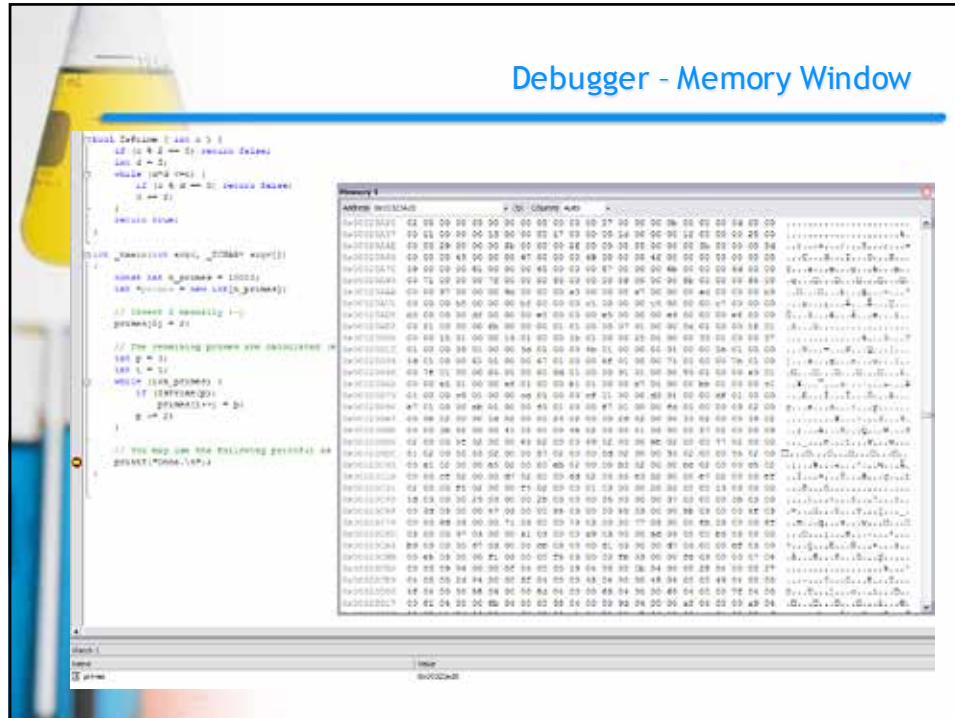
Operating Systems

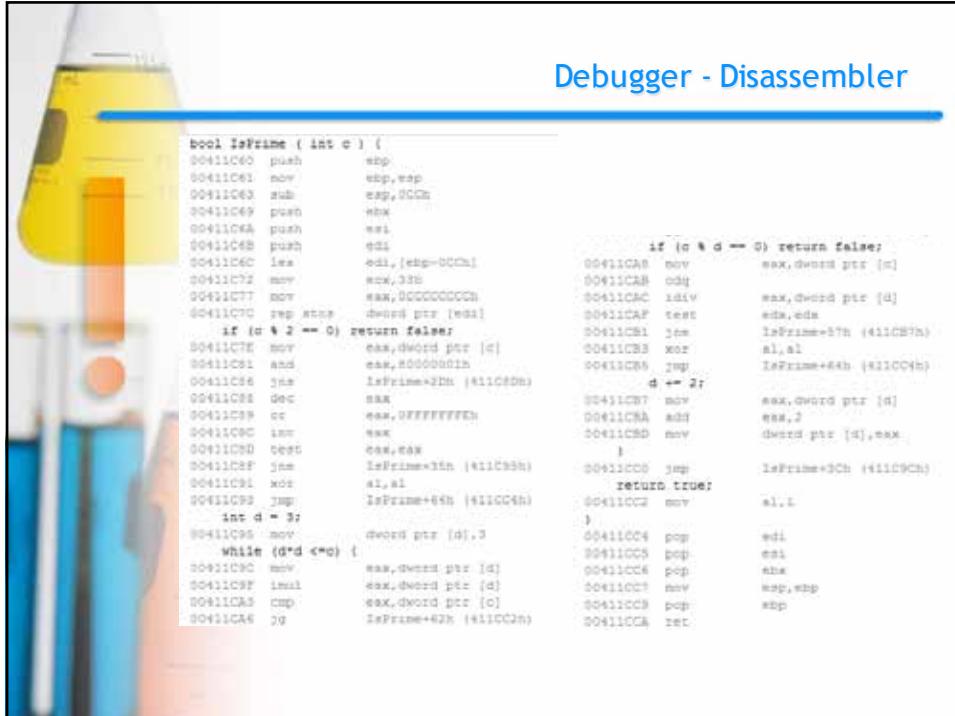
What's inside the Address Space?



- Program instructions
 - Machine instructions executed by CPU directly
 - Intermediate instructions executed by some virtual machine
 - Application code, runtime, libraries, operating system
- Program data
 - Static and dynamic data structures (the latter are so-called heaps)
 - program state reflecting the nested call of methods resp. procedures (stack)

Operating Systems





The screenshot shows a debugger interface with a yellow and blue background. The title bar says "Debugger - Disassembler". Below it is a code editor window displaying assembly code for the `IsPrime` function. The code is color-coded with various registers (eax, ebx, ecx, edx, etc.) highlighted in different colors. The assembly instructions are listed on the left, and their corresponding opcodes and comments are on the right.

```

bool IsPrime ( int c ) {
    00411C60 push    esp
    00411C61 mov     esp,esp
    00411C63 sub     esp,20Ch
    00411C69 push    ebx
    00411C6A push    esi
    00411C6B push    edi
    if (c % d == 0) return false;
    00411C6C lea     edi,[ebp-0CCh]
    00411C72 mov     eax,3Bh
    max,3Bh
    00411C77 mov     eax,0CCCCCCCCh
    00411C7C rep    stos    dword ptr [edi]
    if (c % 2 == 0) return false;
    00411C7E mov     eax,dword ptr [c]
    00411C81 and     eax,00000002h
    00411C86 jne     IsPrime+2Dh (411C8Dh)
    00411C88 dec    eax
    00411C89 or     eax,0FFFFFFFCh
    00411C90 imr    eax
    00411C91 test   eax,eax
    00411C92 jne     IsPrime+3Ch (411C93h)
    00411C93 not    al,al
    00411C94 jmp     IsPrime+6Ch (411C95h)
    int d = 3;
    00411C95 mov     dword ptr [d],3
    while (d*d <=c) {
        00411C96 mov     eax,dword ptr [d]
        00411C97 imul  eax,dword ptr [d]
        00411C98 cmp    eax,dword ptr [c]
        00411C99 jge     IsPrime+62h (411C9Ch)
        00411C9A add    d,d
        00411C9B mov     eax,dword ptr [d]
        00411C9C imul  eax,dword ptr [d]
        00411C9D mov     eax,dword ptr [c]
        00411C9E jge     IsPrime+62h (411C9Ch)
        00411C9F ret
}

```

Reference String

- From the “outside”, the execution of instructions is recognized as a series of memory references
 - References to instructions
 - References to data stored in main memory
 - References to data stored in stack
- Reference string = series of addresses the CPU is referencing while executing a program
- Reference string are huge
 - Caches must be disabled to observe any memory reference
 - With a modern CPU some billion references per second possible
 - Several GByte collected addresses per second
 - Dedicated hardware required to catch the string during program execution

Size of a reference string



- 350 Addresses per sheet of paper
- 10 sheets are approx. 1 mm high
- 2 billion references per second in a_loop example
- 8 billion refs in 4 seconds
= 2285 meter for the paper tower
- Eiffel tower = 312 meter



Reference Locality

- Most programs running on classical “von Neumann” architectures exhibit some “regular” access pattern to memory
 - Repeated updates to counting variables inside loops
 - Repeated execution of instructions forming a loop
 - Access to neighboring elements in an array
- Consequences
 - If some address a is referenced, there is a high probability that a will be referenced in the near future again
 - If some address a is referenced, there is a high probability that addresses around a will be referenced in the near future
- Many hardware and software mechanisms depend on this behavior
 - Primary reason that caches have hit rates close to 1
 - Virtual memory can be realized efficiently

Operating Systems

Reference Locality

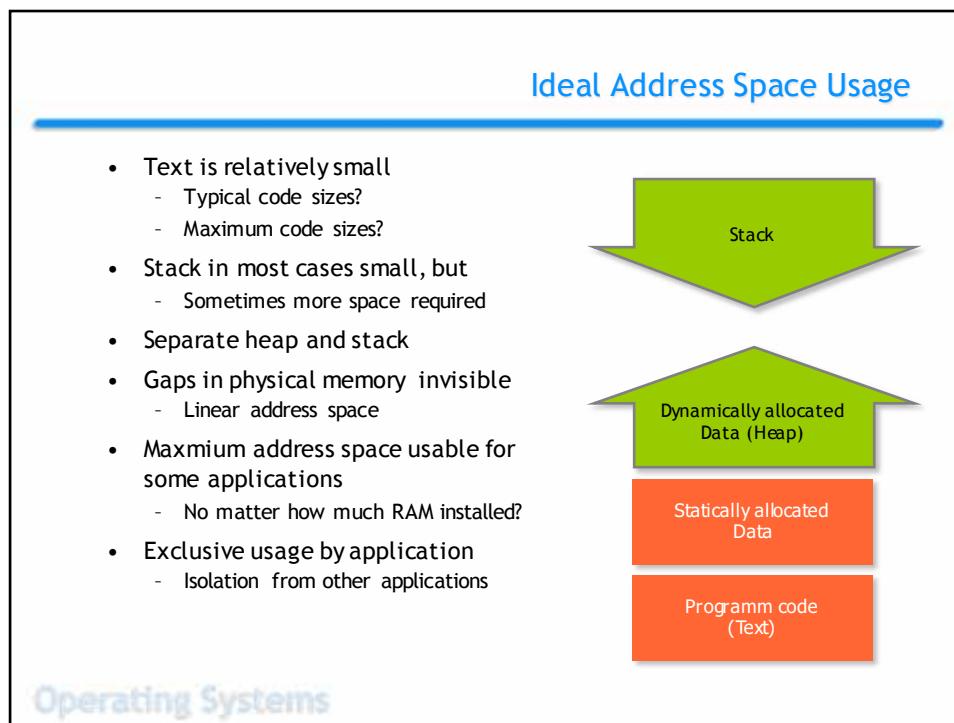
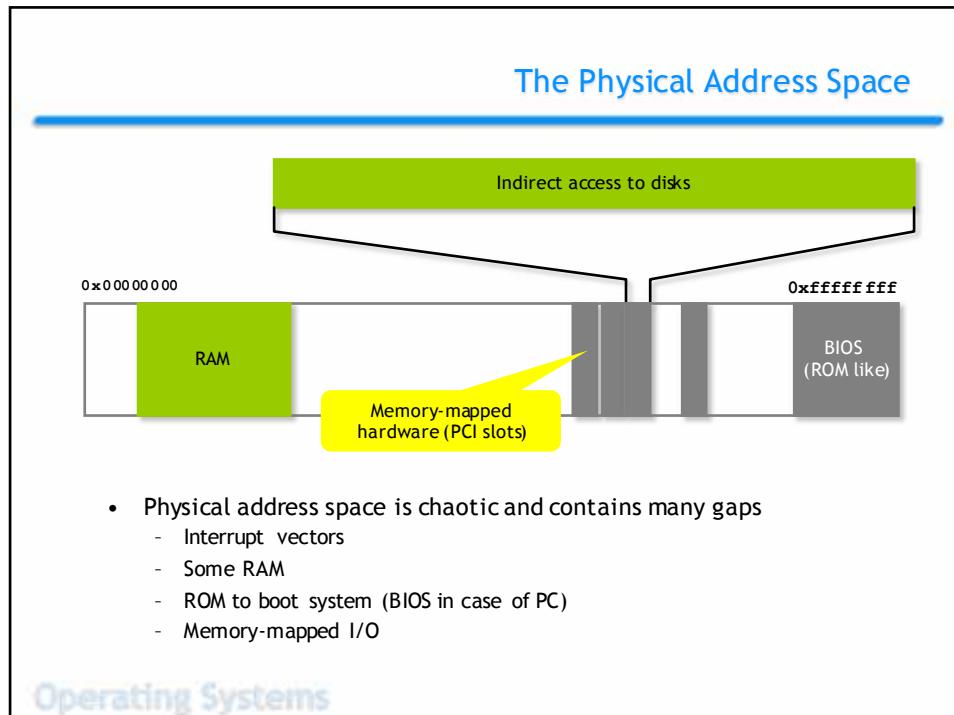
- Instructions are repeated
- Variables are referenced repeatedly
- If one address is referenced, nearby addresses will be referenced with high probability
- Reasons:
 - Sequential program execution
 - Procedure nesting is modest
 - Short loops
 - Frequent processing of arrays, list, and records

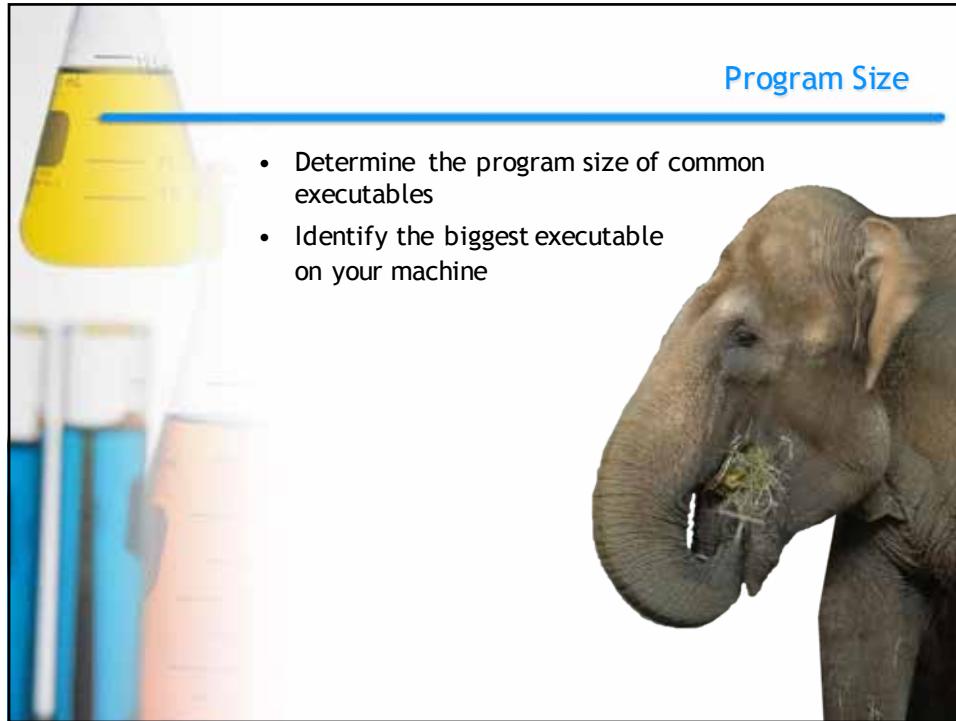
Operating Systems

The Ideal Translation Table

MEMORY

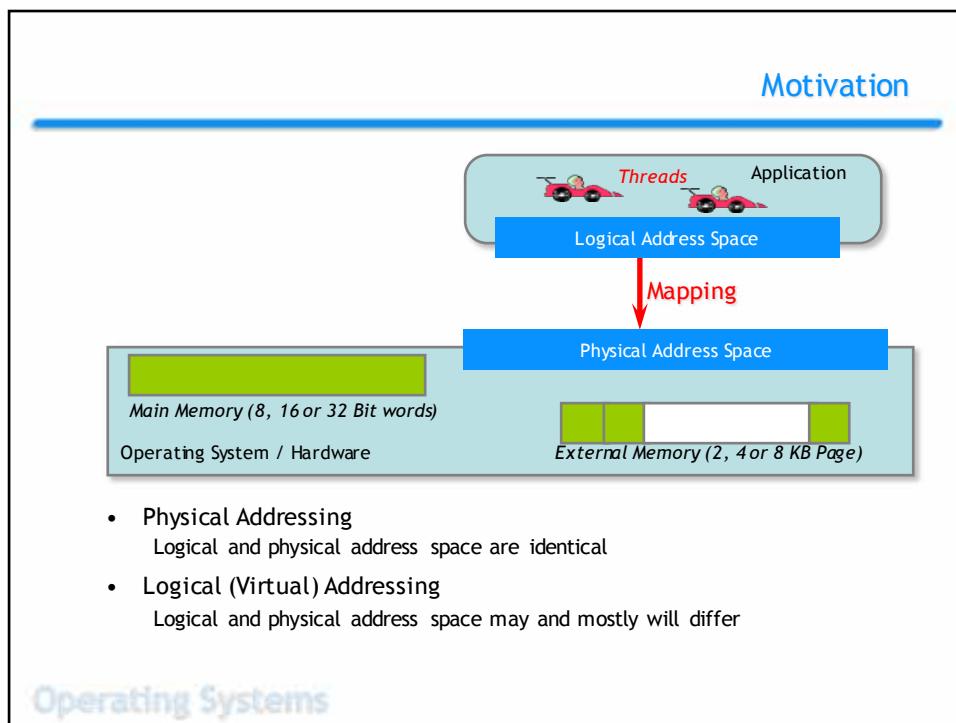
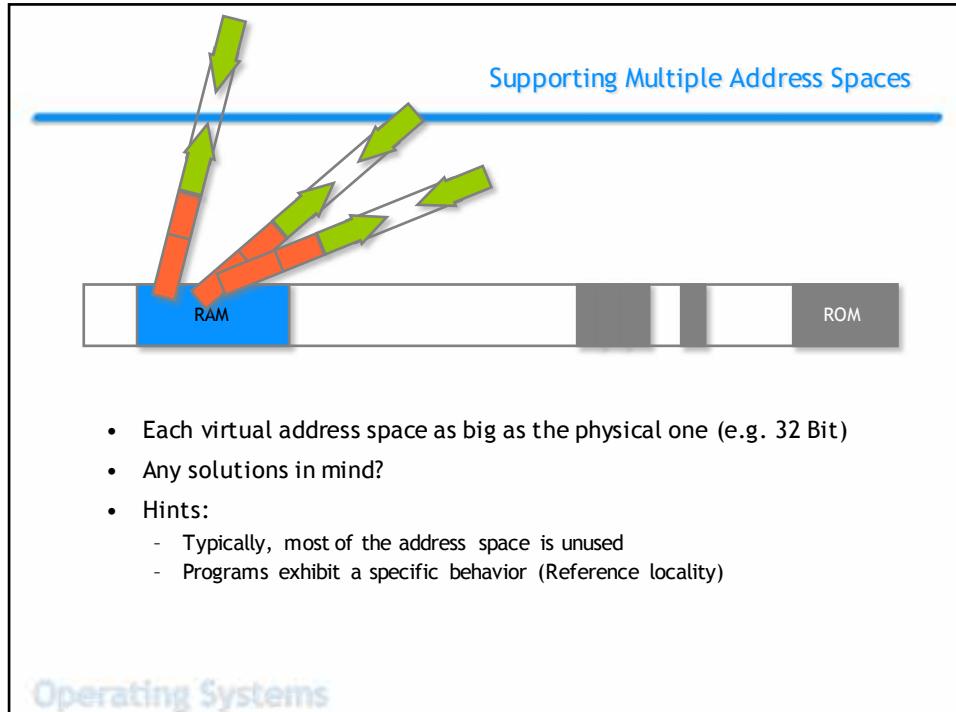
Operating Systems





Program Sizes

| • | Quake.exe | C:\Program... | 24,552 KB | Application | 3/15/2004 7:08 PM | 12/5/2... | A |
|---|-------------------------------|----------------|-----------|-------------|--------------------|-----------|---|
| • | ImageReady.exe | C:\Program... | 29,188 KB | Application | 10/15/2003 2:30 AM | 12/5/2... | A |
| • | AcroFWriter.exe | C:\Program... | 38,013 KB | Application | 3/14/2005 12:01 PM | 4/28/2... | A |
| • | photoshop.exe | C:\Program... | 17,532 KB | Application | 10/17/2003 1:59 AM | 12/5/2... | A |
| • | MSSpry.exe | C:\Program... | 17,460 KB | Application | 3/9/2005 4:55 AM | 3/9/20... | A |
| • | Urui.exe | C:\Program... | 15,267 KB | Application | 3/12/2005 9:56 AM | 3/12/2... | |
| • | nero.exe | C:\Program... | 14,889 KB | Application | 1/5/2005 2:47 PM | 2/19/2... | A |
| • | iPod Updater 2004-04-28.exe | C:\Program... | 14,288 KB | Application | 4/24/2004 5:20 PM | 4/24/2... | A |
| • | SA2_EAXDRIV_0310311.exe | C:\Program... | 13,896 KB | Application | 1/21/2004 9:32 PM | 1/21/2... | |
| • | Adobe20sd2sp1.exe | C:\Program... | 13,329 KB | Application | 8/29/2002 11:36 PM | 8/29/2... | |
| • | Adobe20sd3sp1.exe | C:\Program... | 13,033 KB | Application | 8/29/2002 11:18 PM | 8/29/2... | |
| • | ace.exe | C:\Windows... | 12,877 KB | Application | 3/21/2003 12:00 AM | 3/21/2... | |
| • | stylevision.exe | C:\Program... | 12,368 KB | Application | 3/9/2005 4:43 AM | 3/9/20... | |
| • | ace.exe | C:\Windows... | 11,972 KB | Application | 3/21/2003 12:00 AM | 3/21/2... | |
| • | RTWORD.EXE | C:\Program... | 11,766 KB | Application | 6/10/2003 10:29 PM | 6/10/2... | A |
| • | RTWORD.EXE | C:\Windows... | 11,756 KB | Application | 8/6/2003 2:24 PM | 8/6/20... | R |
| • | Illustrator.exe | C:\Program... | 11,473 KB | Application | 10/14/2003 6:23 AM | 12/5/2... | A |
| • | 51_76_vinva3k_english_vhd.exe | C:\Drivers\... | 11,324 KB | Application | 7/23/2004 2:47 PM | 7/23/2... | |
| • | Recode.exe | C:\Program... | 10,889 KB | Application | 12/8/2004 9:15 PM | 7/12/2... | A |
| • | 51_72_vinva3k_english_vhd.exe | C:\Drivers\... | 10,481 KB | Application | 7/10/2004 6:15 PM | 7/10/2... | |
| • | Acrobat.exe | C:\Program... | 9,981 KB | Application | 11/4/2003 12:58 AM | 11/4/2... | A |
| • | ImportLogon.exe | C:\Program... | 9,899 KB | Application | 4/25/2003 8:19 PM | 4/25/2... | A |
| • | Excel.exe | C:\Program... | 9,845 KB | Application | 5/19/2004 2:56 AM | 5/19/2... | A |
| • | Excel.exe | C:\Windows... | 9,839 KB | Application | 8/13/2003 3:34 AM | 8/13/2... | R |
| • | WordBuilder.exe | C:\Program... | 9,513 KB | Application | 8/30/2003 9:50 PM | 8/30/2... | A |
| • | CHSP021216.exe | C:\Program... | 9,463 KB | Application | 7/22/2003 4:31 PM | 7/22/2... | |
| • | 43_45_vin3xp_english.exe | C:\Drivers\... | 8,964 KB | Application | 4/25/2003 8:44 PM | 4/28/2... | |
| • | iTunes.exe | C:\Program... | 8,823 KB | Application | 5/4/2005 5:10 PM | 5/4/20... | A |
| • | patchdisk.exe | C:\Program... | 8,780 KB | Application | 1/15/2003 1:50 PM | 1/15/2... | |
| • | RunExplorer.exe | C:\Program... | 8,456 KB | Application | 3/12/2003 10:05 AM | 3/12/2... | |
| • | ChessProgram.exe | C:\Program... | 7,784 KB | Application | 3/5/2003 8:11 PM | 3/5/20... | A |
| • | Thunks.exe | C:\Program... | 7,689 KB | Application | 11/24/2003 7:39 PM | 2/6/20... | A |

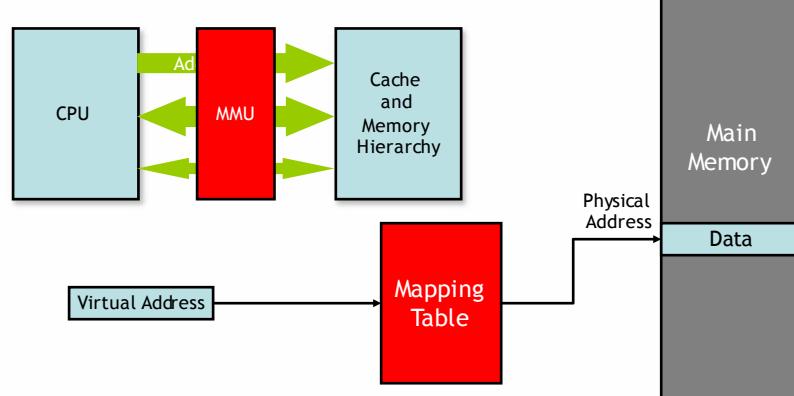


Goal

- Support for multiple address spaces
 - With multiple threads (multi-threaded programming)
- Internal protection (inside the address space)
 - Protection against stack overflows and other attacks
 - Protection against program flaws
- External protection (between address spaces)
 - Protect OS from malicious or faulty applications
 - Protect application from other applications
- Increase virtual memory beyond installed main memory
 - Integrate external persistent memory (disks)
 - Paging /Swaping
- Support for hardware-related servers
 - Memory-mapped I/O accessible in some virtual address spaces

Operating Systems**Basic Idea**

- Introduction of a mapping table f:
 - Physical Address = f (Virtual Address)
 - Memory management unit (MMU)

**Operating Systems**

Emerging Problems

- Size of table
 - Mapping 32 bit virtual address to 32 bit physical address:
 - 2^{32} entries in table required
 - 4 Byte for each entry (32 bit physical address)
- 16 Gbyte memory consumption
- Sparse array, so compression strategies are possible
- Performance
 - 1 additional memory access (index-based access)
 - Complex access algorithm in case of compressed sparse array
- Conclusion: Good idea, but too expensive!

$$2^{32} \cdot 4 = 2^{34}$$

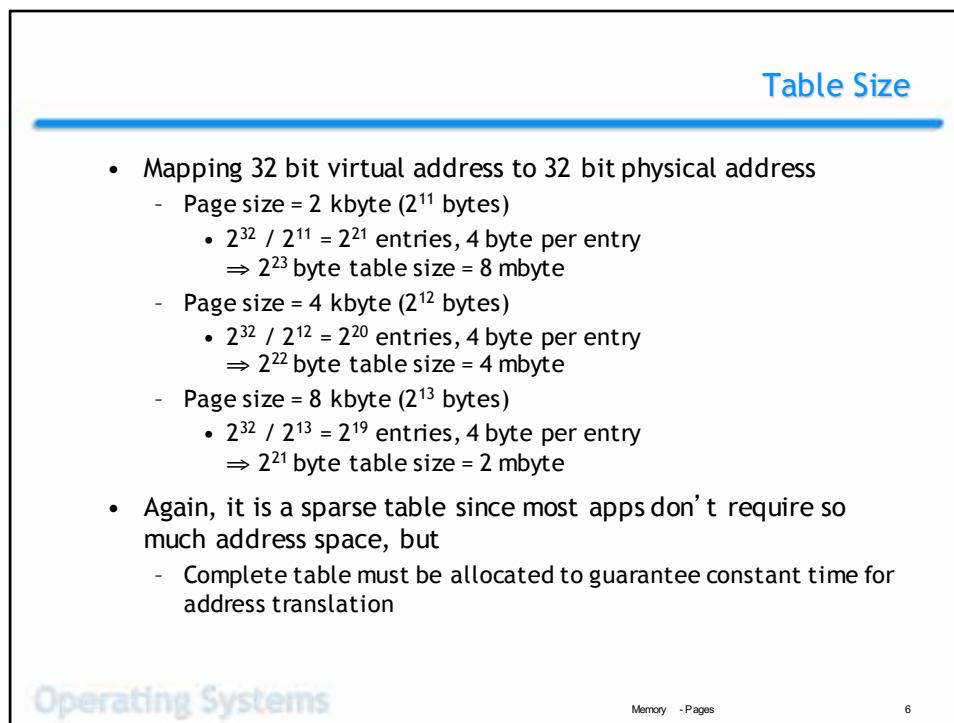
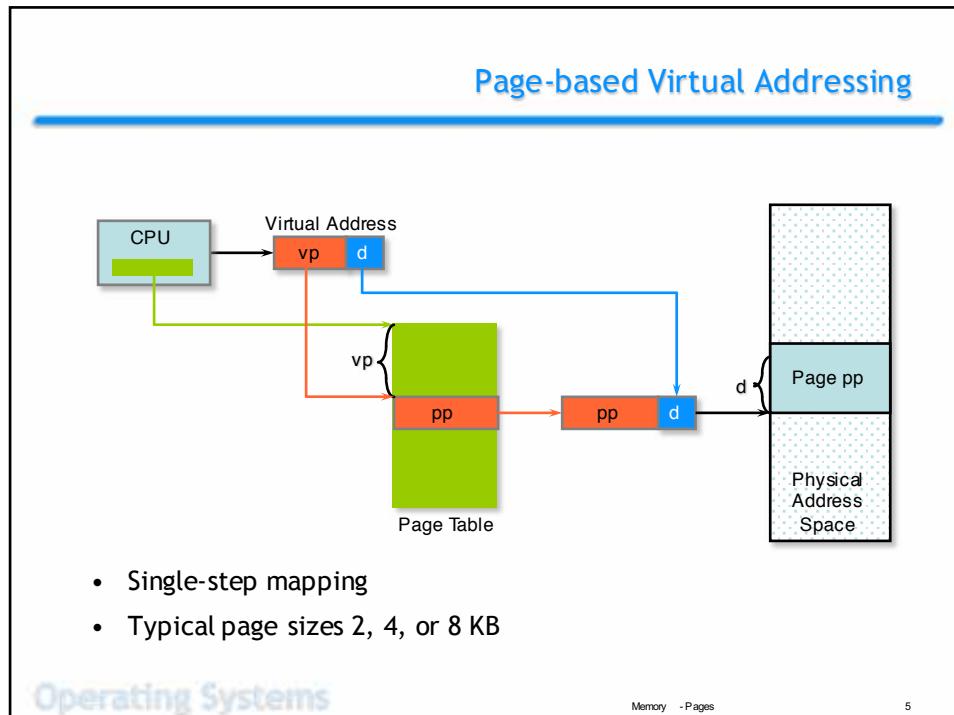
Operating Systems

Operating Systems

Memory
Pages

System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

-
- The diagram illustrates the MMU architecture. It shows the CPU sending a Virtual Address to the Memory Management Unit (MMU). The MMU then sends a Physical Address to the Main Memory. A green line connects the MMU to a Translation Table stored in Main Memory, indicating that the MMU uses a mapping table for address translation.
- MMU responsible for virtual address translation
 - Mapping table also stored in main memory
 - Size of table must be minimized
 - Each table entry is valid for a contiguous set of addresses called page



Memory Sizes of Applications

- Determine the memory consumption of common applications
- Hint
 - Start some applications and check the Windows task manager

Memory - Pages 7

Memory consumption for Selected Applications

| Image Name | PID | User Name | Session ID | CPU | Mem Usage | Page Flts. | Threads |
|-------------------|------|-----------------|------------|-----|-----------|------------|---------|
| PhotoShop.exe | 2058 | peter | 0 | 00 | 234,296 K | 29,558 | 7 |
| VersionCut.exe | 2672 | SYSTEM | 0 | 00 | 64,032 K | 18,123 | 28 |
| WPS Office.exe | 2158 | peter | 0 | 00 | 12,074 K | 26,350 | 8 |
| PaintShop Pro.exe | 2558 | peter | 0 | 00 | 32,195 K | 11,767 | 5 |
| surfheat.exe | 1576 | cverma | 0 | 00 | 30,152 K | 14,178 | 68 |
| POINTERITF.exe | 3716 | peter | 0 | 95 | 36,260 K | 233,776 | 9 |
| icloudsyncd.exe | 208 | peter | 0 | 00 | 24,632 K | 7,311 | 10 |
| spooler.exe | 1548 | SYSTEM | 0 | 00 | 18,692 K | 26,647 | 37 |
| IEFrame.dll | 2448 | SYSTEM | 0 | 00 | 1,072 K | 4,228 | 4 |
| regsvr32.exe | 2852 | peter | 0 | 00 | 18,948 K | 49,226 | 8 |
| win32k.exe | 2012 | SYSTEM | 0 | 00 | 13,240 K | 3,311 | 23 |
| larmmy.exe | 2036 | SYSTEM | 0 | 00 | 10,900 K | 229,940 | 37 |
| explorer.exe | 2972 | peter | 0 | 00 | 10,380 K | 359,471 | 17 |
| dbgrid.exe | 2464 | peter | 0 | 00 | 0,988 K | 2,849 | 19 |
| andreas.exe | 3029 | SYSTEM | 0 | 95 | 8,294 K | 2,188 | 36 |
| hostess.exe | 2558 | peter | 0 | 00 | 3,044 K | 1,371 | 4 |
| trac32.exe | 2208 | peter | 0 | 00 | 7,735 K | 2,361 | 2 |
| pcm32.exe | 2208 | peter | 0 | 95 | 7,932 K | 2,122 | 4 |
| extinct.exe | 1988 | MR/WORK SERVICE | 0 | 00 | 7,776 K | 1,945 | 35 |
| CRIS.exe | 372 | SYSTEM | 0 | 00 | 7,184 K | 12,558 | 11 |
| TurnerAptex.exe | 2294 | peter | 0 | 95 | 7,056 K | 1,768 | 4 |
| appletalk.dll | 2672 | SYSTEM | 0 | 00 | 6,988 K | 20,589 | 6 |
| AjaxPub.exe | 2691 | peter | 0 | 00 | 6,260 K | 4,201 | 2 |
| CTHLPWR.EXE | 2072 | peter | 0 | 00 | 6,000 K | 1,596 | 8 |
| hostpath.exe | 2444 | peter | 0 | 00 | 5,848 K | 8,834 | 2 |
| MinPaint.exe | 2558 | peter | 0 | 00 | 5,752 K | 1,331 | 2 |
| CTStylit.exe | 2064 | peter | 0 | 00 | 5,702 K | 1,303 | 2 |
| msasn1.dll | 2472 | peter | 0 | 00 | 5,672 K | 6,791 | 12 |
| shobj.dll | 1240 | LOCAL SERVICE | 0 | 00 | 5,128 K | 1,238 | 14 |
| services.exe | 340 | SYSTEM | 0 | 00 | 5,028 K | 7,504 | 15 |
| KAV.exe | 2240 | peter | 0 | 00 | 4,735 K | 9,343 | 13 |
| OUTLOOK.EXE | 354 | peter | 0 | 00 | 4,724 K | 34,393 | 17 |

Memory - Pages 8

Another Reduction in Size

- Most applications require only a fraction of the possible address space
- With the exception of Photoshop - and of course Photoshop is different ☺ - most applications require less than 64 mbyte
 - 64 mbyte is only 1/64 of 4 gbyte
- Observation
 - Several contiguous memory areas are used by most apps
 - Code
 - Static data
 - Dynamic data = heap
 - Stack
- Possible improvement
 - Introduction of a multi-stage, tree-like mapping table

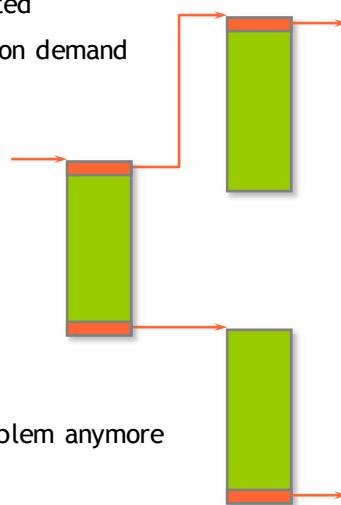
Operating Systems

Memory - Pages

9

Gain

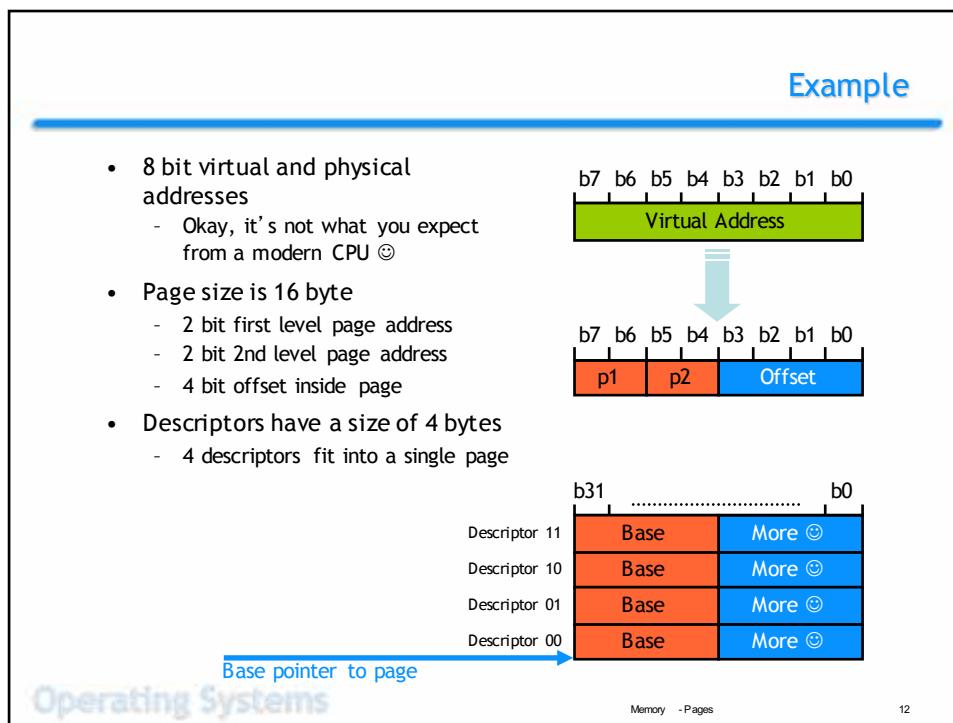
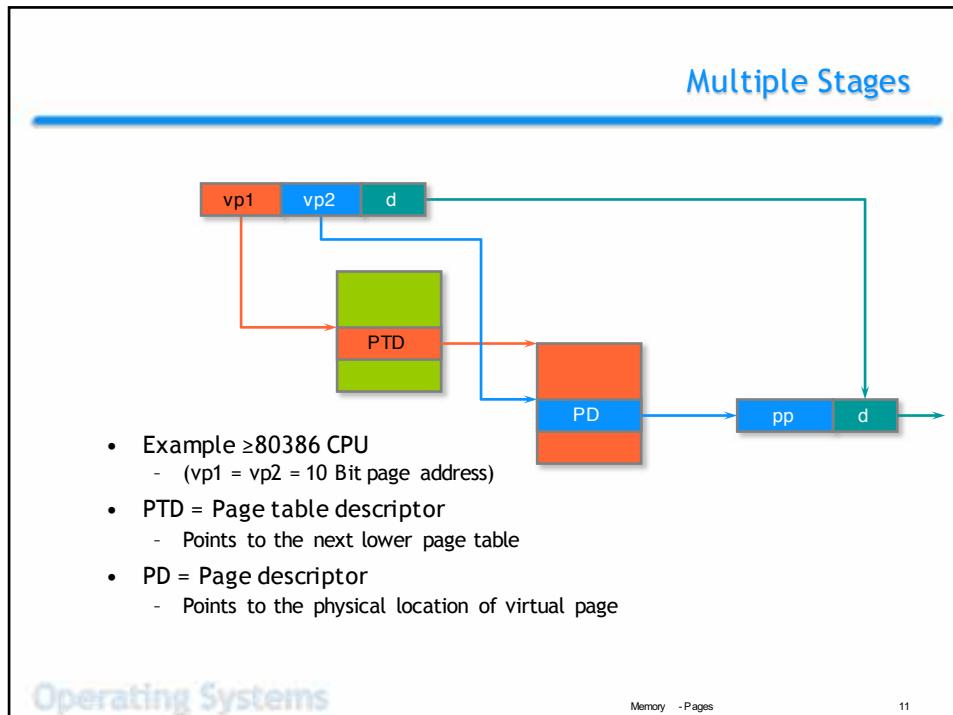
- Root page table MUST be allocated
- Page tables of later stages only on demand
 - Example: First and last byte in address space used:
 - $vp1=vp2=10$
 - $3 * 2^{10} * 4 \text{ Byte} = 12 \text{ Kbyte}$
 - Above config also works with
 - $1024 * 4 \text{ KB} = 4 \text{ MB}$ for text, data, and heap
 - $1024 * 4 \text{ KB} = 4 \text{ MB}$ stack
- 12 kbyte are sufficient for many small executables
- Conclusion: Size of table no problem anymore

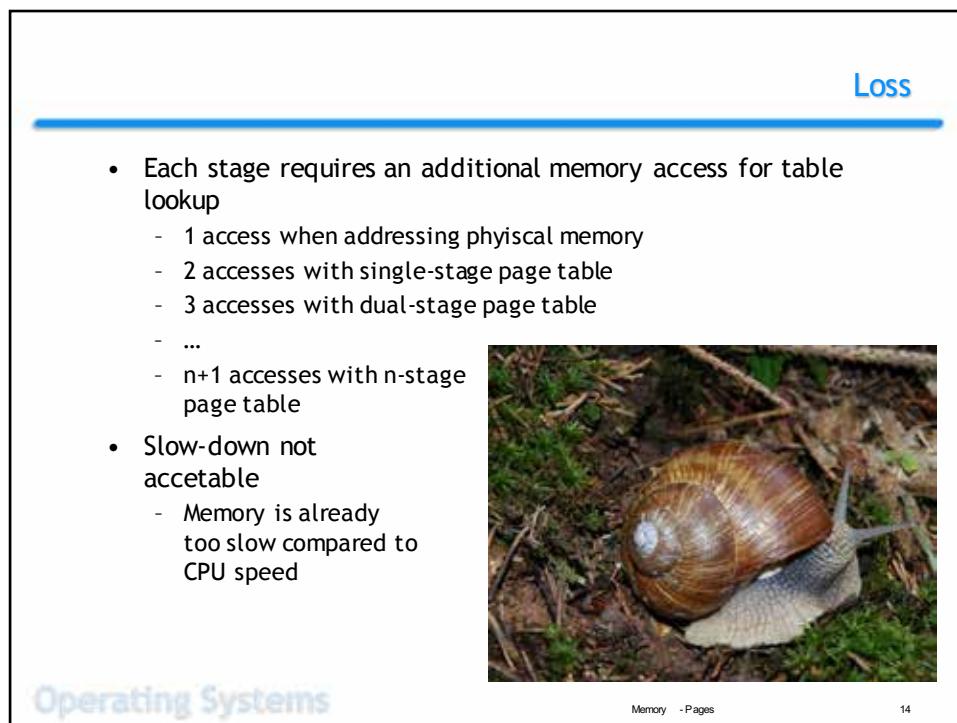
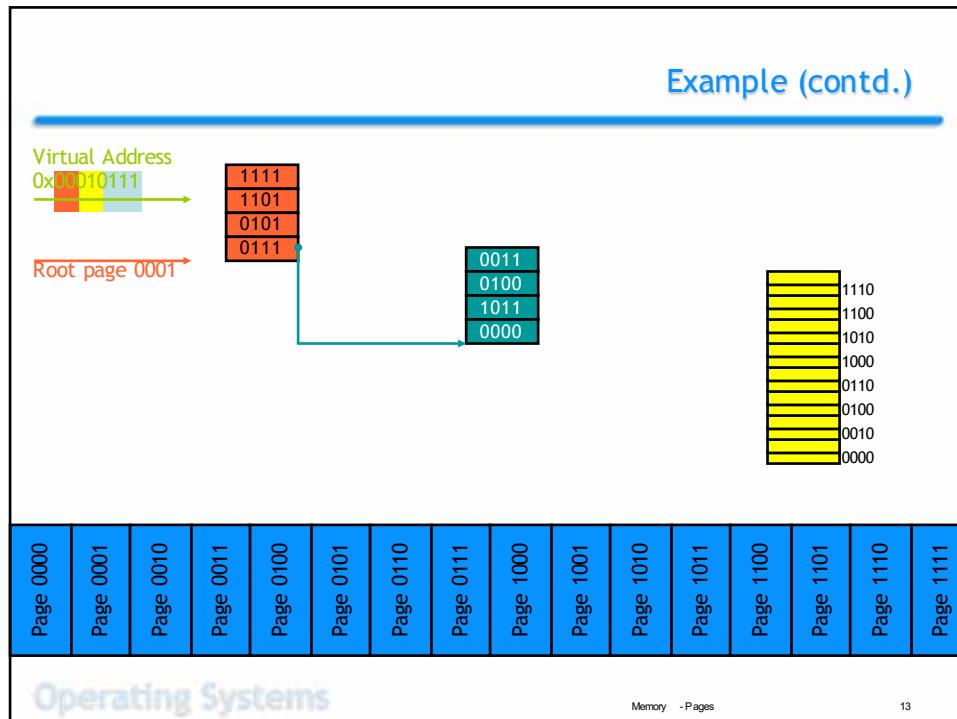


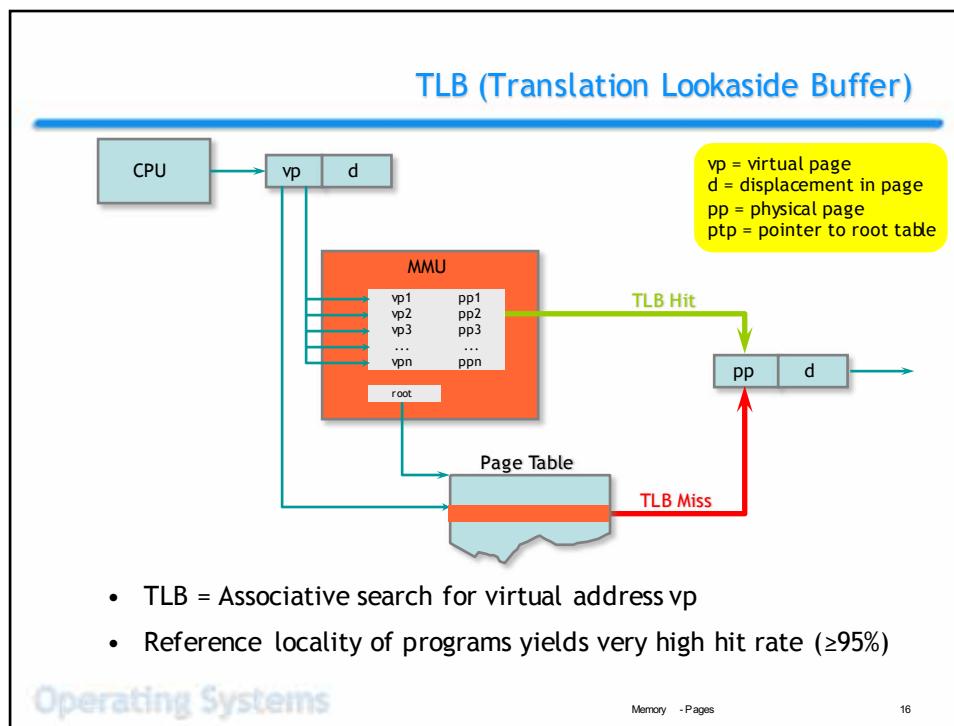
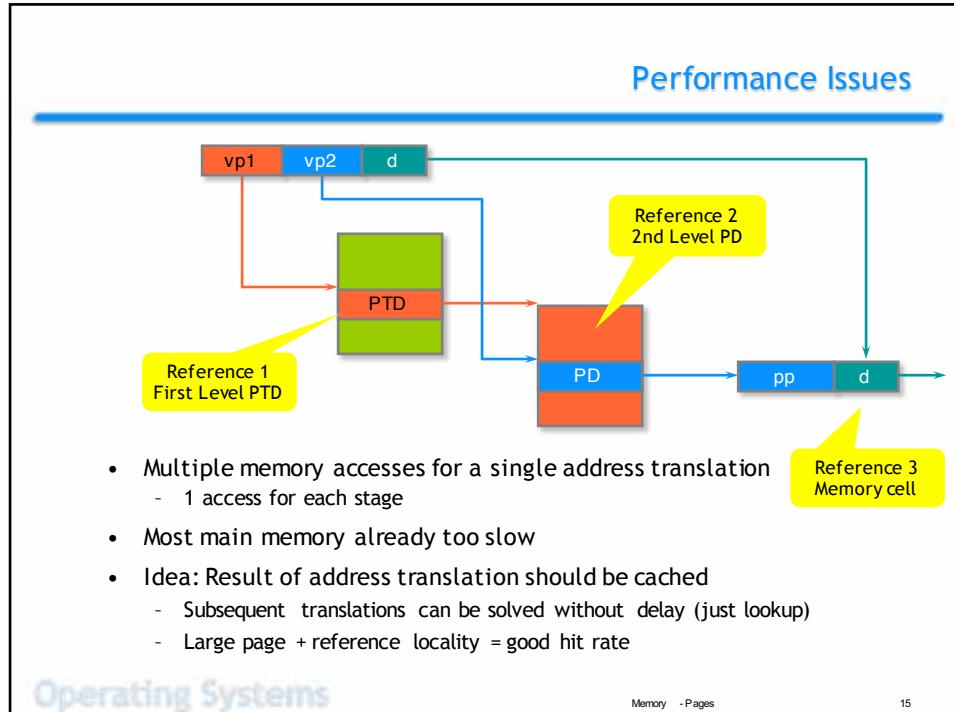
Operating Systems

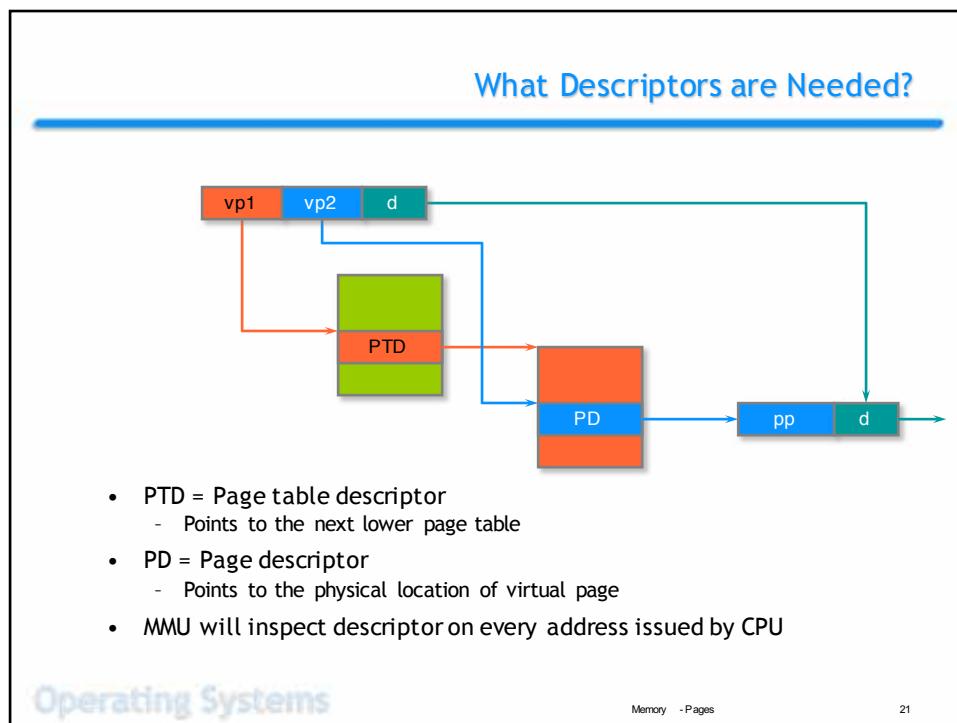
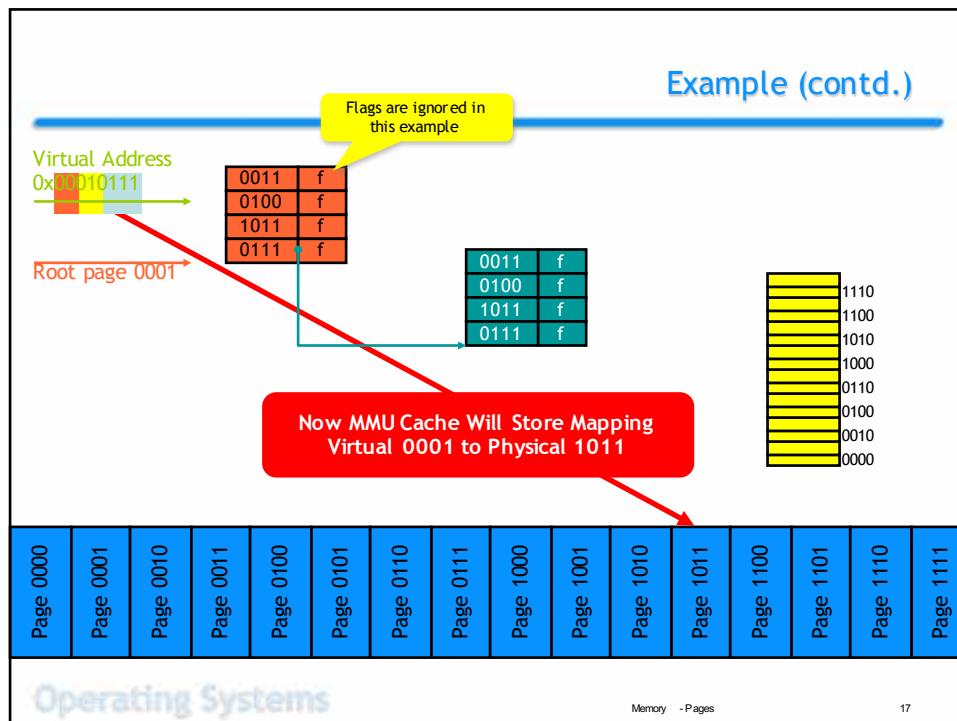
Memory - Pages

10









What Must/Can be Stored in Descriptor?

- Must be stored
 - Page location
 - Base address of virtual page in physical memory
 - Type flag to distinguish page descriptors from page table descriptors (could also be defined implicitly by translation stage)
- Can be stored
 - Access permissions
 - Read, Write, Execute for user-mode access
 - Read, Write, Execute for privileged access
 - Statistics
 - Has some address inside page been referenced?
 - Has some address inside page been modified?
 - Page location again
 - Page is currently not in main memory but on disk

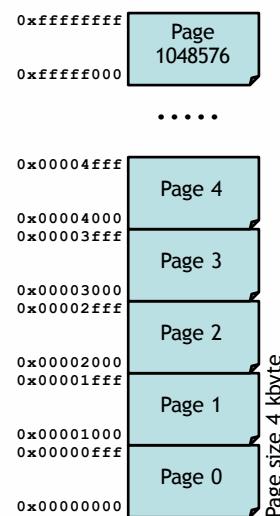
Operating Systems

Memory - Pages

22

Size of Descriptor

- Address of page in physical memory is 32 bit, but
 - by limiting the position of a page to multiples of page size k , the lower k bits are always 0 and can be ignored
- Assuming a page size of 4 kbyte, 12 bit are free in each descriptor
 - Space for permissions, statistics, etc.
 - Additional bits can be used to implement advanced memory techniques
- 32 bit descriptor size sufficient with all extras ☺



Operating Systems

Memory - Pages

23

Page Table Descriptor

PTD: A horizontal bar divided into four colored segments: grey (T), green (P), orange (Base of Page Table), and light blue (Free).

- T-Bit required to identify descriptor type:
 - Page table descriptor (PTD)
 - Page descriptor (PD)
- P-Bit:
 - Present bit:
 - P=1: Page accessible in main memory
 - Base address points to main memory
 - P=0: Page stored on external memory (e.g. disk)
 - Base address identifies external memory block
 - Some bits in descriptor free for OS usage

Operating Systems Memory - Pages 24

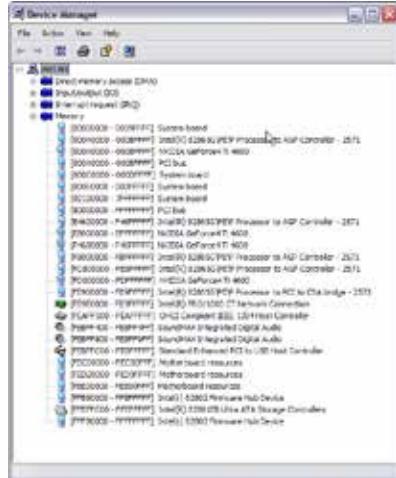
Page Descriptor

PD: A horizontal bar divided into nine colored segments: grey (T), green (P), blue (R), blue (D), red (C), red (Access), orange (Base Address of Page), and light blue (Free).

- Access control bits
 - Define the allowed access modes (read/write/exec, user/supervisor)
 - Denying any access sometimes useful
 - Access violation leads to interrupt and OS intervention
- C-Bit (Cache Disable Bit)
- R-Bit (Referenced Bit)
 - R=1 \Leftrightarrow Page has been referenced since last reset of bit
 - Must be reset explicitly by OS
- D-Bit (Dirty Bit)
 - D=1 \Leftrightarrow Page has been written since last reset of bit (\Rightarrow R=1)
 - Must be reset explicitly by OS
- Again a few bits freely available to OS

Operating Systems Memory - Pages 25

Cache Disable Bit



- Some virtual pages map to physical pages with memory-mapped IO
 - Make PCI memory areas available to privileged services
 - Access to command and status registers for some peripherals
- Changes to these memory positions happen as side-effects
- Caches shouldn't store the content they've seen the last time
 - Otherwise caches would signal a hit after the first access (hit 'a' once on the keyboard, means 'a' forever ☺)

Operating Systems

Memory - Pages 26

Conclusions

- Virtual memory can be implemented efficiently
 - Space efficiency
 - tree of page tables
 - Time efficiency
 - TLB cache in MMU
- The advantages outweigh the space and time requirements by far
 - Homogeneous address space
 - Internal and external protection
- Many technical details still to be solved
 - Modification in some page table requires partial or sometimes complete flush of TLB

Operating Systems

Memory - Pages 31

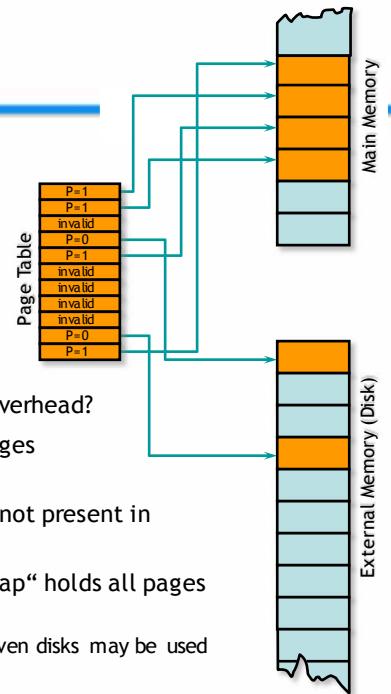
Operating Systems

Memory
Paging: Pages are present or not

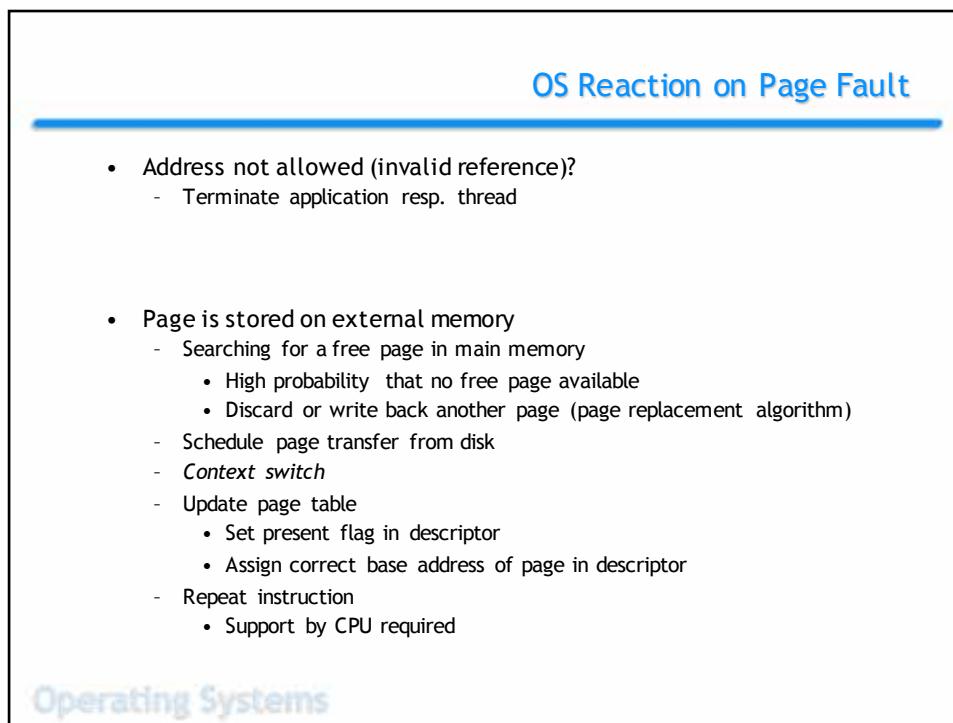
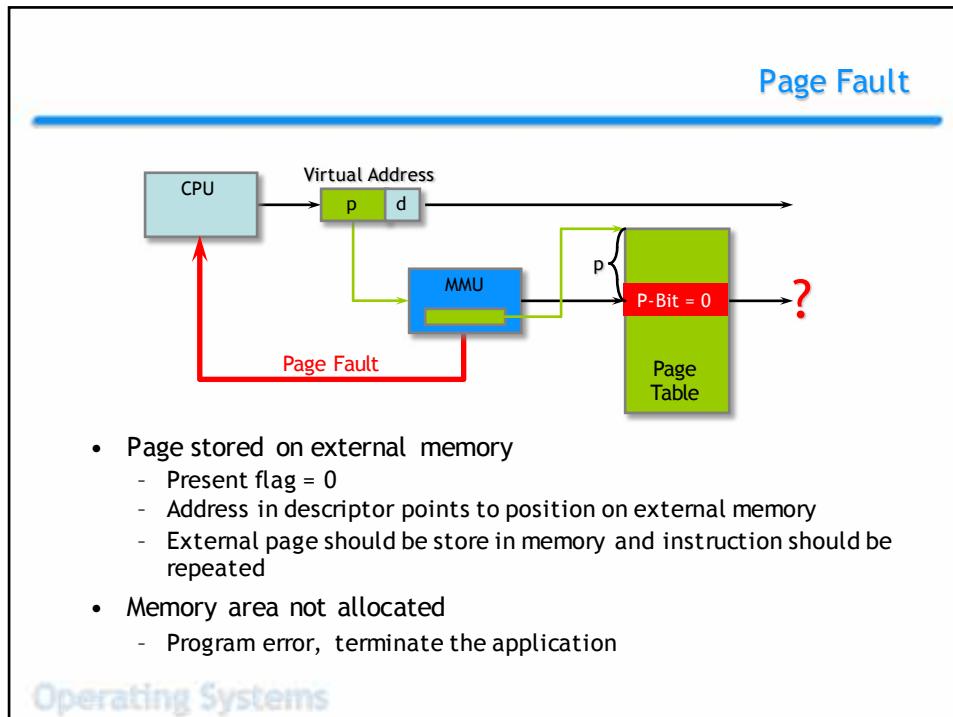
System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

Paging

- Advantages
 - Increasing virtual address space
 - Virtual address space limited by capacity of disk only
 - Enough space to support multiple address spaces in parallel
- Reference locality guarantees low overhead?
- Reading and writing of individual pages
 - Present flag in MMU descriptor
- Page fault in case of access to page not present in memory
- Special file called „pagefile“ or „swap“ holds all pages not in main memory
 - In some OS dedicated partitions or even disks may be used



Operating Systems



Efficiency of Page Replacement

- Effective access time depends on the probability of a page fault

$$t_{\text{effective}} = (1 - p) * t_{\text{Memory}} + p * t_{\text{PageFault}}$$

- p = Page fault probability
- t_{Memory} = Access time main memory (few ns)
- $t_{\text{PageFault}}$ = Time to handle the page fault

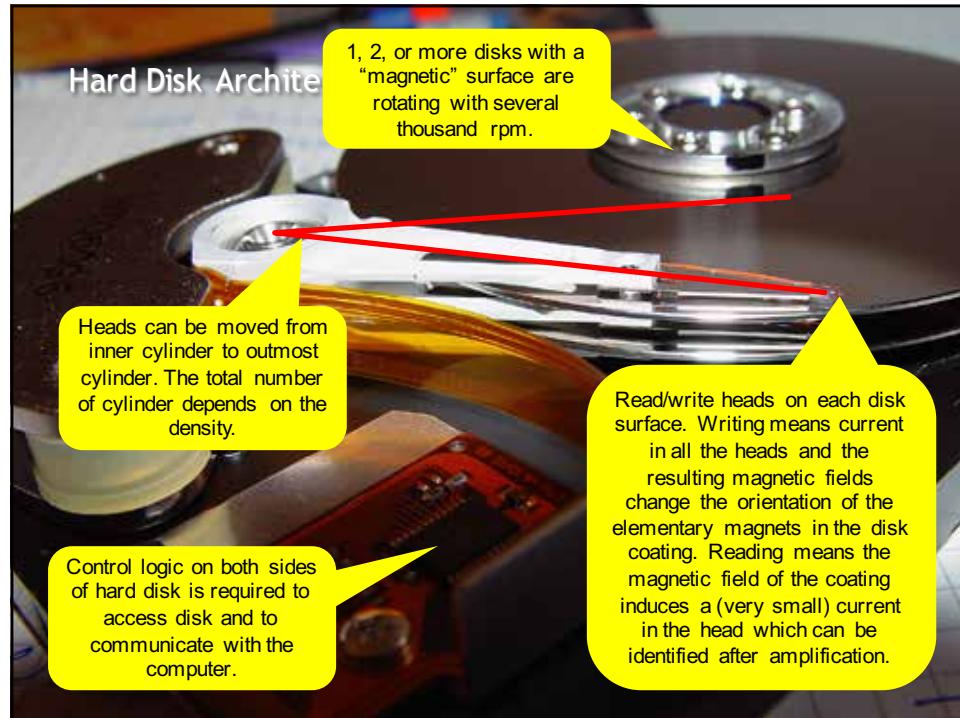
- $t_{\text{PageFault}}$

Interrupt latency and handling Search for free page in main memo Write back candidate if dirty Read requested page from disk Repeat instruction

$$t_{\text{Page Fault}} = t_{\text{Interrupt}} + t_{\text{Search}} + t_{\text{Write}} + t_{\text{Read}} + t_{\text{Instr}}$$

$$\approx 2 * t_{\text{Read/Write}}$$

Operating Systems



Example

- Assumption
 - Memory access time: 100 nsec
 - Disk access time: 10 msec

$$\begin{aligned}t_{\text{effective}} &= (1-p)*10^{-7} + 2p*10^{-2} \\&\approx 2p*10^{-2}\end{aligned}$$

- Required: Less than 10% overhead

$$1.1*10^{-7} > 10^{-7} + 2p*10^{-2}$$

$$p < 5*10^{-7}$$

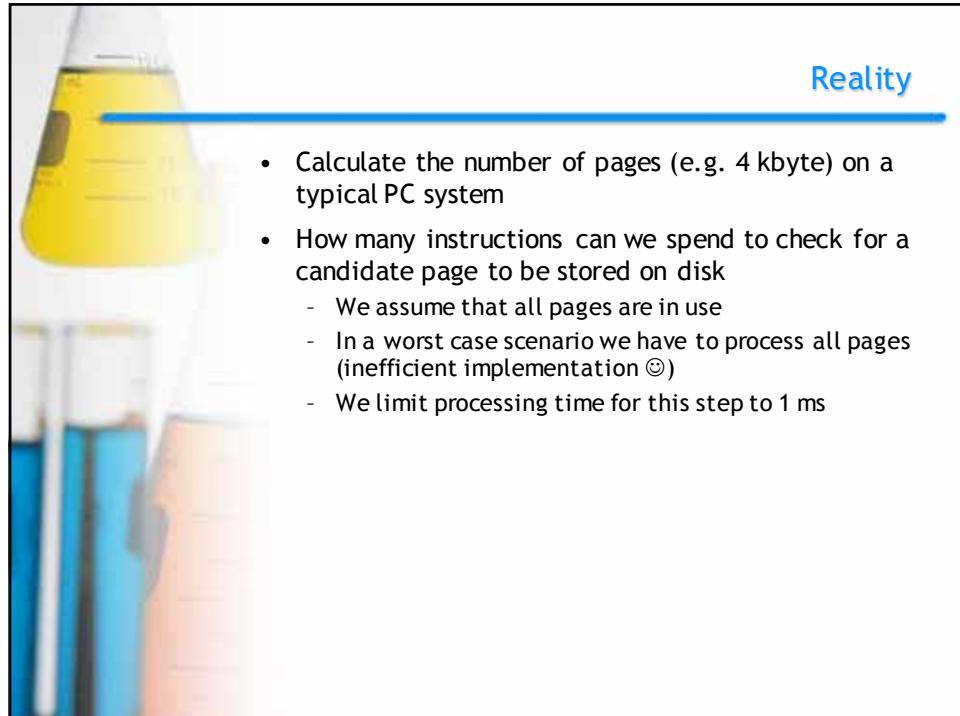
- At most 1 page fault for each 2 million instructions

Operating Systems

Time of Page Fault

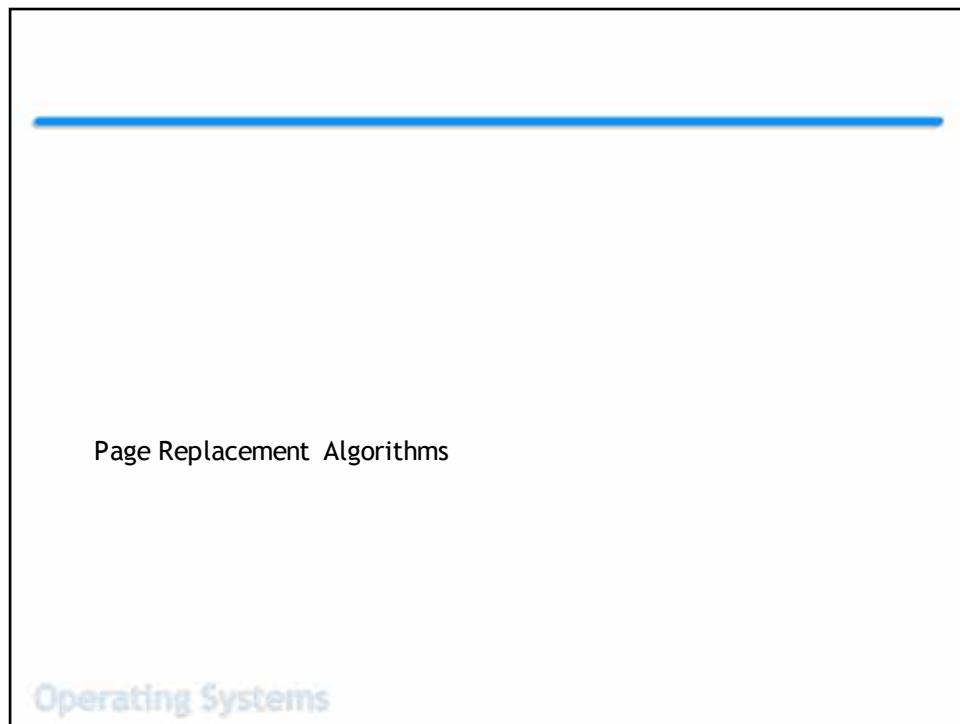
- Page fault handling dominated by at most 2 accesses to disk
 - Access time for disks in the order of 3-5 ms
- How can we improve the effective access time
 - Keep page faults at minimum
 - Always have the required pages in main memory ☺
 - Not very much possible with respect to interrupt handling
 - Faster CPU, better OS, but the gain very limited
 - Find the candidate page to leave memory very fast
 - Crucial since we have a lot of pages
 - Read and write page
 - Buy bigger and/or faster disk
(bigger also means faster in most cases, since data are stored more densely)
 - No substantial improvement possible when repeating instructions

Operating Systems



Reality

- Calculate the number of pages (e.g. 4 kbyte) on a typical PC system
- How many instructions can we spend to check for a candidate page to be stored on disk
 - We assume that all pages are in use
 - In a worst case scenario we have to process all pages (inefficient implementation ☺)
 - We limit processing time for this step to 1 ms

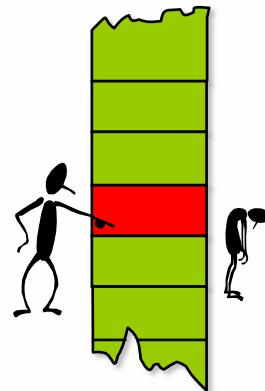


Page Replacement Algorithms

Operating Systems

Page Replacement Algorithm

- Assumption
No free page available when page fault occurs
- Possible consequences
 - Terminate application
 - Swap another application completely
 - Replace a memory page used otherwise
- But which page to replace?
 - Lots of pages available
 - 512 MB main memory, page size 4 kb
 $\Rightarrow 131072$ pages
 - Goal is to minimize the number of future page faults



Operating Systems

Reference Strings

- Used to evaluate resp. benchmark replacement algorithms
- Calculated from the sequence of memory references executed by the CPU
 - Hard to obtain! Why?
 - Only page information of address required
 - Only first reference to page represented
 - Subsequent references inside same page can be omitted
- Applying reference string to a replacement algorithm with a fixed set of memory pages
 - Count the number of page faults

Operating Systems



Create a Reference String

- Write a function that creates a reference string of given size with a parameterized reference locality
- You may inspect a running version of a program that simulates [page replacement algorithms](#)



Possible Solution

```
public ReferenceString ( int size, int n_pages, int alpha, int window, int seed )
{
    Random rgen = new Random(seed);
    refstring = new int[size];
    for (int i=0; i<size; i++)
    {
        do
        {
            bool select_old = (rgen.Next(1000) < alpha) && (i>1);
            if (select_old)
            {
                int act_window = (i>window) ? window : i;
                int old = rgen.Next(act_window);
                refstring[i] = refstring[i-old-1];
            }
            else
                refstring[i] = rgen.Next(n_pages);
        }
        while ((i>0) && (refstring[i-1] == refstring[i]));
    }
    current = 0;
}
```

- **alpha** = probability (between 0 and 1000) to select a previously referenced page from a window of last referenced pages with size **window**

Optimal Replacement by Belady

- Replace the page which will not be used for the longest time in future
 - Quite hard to implement
 - Used to compare with other strategies
- Why is Belady optimal?



Operating Systems

Implement Belady

- Extend your program creating a reference string with an implementation of Belady
 - By using the reference string, we know the future
- Inspect the complete program

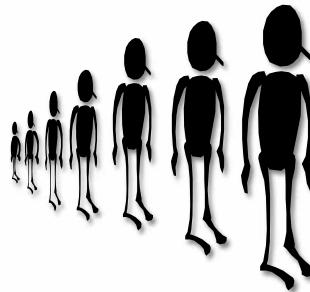


| Belady Example | | | | | | |
|--------------------------|-----------------|--|----------------------------------|--|--|--|
| Number of reference | Referenced page | Virtual pages currently in main memory (-1 means free) | What happens? | | | |
| 0 | -1 | -1 | page fault, replace empty page 0 | | | |
| 1 | ref= 5 | 1 -1 | page fault, replace empty page 1 | | | |
| 2 | ref= 1 | 1 5 -1 | page found at 0 | | | |
| 3 | ref= 5 | 1 5 5 -1 | page found at 0 | | | |
| 4 | ref= 1 | 1 5 5 2 -1 | page found at 0 | | | |
| 5 | ref= 5 | 1 5 5 2 1 -1 | page found at 1 | | | |
| 6 | ref= 0 | 1 5 5 2 1 -1 | page fault, replace empty page 2 | | | |
| 7 | ref= 1 | 1 5 5 0 -1 | page found at 0 | | | |
| 8 | ref= 3 | 1 5 5 0 3 -1 | page fault, replace empty page 3 | | | |
| 9 | ref= 5 | 1 5 5 0 3 -1 | page found at 1 | | | |
| 10 | ref= 1 | 1 5 5 0 3 -1 | page found at 0 | | | |
| 11 | ref= 5 | 1 5 5 0 3 -1 | page found at 1 | | | |
| 12 | ref= 6 | 1 5 5 0 6 -1 | page fault, replace used page 3 | | | |
| 13 | ref= 5 | 1 5 5 0 6 -1 | page found at 1 | | | |
| 14 | ref= 1 | 1 5 5 0 6 -1 | page found at 0 | | | |
| 15 | ref= 6 | 1 5 5 0 6 -1 | page found at 3 | | | |
| 16 | ref= 5 | 1 5 5 0 6 -1 | page found at 1 | | | |
| 17 | ref= 0 | 1 5 5 0 6 -1 | page found at 2 | | | |
| 18 | ref= 9 | 1 5 5 0 6 -1 | page fault, replace used page 3 | | | |
| 19 | ref= 0 | 1 5 5 0 9 -1 | page found at 9 | | | |
| 20 | ref= 5 | 1 5 5 0 9 -1 | page found at 1 | | | |
| 21 | ref= 9 | 1 5 5 0 9 -1 | page found at 3 | | | |
| 22 | ref= 0 | 1 5 5 0 9 -1 | page found at 2 | | | |
| 23 | ref= 9 | 1 5 5 0 9 -1 | page found at 3 | | | |
| 24 | ref= 3 | 1 5 5 0 9 -1 | page fault, replace used page 1 | | | |
| 25 | ref= 1 | 1 3 5 0 9 -1 | page found at 0 | | | |
| 26 | ref= 9 | 1 3 5 0 9 -1 | page found at 3 | | | |
| 27 | ref= 3 | 1 3 5 0 9 -1 | page found at 1 | | | |
| 28 | ref= 9 | 1 3 5 0 9 -1 | page found at 3 | | | |
| 29 | ref= 6 | 1 3 5 0 9 -1 | page fault, replace used page 1 | | | |
| 30 | ref= 2 | 1 6 5 0 9 -1 | page fault, replace used page 1 | | | |
| 31 | ref= 9 | 1 6 5 0 9 -1 | page found at 3 | | | |
| 32 | ref= 0 | 1 6 5 0 9 -1 | page found at 2 | | | |
| 33 | ref= 9 | 1 6 5 0 9 -1 | page found at 3 | | | |
| 34 | ref= 0 | 1 6 5 0 9 -1 | page found at 2 | | | |
| 35 | ref= 5 | 1 6 5 0 9 -1 | page found at 1 | | | |
| 36 | ref= 0 | 1 2 0 -1 | page fault, replace empty page 0 | | | |
| 37 | ref= 5 | 1 2 0 -1 | page found at 1 | | | |
| 38 | ref= 0 | 1 2 0 -1 | page found at 2 | | | |
| 39 | ref= 2 | 1 2 0 -1 | page found at 1 | | | |
| 40 | ref= 8 | 1 2 0 -1 | page fault, replace empty page 8 | | | |
| 41 | ref= 0 | 1 2 0 -1 | page found at 2 | | | |
| 42 | ref= 6 | 1 2 0 -1 | page fault, replace used page 3 | | | |
| 43 | ref= 0 | 1 2 0 -1 | page found at 2 | | | |
| 44 | ref= 2 | 1 2 0 -1 | page found at 1 | | | |
| 45 | ref= 7 | 1 2 0 -1 | page fault, replace used page 1 | | | |
| 46 | ref= 0 | 1 7 0 -1 | page found at 6 | | | |
| 47 | ref= 7 | 1 7 0 -1 | page found at 1 | | | |
| 48 | ref= 1 | 1 7 0 -1 | page found at 0 | | | |
| 49 | ref= 7 | 1 7 0 -1 | page found at 1 | | | |
| Number of hits 38 | | | | | | |
| Number of page faults 12 | | | | | | |

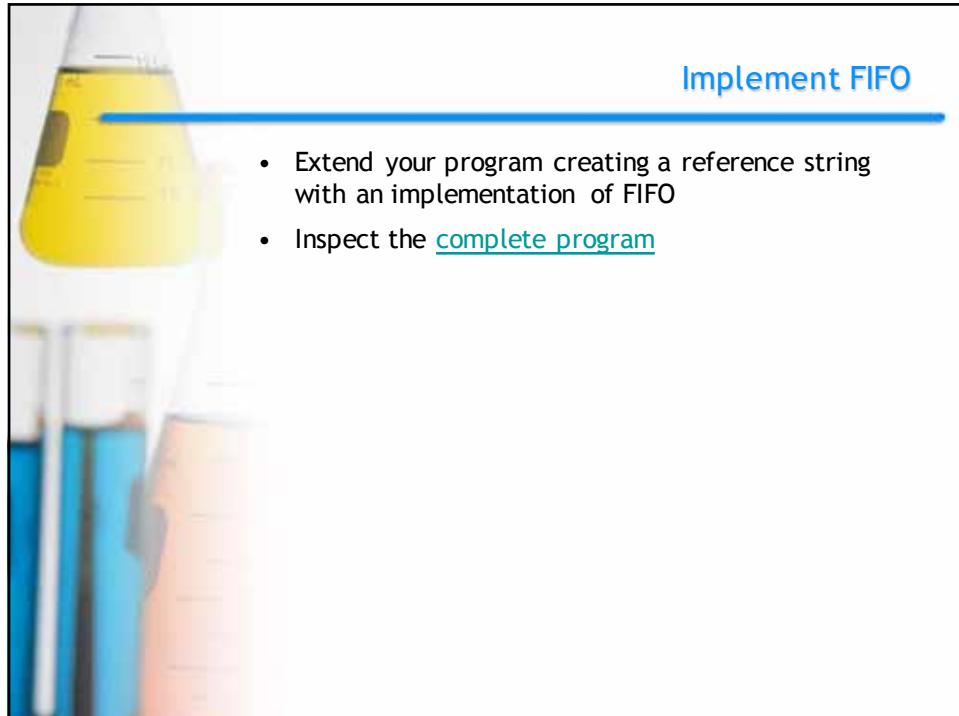
This and all subsequent examples: a set of 50 references are executed on main memory consisting of 4 pages in total. The virtual address space has a size of 10 pages (0..9)

FIFO

- Replace the page which is in main memory for the longest time
- Does not take reference locality into account
 - Frequently used pages (hot spots) are paged out
 - Rarely referenced pages stay in main memory
 - Start up code of some application and others
- Easy to implement via list



Operating Systems

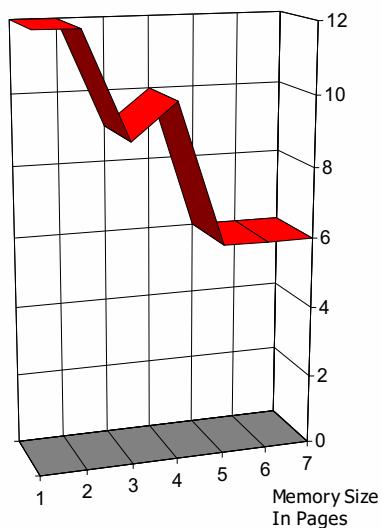


- Extend your program creating a reference string with an implementation of FIFO
- Inspect the [complete program](#)

| FIFO Example | | | | | | |
|--------------------------|---|----|----|----|----------------------------------|---------------------------------|
| Algorithm FIFO | | | | | | |
| 0 ref= | 1 | -1 | -1 | -1 | page fault, replace empty page 0 | |
| 1 ref= | 5 | 1 | -1 | -1 | page fault, replace empty page 1 | |
| 2 ref= | 1 | 1 | 5 | -1 | page found at 0 | |
| 3 ref= | 5 | 1 | 5 | -1 | page found at 1 | |
| 4 ref= | 2 | 1 | 5 | -1 | page found at 2 | |
| 5 ref= | 5 | 1 | 5 | -1 | page found at 1 | |
| 6 ref= | 0 | 1 | 5 | -1 | page found at 0 | |
| 7 ref= | 1 | 1 | 5 | 0 | page fault, replace empty page 2 | |
| 8 ref= | 3 | 1 | 5 | 0 | page found at 0 | |
| 9 ref= | 5 | 1 | 5 | 0 | page found at 1 | |
| 10 ref= | 1 | 1 | 5 | 0 | page found at 0 | |
| 11 ref= | 5 | 1 | 5 | 0 | page found at 1 | |
| 12 ref= | 6 | 1 | 5 | 0 | page fault, replace used page 3 | |
| 13 ref= | 5 | 5 | 0 | 3 | page found at 0 | |
| 14 ref= | 1 | 5 | 0 | 3 | page fault, replace used page 3 | |
| 15 ref= | 6 | 0 | 3 | 6 | page found at 2 | |
| 16 ref= | 5 | 0 | 3 | 6 | page fault, replace used page 3 | |
| 17 ref= | 0 | 3 | 6 | 1 | page fault, replace used page 3 | |
| 18 ref= | 9 | 6 | 1 | 5 | page fault, replace used page 3 | |
| 19 ref= | 0 | 1 | 5 | 0 | page found at 0 | |
| 20 ref= | 5 | 1 | 5 | 0 | page found at 1 | |
| 21 ref= | 9 | 1 | 5 | 0 | page found at 2 | |
| 22 ref= | 0 | 1 | 5 | 0 | page found at 0 | |
| 23 ref= | 9 | 1 | 5 | 0 | page found at 1 | |
| 24 ref= | 3 | 1 | 5 | 0 | page fault, replace used page 3 | |
| 25 ref= | 1 | 5 | 0 | 9 | page fault, replace used page 3 | |
| 26 ref= | 9 | 0 | 9 | 3 | page found at 1 | |
| 27 ref= | 3 | 0 | 9 | 1 | page found at 2 | |
| 28 ref= | 9 | 0 | 9 | 1 | page found at 1 | |
| 29 ref= | 6 | 0 | 9 | 3 | page fault, replace used page 3 | |
| 30 ref= | 2 | 9 | 3 | 1 | page fault, replace used page 3 | |
| 31 ref= | 9 | 2 | 9 | 1 | page fault, replace used page 3 | |
| 32 ref= | 0 | 3 | 1 | 6 | page fault, replace used page 3 | |
| 33 ref= | 9 | 0 | 6 | 2 | page found at 0 | |
| 34 ref= | 0 | 6 | 2 | 9 | page found at 3 | |
| 35 ref= | 1 | 6 | 2 | 9 | 0 | page fault, replace used page 3 |
| 36 ref= | 0 | 2 | 9 | 0 | 1 | page found at 2 |
| 37 ref= | 9 | 2 | 9 | 0 | 1 | page found at 1 |
| 38 ref= | 0 | 2 | 9 | 0 | 1 | page found at 2 |
| 39 ref= | 2 | 2 | 9 | 0 | 1 | page found at 0 |
| 40 ref= | 8 | 2 | 9 | 0 | 1 | page fault, replace used page 3 |
| 41 ref= | 0 | 9 | 0 | 1 | 8 | page found at 1 |
| 42 ref= | 6 | 9 | 0 | 1 | 8 | page fault, replace used page 3 |
| 43 ref= | 0 | 0 | 1 | 8 | 6 | page found at 0 |
| 44 ref= | 2 | 0 | 1 | 8 | 6 | page fault, replace used page 3 |
| 45 ref= | 7 | 1 | 8 | 6 | 2 | page fault, replace used page 3 |
| 46 ref= | 0 | 8 | 6 | 2 | 7 | page fault, replace used page 3 |
| 47 ref= | 7 | 6 | 2 | 7 | 0 | page found at 3 |
| 48 ref= | 1 | 6 | 2 | 7 | 0 | page fault, replace used page 3 |
| 49 ref= | 7 | 2 | 7 | 0 | 1 | page found at 1 |
| Number of hits 28 | | | | | | |
| Number of page faults 22 | | | | | | |

Belady's Anomaly with FIFO

- Reference string
 - 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5
- Higher number of page faults with more main memory!

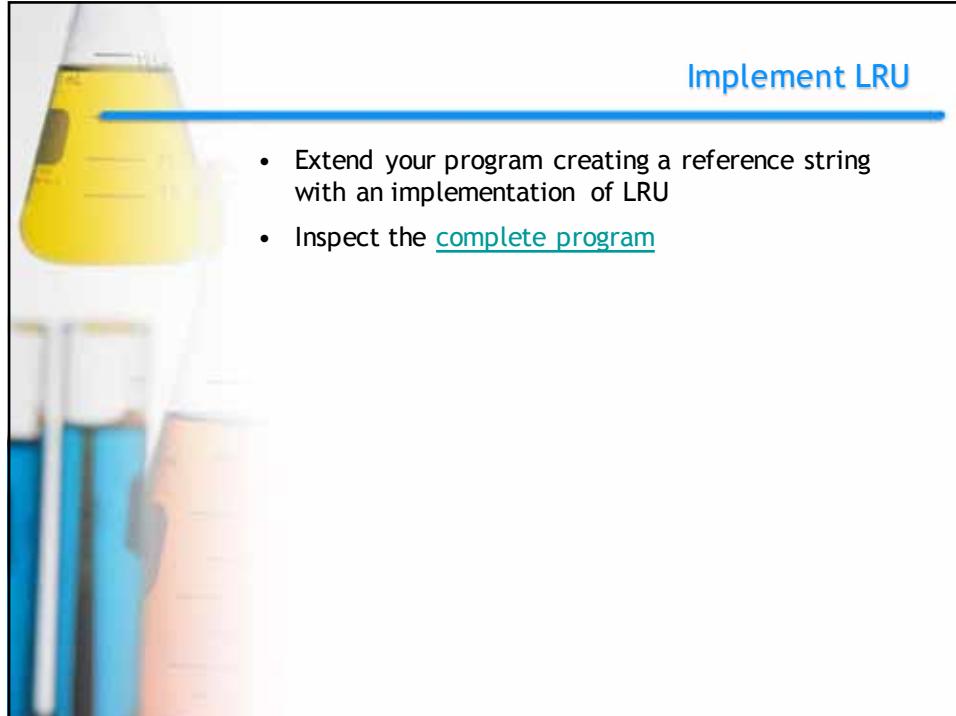


Operating Systems

Least Recently Used (LRU)

- Replace the page which has not been referenced for the longest time in the past
 - Again reference locality
 - What you don't reference in the past for a long time you won't reference in the near future
 - Same with your CDs, cloths, etc.
 - Project the past onto the future
- LRU is yet the best approximation to Belady
- Hard to implement

Operating Systems



- Extend your program creating a reference string with an implementation of LRU
- Inspect the [complete program](#)

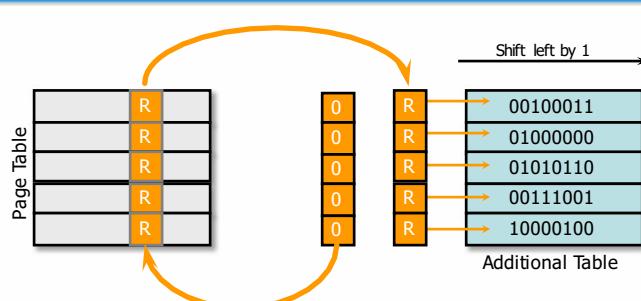
| LRU Example | | | | | | |
|--------------------------|---|----|----|----|----|----------------------------------|
| <pre>Algorithm LRU</pre> | | | | | | |
| 0 ref= | 1 | -1 | -1 | -1 | | page fault, replace empty page 0 |
| 1 ref= | 5 | 1 | -1 | -1 | | page fault, replace empty page 1 |
| 2 ref= | 1 | 1 | 5 | -1 | | page found at 0 |
| 3 ref= | 5 | 1 | 5 | -1 | | page found at 1 |
| 4 ref= | 1 | 1 | 5 | -1 | | page found at 0 |
| 5 ref= | 5 | 1 | 5 | -1 | | page found at 1 |
| 6 ref= | 0 | 1 | 5 | -1 | | page fault, replace empty page 2 |
| 7 ref= | 1 | 1 | 5 | 0 | -1 | page found at 0 |
| 8 ref= | 3 | 1 | 5 | 0 | -1 | page fault, replace empty page 3 |
| 9 ref= | 5 | 1 | 5 | 0 | -1 | page found at 0 |
| 10 ref= | 1 | 1 | 5 | 0 | 3 | page found at 0 |
| 11 ref= | 5 | 1 | 5 | 0 | 3 | page found at 1 |
| 12 ref= | 6 | 1 | 5 | 0 | 3 | page fault, replace used page 2 |
| 13 ref= | 5 | 1 | 5 | 6 | 3 | page found at 1 |
| 14 ref= | 1 | 1 | 5 | 6 | 3 | page found at 0 |
| 15 ref= | 6 | 1 | 5 | 6 | 3 | page found at 2 |
| 16 ref= | 5 | 1 | 5 | 6 | 3 | page found at 1 |
| 17 ref= | 0 | 1 | 5 | 6 | 3 | page fault, replace used page 3 |
| 18 ref= | 9 | 1 | 5 | 6 | 0 | page fault, replace used page 0 |
| 19 ref= | 0 | 9 | 5 | 6 | 0 | page found at 3 |
| 20 ref= | 5 | 9 | 5 | 6 | 0 | page found at 1 |
| 21 ref= | 9 | 9 | 5 | 6 | 0 | page found at 0 |
| 22 ref= | 0 | 9 | 5 | 6 | 0 | page found at 0 |
| 23 ref= | 9 | 9 | 5 | 6 | 0 | page found at 0 |
| 24 ref= | 3 | 9 | 5 | 6 | 0 | page fault, replace used page 2 |
| 25 ref= | 1 | 9 | 5 | 0 | 0 | page fault, replace used page 1 |
| 26 ref= | 9 | 9 | 1 | 0 | 0 | page found at 0 |
| 27 ref= | 3 | 9 | 1 | 0 | 0 | page found at 2 |
| 28 ref= | 9 | 9 | 1 | 0 | 0 | page found at 0 |
| 29 ref= | 6 | 9 | 1 | 0 | 0 | page fault, replace used page 3 |
| 30 ref= | 2 | 9 | 1 | 6 | 0 | page fault, replace used page 1 |
| 31 ref= | 9 | 9 | 2 | 6 | 0 | page found at 0 |
| 32 ref= | 0 | 9 | 2 | 6 | 0 | page fault, replace used page 2 |
| 33 ref= | 9 | 9 | 2 | 0 | 6 | page found at 0 |
| 34 ref= | 0 | 9 | 2 | 0 | 6 | page found at 2 |
| 35 ref= | 1 | 9 | 2 | 0 | 6 | page fault, replace used page 3 |
| 36 ref= | 0 | 9 | 2 | 0 | 1 | page found at 0 |
| 37 ref= | 9 | 9 | 2 | 0 | 1 | page found at 0 |
| 38 ref= | 0 | 9 | 2 | 0 | 1 | page found at 2 |
| 39 ref= | 2 | 9 | 2 | 0 | 1 | page found at 1 |
| 40 ref= | 8 | 9 | 2 | 0 | 1 | page fault, replace used page 3 |
| 41 ref= | 0 | 9 | 2 | 0 | 8 | page found at 2 |
| 42 ref= | 6 | 9 | 2 | 0 | 8 | page fault, replace used page 0 |
| 43 ref= | 0 | 6 | 2 | 0 | 8 | page found at 2 |
| 44 ref= | 2 | 6 | 2 | 0 | 8 | page found at 1 |
| 45 ref= | 7 | 6 | 2 | 0 | 8 | page fault, replace used page 3 |
| 46 ref= | 0 | 6 | 2 | 0 | 7 | page found at 2 |
| 47 ref= | 7 | 6 | 2 | 0 | 7 | page found at 3 |
| 48 ref= | 3 | 6 | 2 | 0 | 7 | page fault, replace used page 0 |
| 49 ref= | 7 | 1 | 2 | 0 | 7 | page found at 3 |
| Number of hits 33 | | | | | | |
| Number of page faults 17 | | | | | | |

Implementing LRU

- LRU is hard to implement
- Problem
 - Sort all pages by last access time
 - Every reference (write and read) has to be considered
- Counter-based implementation
 - Add counter field to page descriptor
 - Add logical clock to CPU (Increment with every reference)
 - MMU updates descriptor with CPU clock for every reference
 - Page with the smallest clock value is the victim
- Stack-based implementation
 - Stack of referenced pages
 - Referenced pages are pushed on top of stack
 - Page at the stacks' tail is victim

Operating Systems

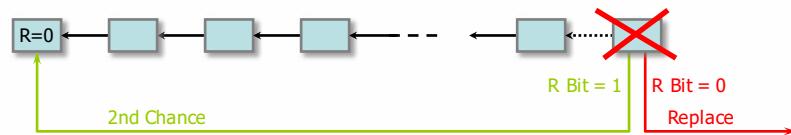
Approaching LRU (1)



- Additional table
 - Utilizing the referenced bit in each descriptor
 - Periodically (e.g. 100 msec), the referenced flag updates the MSBs in the additional table
 - Table entry reflects the last n references to each page
- Page with the smallest numerical value must leave
- Solution by A. Tanenbaum

Operating Systems

Approaching LRU (2): 2nd Chance



- Based on FIFO
- Page selected by FIFO
 - Referenced Bit = 0: Page must leave
 - Referenced Bit = 1: Second chance for page
 - Page is moved to the other end of FIFO again with
 - Referenced bit reset
- If all referenced bits set, 2nd chance degenerates into FIFO

Operating Systems

Implement 2nd Chance

- Extend your program creating a reference string with an implementation of 2nd Chance
- Inspect the [complete program](#)



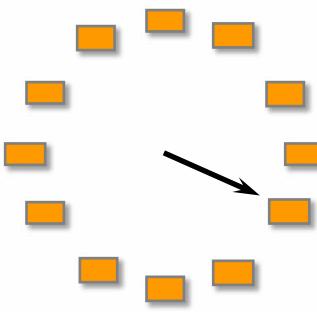
```

Algorithm Secondchance
0 ref= 1 -1- -1- -1- -1- | page fault, replace empty page 0
1 ref= 5 1R -1- -1- -1- -1- | page fault, replace empty page 1
2 ref= 1 1R SR -1- -1- -1- | page found at 1
3 ref= 5 1R SR -1- -1- -1- | page found at 1
4 ref= 1 1R SR -1- -1- -1- | page found at 0
5 ref= 5 1R SR -1- -1- -1- | page found at 1
6 ref= 0 1R SR -1- -1- -1- | page fault, replace empty page 2
7 ref= 1 1R SR OR -1- -1- | page found at 0
8 ref= 3 1R SR OR -1- -1- | page fault, replace empty page 3
9 ref= 5 1R SR OR 3R -1- | page found at 1
10 ref= 1 1R SR OR 3R -1- | page found at 0
11 ref= 5 1R SR OR 3R -1- | page found at 1
12 ref= 6 1R SR OR 3R -1- | page fault, replace used page 0
13 ref= 5 6R 5- 0- 3- -1- | page found at 1
14 ref= 1 6R 5- 0- 3- -1- | page fault, replace used page 0
15 ref= 6 1R 3- 6- -2- -1- | page found at 3
16 ref= 5 1R 3- 6R 5- -1- | page found at 3
17 ref= 0 1R 3- 6R SR -1- | page fault, replace used page 0
18 ref= 9 0R 6R 1- -1- -1- | page fault, replace used page 0
19 ref= 0 9R 0- 6- -5- -1- | page found at 1
20 ref= 5 9R OR 6- -5- -1- | page found at 3
21 ref= 9 9R OR 6- SR -1- | page found at 0
22 ref= 0 9R OR 6- SR -1- | page found at 1
23 ref= 9 9R OR 6- SR -1- | page found at 0
24 ref= 3 9R OR 6- SR -1- | page fault, replace used page 0
25 ref= 1 3R SR 9- 0- -1- | page fault, replace used page 0
26 ref= 9 1R 0- 3- 5- -1- | page fault, replace used page 0
27 ref= 3 9R 3- 5- 1- -1- | page found at 1
28 ref= 9 9R 3R 5- 1- -1- | page fault, replace used page 0
29 ref= 6 9R 3R 5- 1- -1- | page fault, replace used page 0
30 ref= 2 6R 1- 9- 3- -1- | page found at 1
31 ref= 9 2R 9- 3- 6- -1- | page found at 1
32 ref= 0 2R 9R 3- 6- -1- | page fault, replace used page 0
33 ref= 9 0R 6- 2- 9- -1- | page found at 3
34 ref= 0 0R 6- 2- 9R -1- | page found at 0
35 ref= 1 0R 6- 2- 9R -1- | page fault, replace used page 0
36 ref= 0 1R 2- 9R 0- -1- | page found at 3
37 ref= 9 1R 2- 9R OR -1- | page found at 2
38 ref= 0 1R 2- 9R OR -1- | page found at 3
39 ref= 2 1R 2- 9R OR -1- | page found at 1
40 ref= 8 1R 2R 9R OR -1- | page fault, replace used page 0
41 ref= 0 8R 2- 9R 0- -1- | page found at 3
42 ref= 6 8R 2- 9R OR -1- | page fault, replace used page 0
43 ref= 0 0R 9- 2- 0R -1- | page found at 3
44 ref= 2 6R 9- 0R 8- -1- | page fault, replace used page 0
45 ref= 7 2R OR 8- 6- -1- | page fault, replace used page 0
46 ref= 0 7R 6- 2- 0R -1- | page found at 3
47 ref= 7 7R 6- 2- 0R -1- | page found at 0
48 ref= 1 7R 6- 2- 0R -1- | page fault, replace used page 0
49 ref= 7 1R 2- OR 7- -1- | page found at 3
Number of hits 30
Number of page faults 20

```

2nd Chance Example

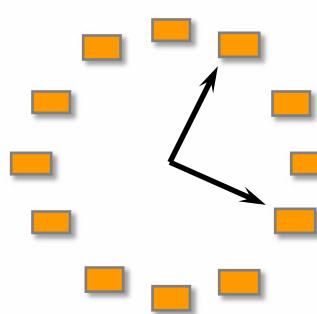
- Minor variation of 2nd chance
 - Updating the FIFO queue is too expensive
 - Using a circular list instead
- Hand of clock points to „oldest“ page
 - Referenced bit = 0
 - Replace page
 - Referenced bit = 1
 - Reset flag
 - Move hand



Operating Systems

2 Hand Version ☺

- Single hand version to clumsy for todays memory sizes
- Two Hands
 - Hand 1 resets flag
 - Hand 2 checks flag
 - Both hands „tick“ synchronously
- Distance between hands determines time period for second chance
 - Small: Only frequently used pages get a second chance
 - Maximum: Identical to single hand version

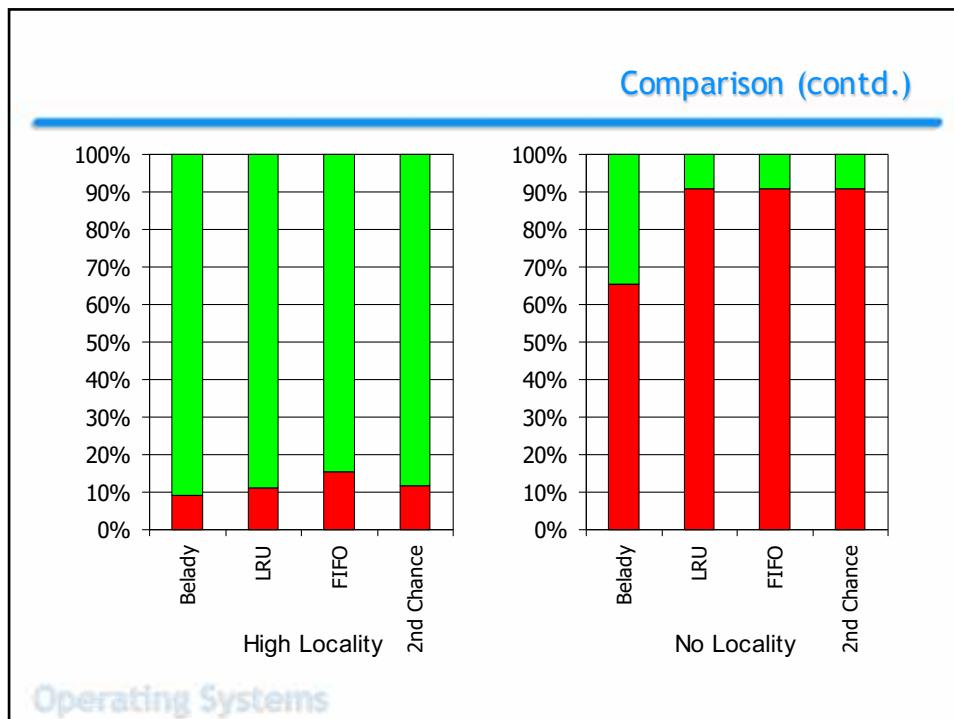


Operating Systems

Comparison Belady, LRU, FIFO, 2nd Chance

- High reference locality
 - 0.9 probability to reference one of the past 20 pages
- No reference locality
 - Next page is chosen randomly
- Configuration
 - 1000000 references
 - References within pages 0 and 100
 - Size of main memory 10 pages

Operating Systems



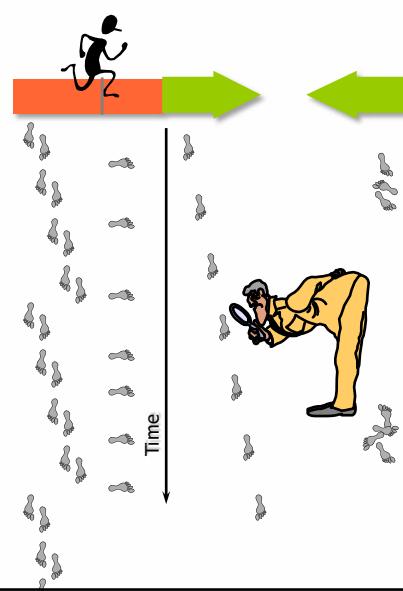
- ### Other Replacement Algorithms
- Random
 - A randomly selected page is the victim
 - Least Frequently Used (LFU)
 - Increment counter with each page reference
 - Page with the smallest value will be discarded
 - Problem with out-of-date hot spots
 - Shift counter periodically to the left
 - Exponential attenuation
 - ...
- Operating Systems*

Isn't main memory nothing but a huge cache?

Operating Systems

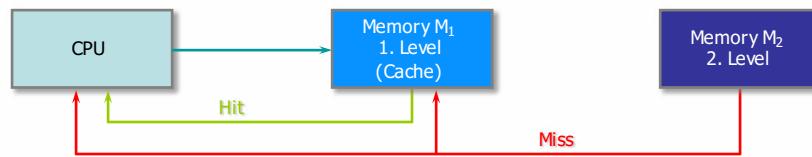
Reference Locality

- Instructions are repeated
- Variables are referenced repeatedly
- If one address is referenced, nearby addresses will be referenced with high probability
- Reasons:
 - Sequential program execution
 - Procedure nesting is modest
 - Short loops
 - Frequent processing of arrays, list, and records



Operating Systems

Caching: Utilizing Reference Locality



- 2-Level-Storage
- Characteristic parameters of M_k:
 - T_k = Average access time
 - P_k = Price per byte
 - S_k = Size of memory in bytes
- Hit rate H
 - Probability, that a referenced address is in cache

Operating Systems

Average Access Time and Costs

- Average access time:

$$T = HT_1 + (1-H)(T_1 + T_2) = T_1 + (1-H)T_2$$

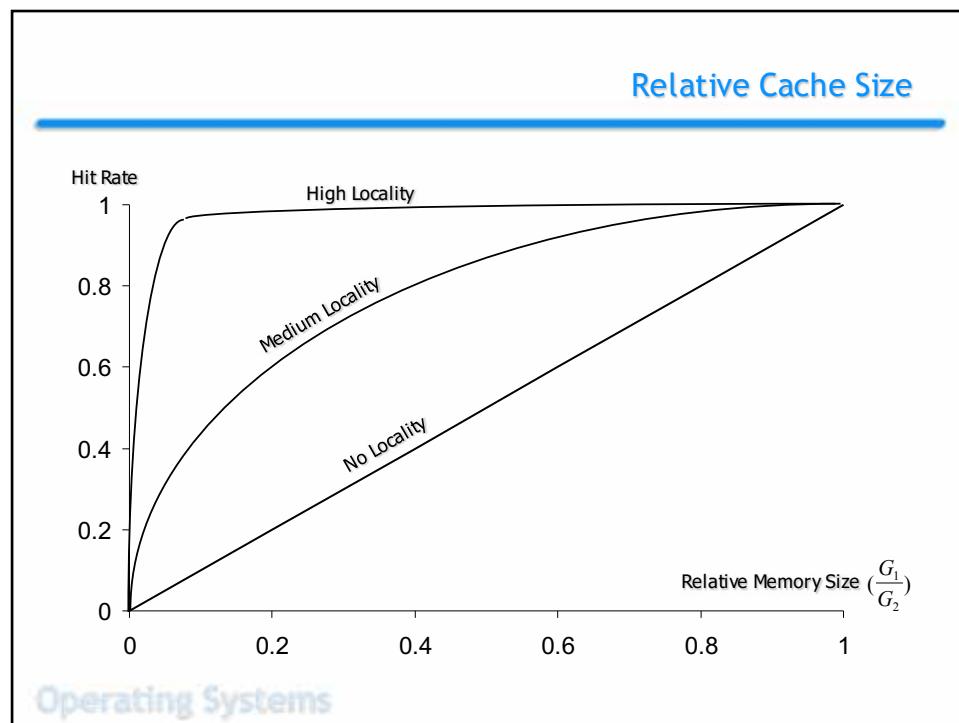
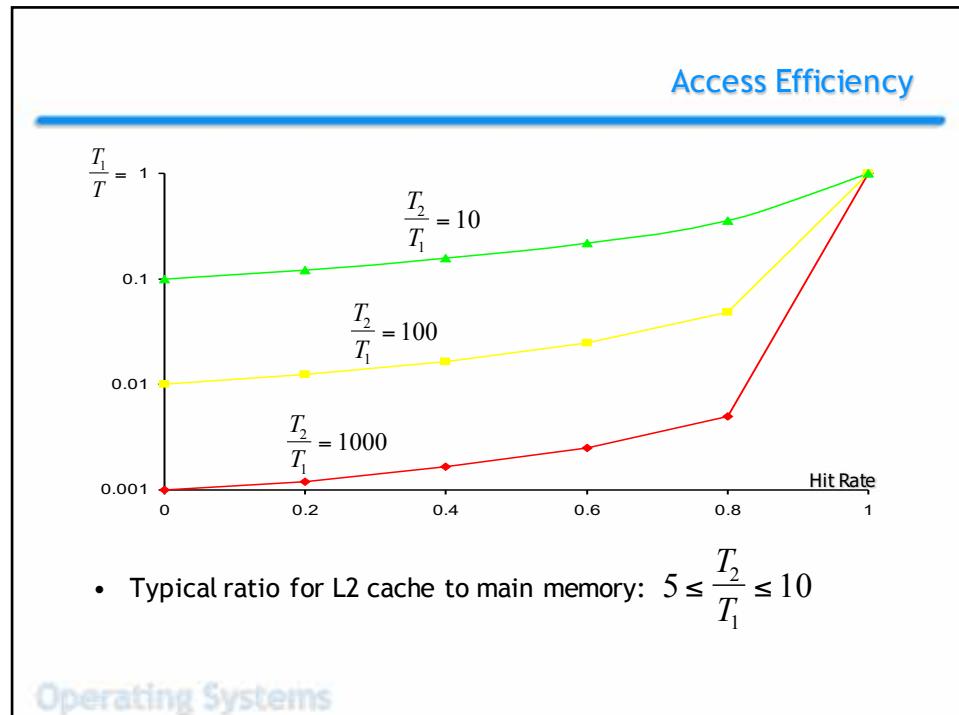
- Costs:

$$P = \frac{P_1 S_1 + P_2 S_2}{S_1 + S_2}$$

- Access efficiency:

$$\frac{T_1}{T} = \frac{1}{H + (1-H)\frac{T_2}{T_1}}$$

Operating Systems



Optimization Goal

- Strive for achieve

$$T_1 \ll T_2 : T \approx T_1$$

-

$$S_1 \ll S_2 : S \approx S_2$$

- Required hit rate H between 0.8 and 0.9
- Very small caches ($S_1 \ll S_2$) yield high hit rates already

Operating Systems

Multi-Level Caching



- Take advantage of memory hierarchies
 - L1 cache and MMU TLB to L2 cache
 - L2 cache to main memory (sometimes also L3 cache)
 - Main memory to external storage (disk)
- Caching between main memory and external storage promising:
 - Access time ratio: 1: 10^5
 - Size ratio: 1: 10^2 - 10^3
 - Cost ratio: 1:200
 - June 2005: 1 GB RAM 120 Dollar, 200 GB Disk 120 Dollar
 - November 2008: 1 GB RAM 20 Dollar, 1000 GB Disk 120 Dollar

Operating Systems

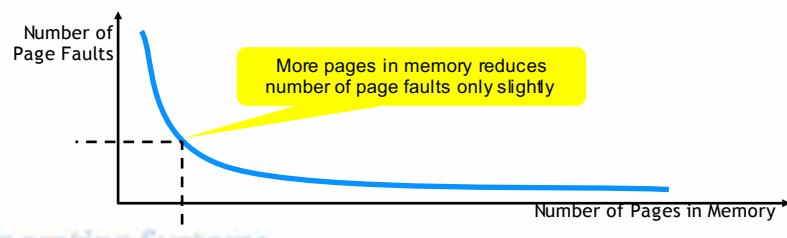
Operating Systems

Paging
Working Set Theory

System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

Motivation

- Most programs don't exploit the complete virtual address space
 - Memory footprint in many cases a few mbyte only
 - Of course, huge server applications such as databases and heavy web servers are different
- Due to reference locality, at a given time interval Δt they require even less pages
- Can we limit the number of pages in main memory for every given application?

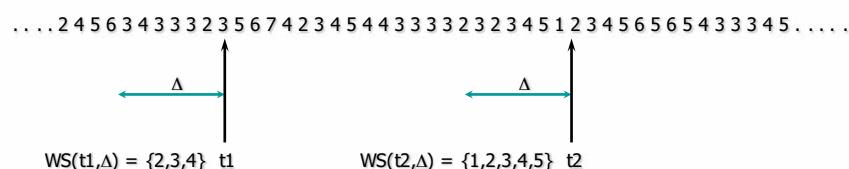


Locality Set and Working Set

- Operating point k
 - Less than k pages in memory \Rightarrow Many page faults
 - More than k pages in memory \Rightarrow No substantial change in number of faults
 - k depends on application
 - Optimal number of k pages = Locality Set
 - Working set (WS) of a thread
- $WS(t, \Delta) = \text{last } \Delta \text{ referenced pages before } t$
- Selection of Δ critical:
 - Too small: Locality set not part of working set
 - Too large: Comprises too many locality sets
 - waste of resources
 - Sensible values were around 10^4 (in 1968)

Operating Systems

Example



- Size of WS relevant for page assignment:

$$D = \sum_p |WS_p(t, \Delta)|$$

- $D >$ Available number of pages: severe memory congestion
 - Swap whole applications (address space)
 - Stop executing some applications (reduced number of faults)
- $D <$ Available number of pages: Resume/Start other threads

Operating Systems

Computing the Working Set

The diagram shows a horizontal row of nine empty square boxes representing memory pages. A double-headed arrow below the row is labeled with the Greek letter Δ , indicating a window or range of pages. A small black stick figure with a star on its head is on the far left, pointing towards the first few boxes. Another small black stick figure is on the far right, pointing towards the last few boxes.

- Must be calculated for every thread individually
- Costly implementation
 - Writes and reads must be considered
 - Referenced pages are put in front of queue
 - Last page may be dropped
 - Every page in queue is part of working set
- Hard to recognize size changes in WS
- Approximations possible (just like LRU)

Operating Systems

Number of Assigned Pages

- As little as possible
 - Supports more address spaces in parallel
 - Free pages are available just in case
 - Disadvantage: Page faults, but free main memory
- As much as possible
 - Very low page fault rate
 - Applications started first are preferred
 - Disadvantage: Maybe waste of resources
- Same number of pages for every application?

Operating Systems

Thrashing

- Threshold k
- More than k pages allocated
 - Acceptable number of page faults
 - Some pages referenced rarely
 - Surplus on pages better used by other applications
- Exactly k pages allocated
 - Optimum
- Less than k pages allocated
 - Thread exhibits high page fault rate
 - OS loads pages continuously
 - Neglectable progress in thread
- How do we determine k ?
 - Application dependent

Relative CPU Utilization

K Assigned Number of Pages

Operating Systems

Thrashing is a Global Effect

Main Memory

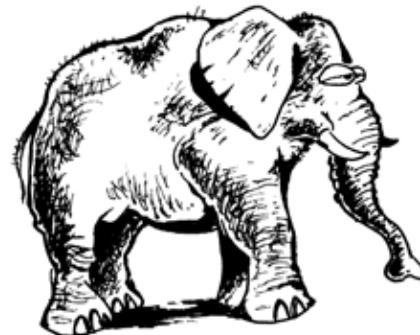
Pagefile on Disk

- Thrashing starts with one or more address spaces
 - Not enough memory assigned to application
- Increasing number of page ins and page outs on disk
- Average access time to disk increases
- Without precautions, the complete system stalls

Operating Systems

Lets Thrash 😊

- Write a small program that occupies a given number of memory pages very strongly?
- How does the operating system behave with an increasing number of program instances?



Operating Systems

Avoid Thrashing



- Tracking page fault rate
- Exceeding a threshold frequency: Assume thrashing
 - Assign additional pages (if possible)
 - Stop applications (write their pages to disk)
- Falling below a minimum threshold frequency
 - Take away pages
 - Continue stopped applications

Operating Systems

Global and Local Replacement

- Are other address spaces involved in page fault?
- Local
 - No, another page of same address space must leave
 - Number of assigned pages remains constant
 - Observing working set useful
- Global
 - Yes, any other page might be replaced
 - Number of assigned pages varies
 - Operating system adapts to specific needs (e.g. priorities)
- Again, hybrid solutions are possible

Operating Systems

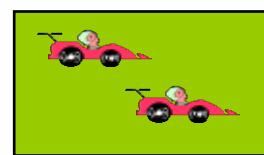
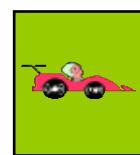
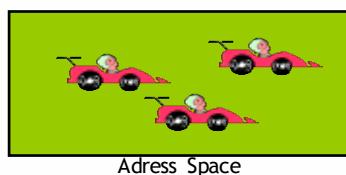
Operating Systems

Threads
Motivation

System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

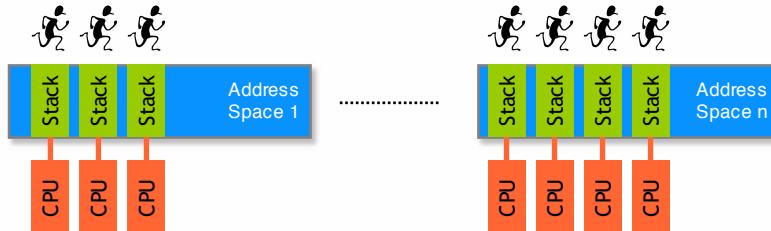
Vision

- Assume that there is an unlimited number of processors available
- Within an application
 - The required number of processors primarily depends on the inherent parallelism
- Between application
 - Every application may have its own processors



Operating Systems

Why Its Useful to Have “Unlimited” CPUs!

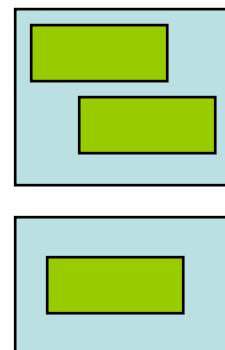


- Abstract from the number of CPUs on a given computer
 - Exploit concurrency inherent to an application
 - Utilize possible parallelism of multiprocessor systems
- Other reasons for multiple threads
 - Improve response time of interactive applications
 - Utilize parallelism available in most modern computer hardware

Operating Systems

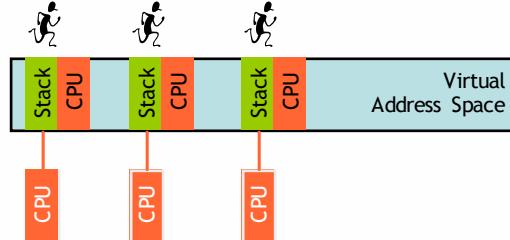
... But 1 Physical CPU Available?

- Easy to map many virtual CPUs to a physical CPU
- A "Thread of Control" consists of
 - Address space
 - State of registers
 - Stack (already part of address space)
- Context Switch
 - Save current thread state
 - Restore another thread state
- When to switch
 - Calls that might block for a longer period of time
 - Implicitly, e.g. in case of page fault or memory congestion



Operating Systems

CPU Multiplexing



- Address space already maintained by operating system
- Each thread will have its runtime stack as part of address space
- Time-Multiplexing the CPU (Context Switch)
 - Storing the register content in memory
 - Restoring the register content of a previously stored thread
 - Switch address space if required (change pointer to root page table)

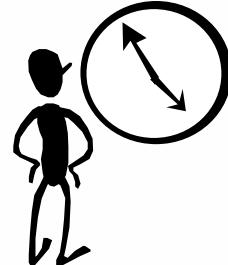
Operating Systems

Context Switches Per Second

- Sleep(0) yields the CPU
- By measuring the time a given number of Sleep(0) requires on an otherwise idle system, the average time for a context switch can be estimated
- Result depends on the CPU speed and other factors
- On a Intel Pentium 4 HT with 3.06 GHz clock speed, we measure approx. 630 ns for a context switch
- The hyper-threading feature leads to 50% system utilization, since only 1 CPU can execute the program

Program Behavior

- Most application programs can't execute continuously for a longer period of time
 - Waiting for user input
 - Waiting for a message
 - Waiting for a free buffer to send a message
 - Waiting for a block from disk
 - waiting for virtual memory to be paged in after page fault
- Some can
 - Number crunching, cryptography, simulations, ...
 - Complex operations on large images (e.g. Photoshop)
 - Realtime games, chess, ...

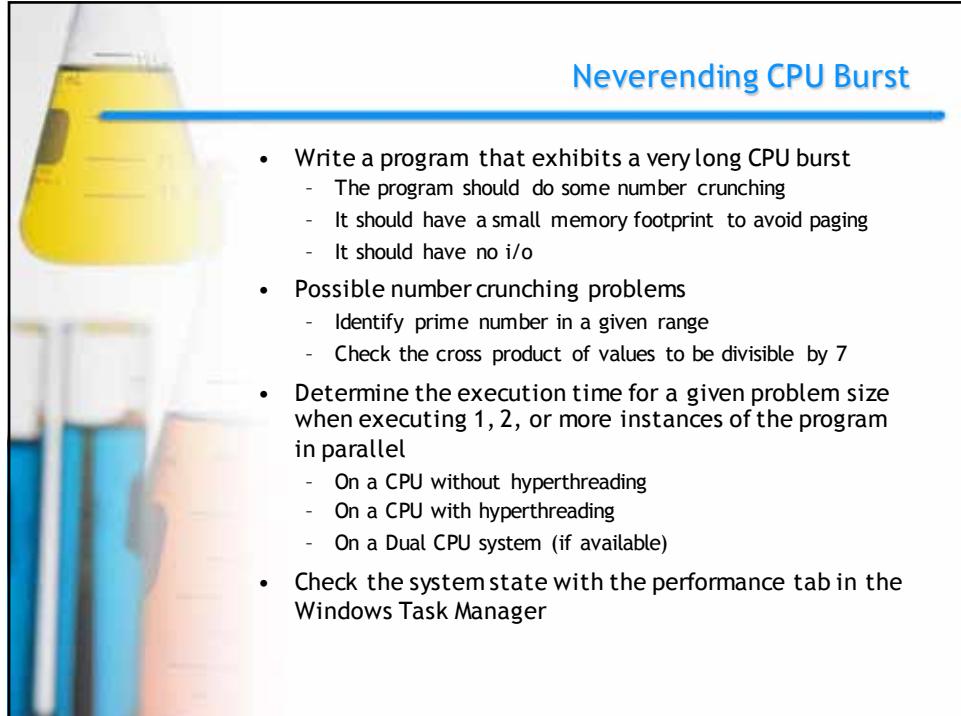


Operating Systems

Bursts

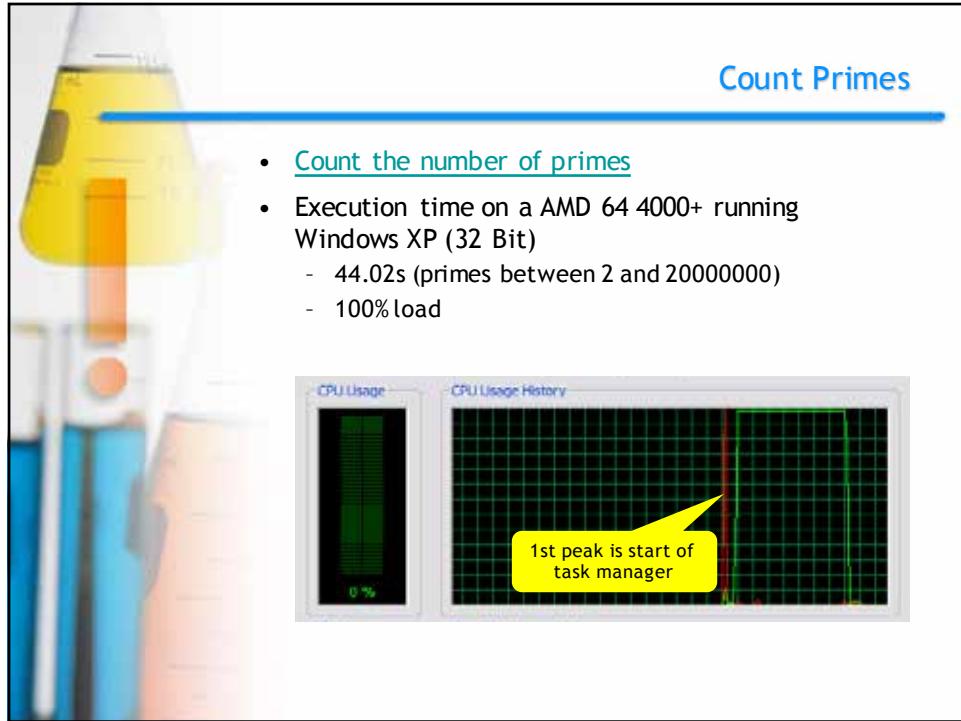
- Different phases of program behavior are called bursts
- CPU burst
 - Thread is running on a CPU
 - Burst length inherent to thread depends on type of application
- IO burst
 - Thread is waiting for some io operation to complete
 - Thread is not running
 - Most io bursts last at least for several milliseconds or longer
- Burst alternation
- Can OS exploit this typical program behavior?

Operating Systems



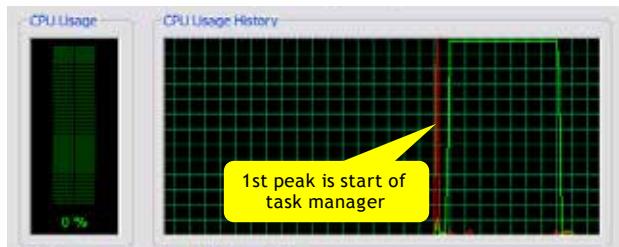
Neverending CPU Burst

- Write a program that exhibits a very long CPU burst
 - The program should do some number crunching
 - It should have a small memory footprint to avoid paging
 - It should have no i/o
- Possible number crunching problems
 - Identify prime number in a given range
 - Check the cross product of values to be divisible by 7
- Determine the execution time for a given problem size when executing 1, 2, or more instances of the program in parallel
 - On a CPU without hyperthreading
 - On a CPU with hyperthreading
 - On a Dual CPU system (if available)
- Check the system state with the performance tab in the Windows Task Manager



Count Primes

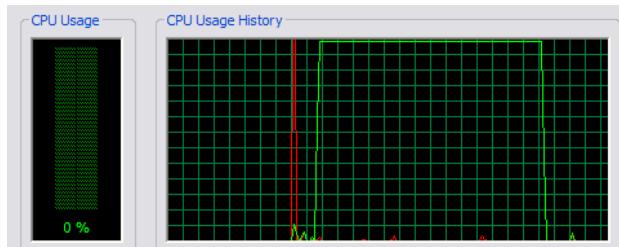
- Count the number of primes
- Execution time on a AMD 64 4000+ running Windows XP (32 Bit)
 - 44.02s (primes between 2 and 20000000)
 - 100% load



1st peak is start of task manager

2 “Count Primes” in Parallel

- Executing 2 of these programs in parallel requires 88 s for both programs to end
- As expected: $2 * 44 \text{ s} = 88 \text{ s}$
- Again 100% load



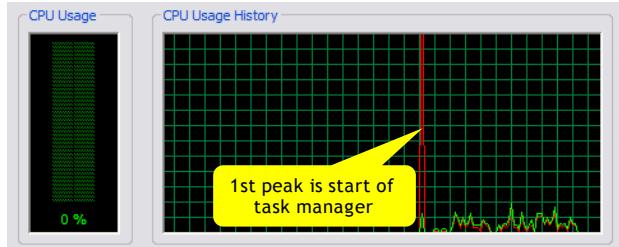
The screenshot shows the Windows Task Manager. On the left, there is a small preview window titled "CPU Usage" showing a grid of green dots. Below it, the text "0 %". On the right, there is a larger window titled "CPU Usage History" showing a graph of CPU usage over time. The graph has a red vertical line at the start and a green vertical line further down, indicating the execution of two parallel tasks.

Mostly IO Burst

- Write a program that exhibits IO bursts only
 - Don't do number crunching at all
 - Write some easy to calculate data to disk
 - Write in a very fine granularity
- Determine the execution time for a given problem size when executing 1, 2, or more instances of the program in parallel
 - On a CPU without hyperthreading
 - On a CPU with hyperthreading
 - On a Dual CPU system (if available)
- Check the system state with the performance tab in the Windows Task Manager

Talk Nonsense

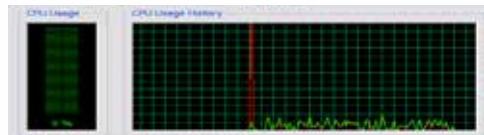
- [Create a file containing nonsense](#)
- Execution time on a AMD 64 4000+ running Windows XP (32 Bit)
 - 48.1s (file size 50000)
 - Between 0% and 4% load



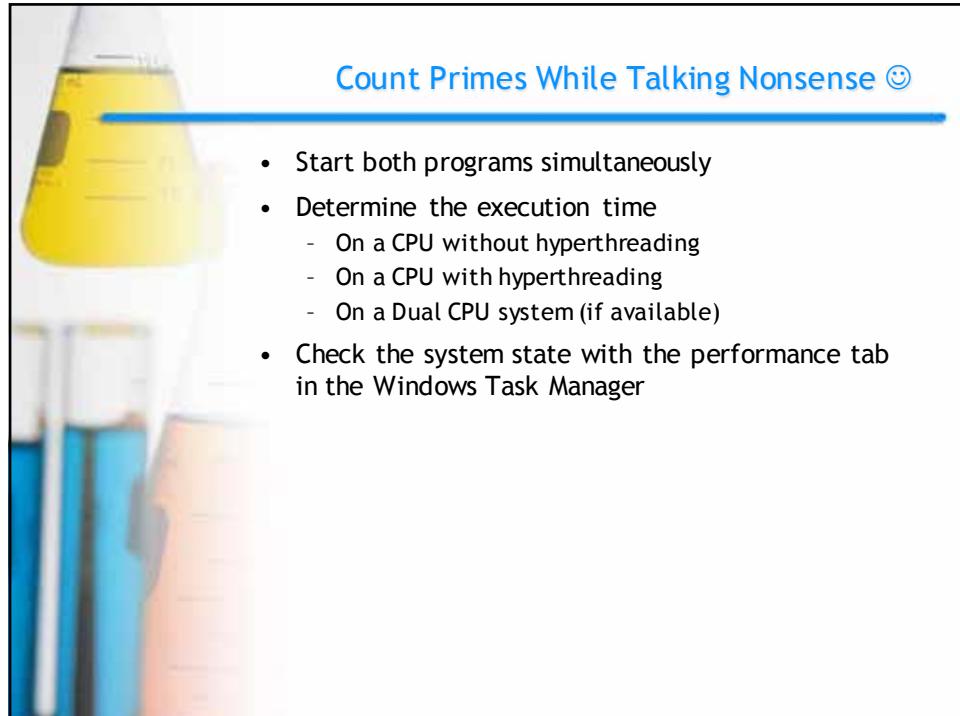
The figure consists of two side-by-side graphs. The left graph is titled 'CPU Usage' and shows a vertical bar chart with '0 %' at the bottom. The right graph is titled 'CPU Usage History' and is a line graph on a grid. A yellow callout box points to a sharp vertical red peak in the history graph, with the text '1st peak is start of task manager'.

2 “Talk Nonsense” in Parallel

- Writing two files on a single disk
 - Both programs finished after 88.7 s
- Writing two files on different disks
 - Both programs finished after 62.2 s

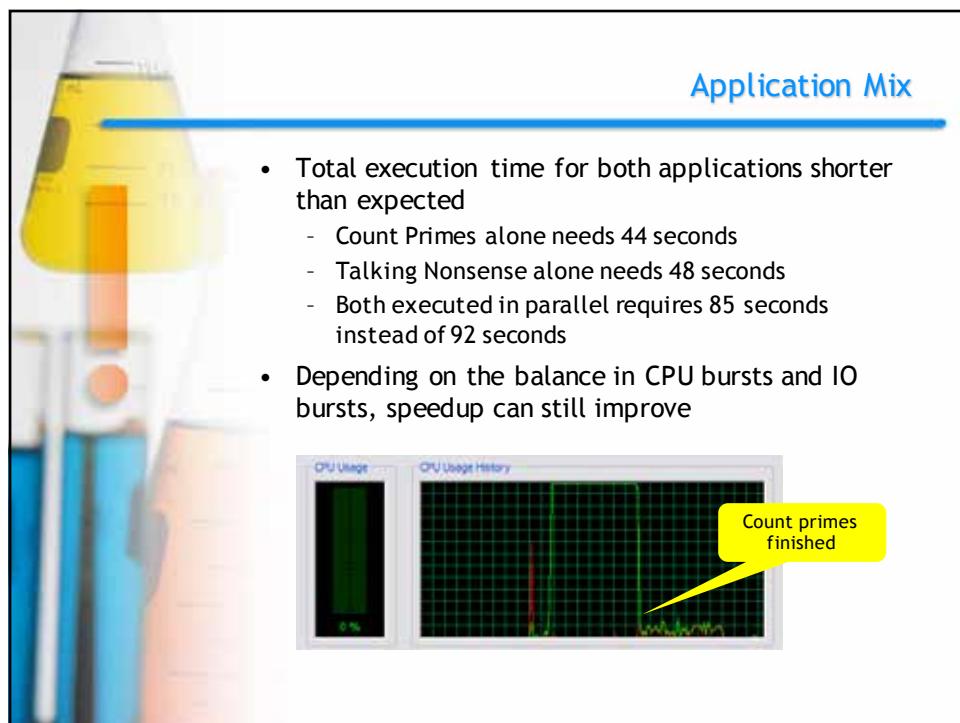


The figure contains two sets of graphs, each pair consisting of a 'CPU Usage' bar chart and a 'CPU Usage History' line graph. The top set of graphs corresponds to writing files on a single disk, showing a single sharp peak in the history graph. The bottom set corresponds to writing files on different disks, showing two distinct peaks in the history graph.



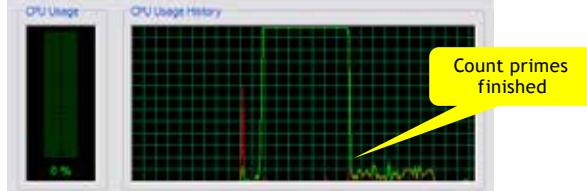
Count Primes While Talking Nonsense ☺

- Start both programs simultaneously
- Determine the execution time
 - On a CPU without hyperthreading
 - On a CPU with hyperthreading
 - On a Dual CPU system (if available)
- Check the system state with the performance tab in the Windows Task Manager



Application Mix

- Total execution time for both applications shorter than expected
 - Count Primes alone needs 44 seconds
 - Talking Nonsense alone needs 48 seconds
 - Both executed in parallel requires 85 seconds instead of 92 seconds
- Depending on the balance in CPU bursts and IO bursts, speedup can still improve



Conclusion

- Parallelism inherent to most modern hardware can be exploited with a mixture of applications
- Speedup possible even on a monoprocessor system
 - Two apps with IO burst, if independent IO hardware is used
 - Applications with CPU bursts and IO bursts may coexist
- In many cases, hyperthreading can be very similar to a dual processor
- The beginning of an IO burst enables an operating system to transfer CPU time to other threads that are willing to execute

Operating Systems

Operating Systems

Threads
Multi-Threading

System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

Single-Threaded Process

- Traditional notion of process
 - An address space
 - Initial thread starting with the execution of `main()`
- Process terminates when initial thread leaves `main()`



Operating Systems

Multiple Threads

- Additional threads must be instantiated explicitly
 - In .NET, additional instances of class Thread must be created
- Reasons to create multiple threads
 - Exploit hardware parallelism
 - Continue computation while another thread blocks
 - React immediately on user input



Operating Systems

Operating Systems

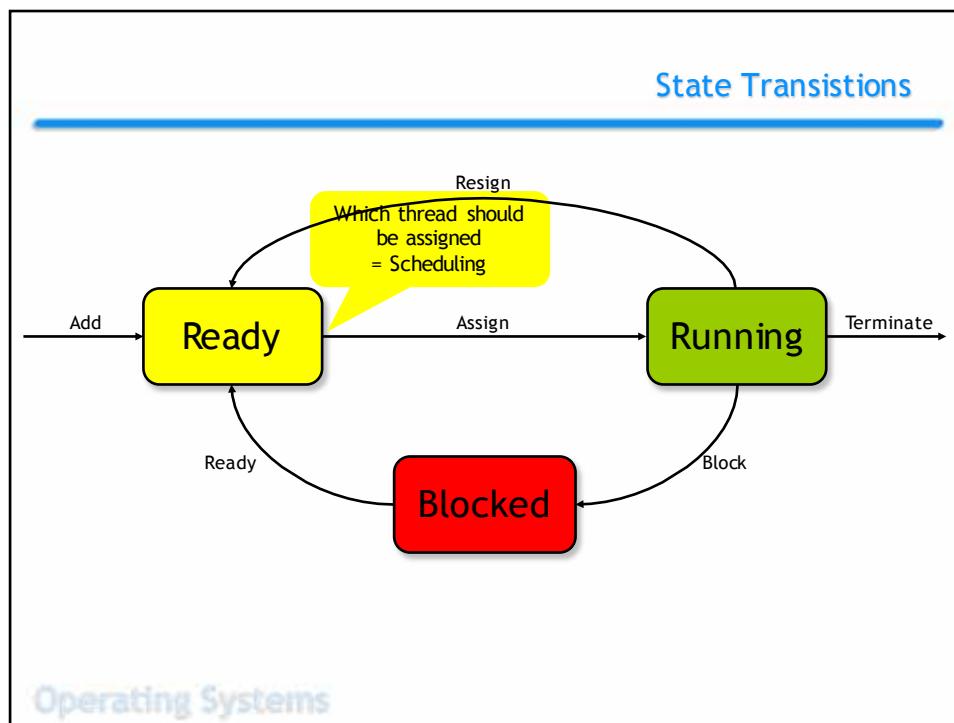
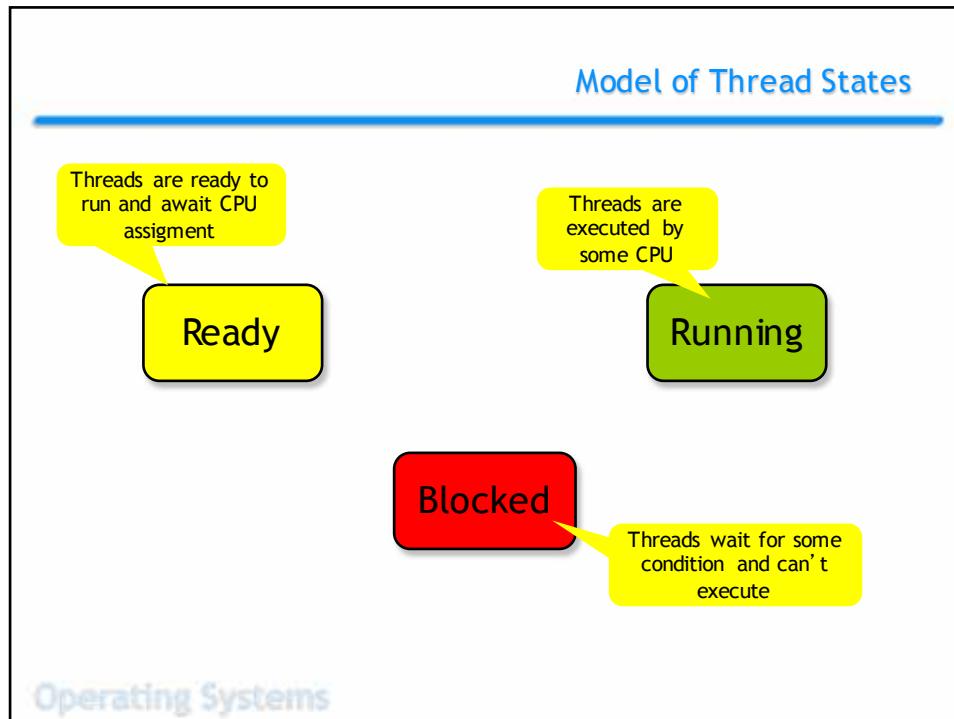
Threads
State Model

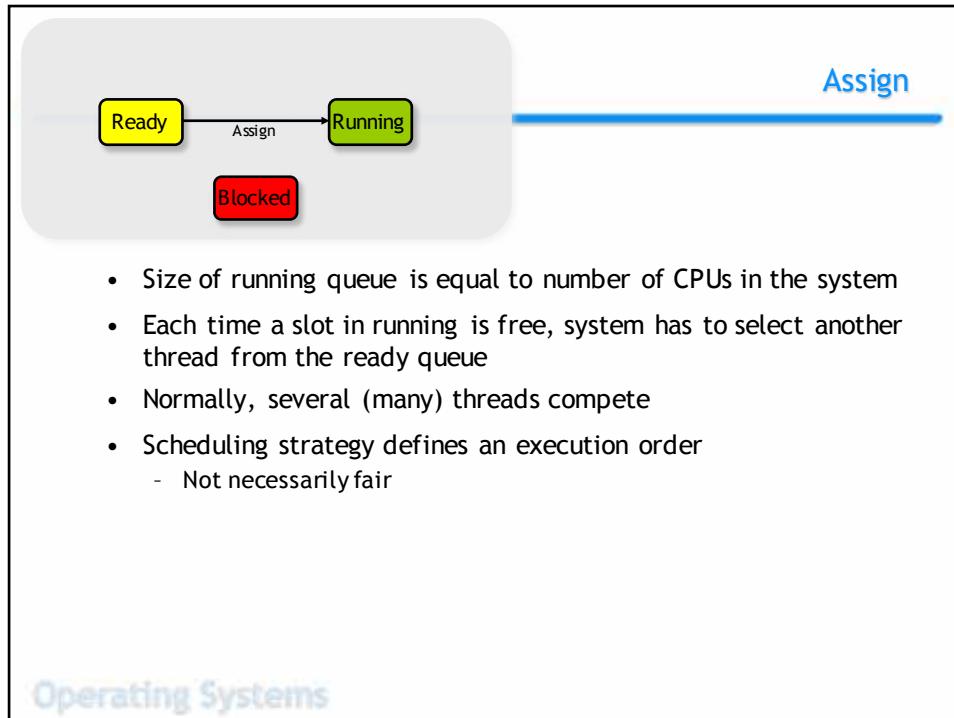
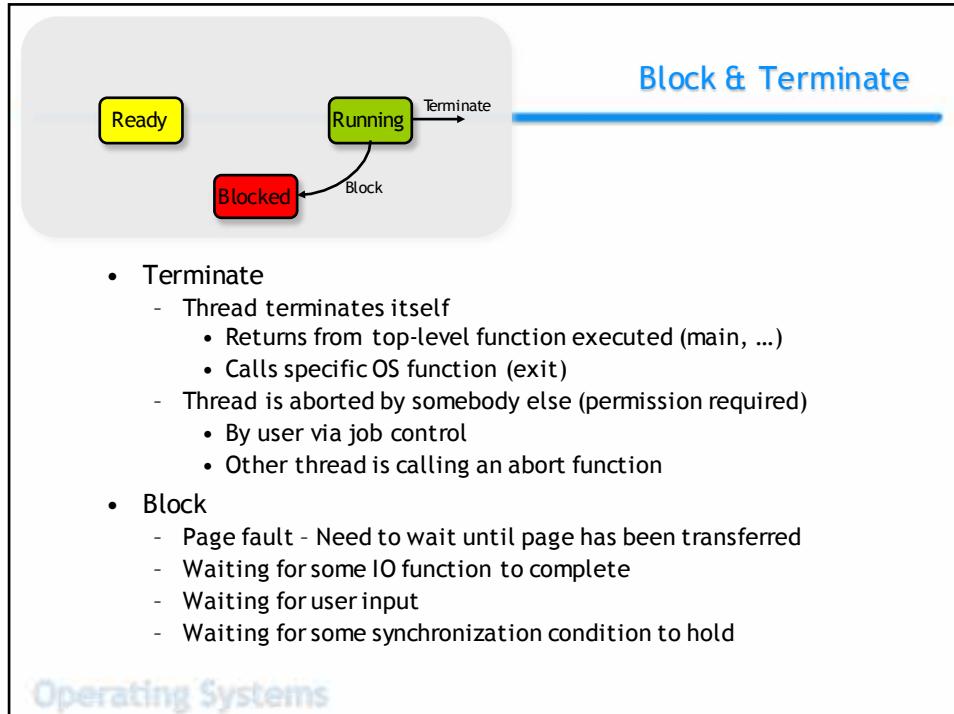
System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

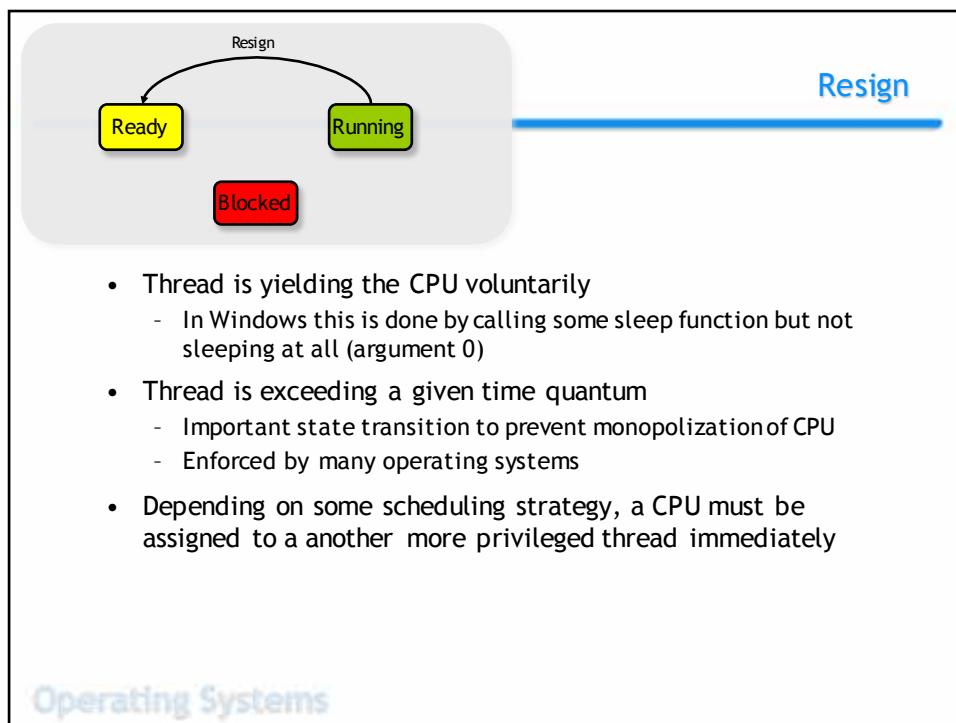
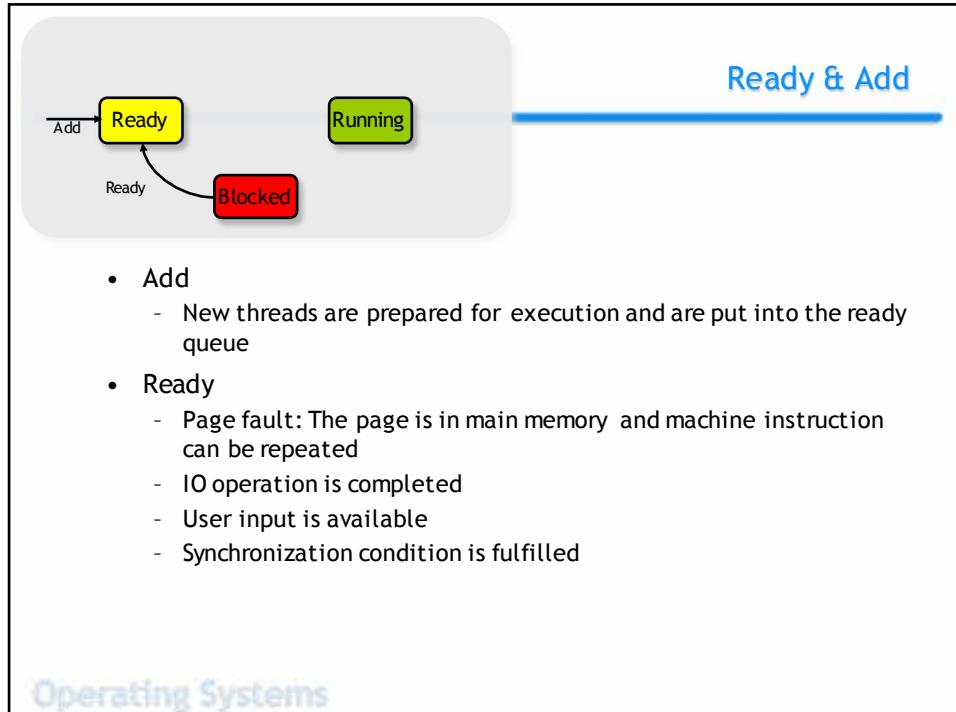
States of a Thread

- Threads can be executed by some CPU
- Threads may be blocked
 - Waiting for a memory page to be read from disk
 - Waiting for a message to arrive
 - Waiting for some user input
- Threads are willing to execute and await CPU assignment

Operating Systems







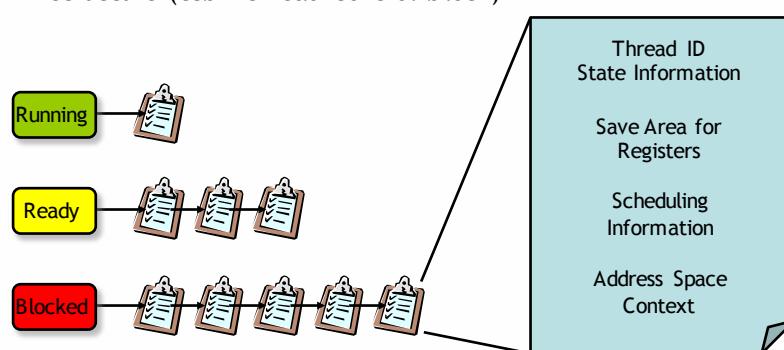
Dispatcher

- Responsible for performing all state transitions
- Several hardware-related issues have to be addressed
 - Saving and restoring registers in an orchestrated way
 - Information about address spaces may change
 - Sensible to asynchronous input such as interrupt
- As a consequence, the dispatcher consists of several assembler routines
- During the assign transition, dispatcher will access scheduler information to select appropriate thread
 - Short-term scheduler selects candidate thread fast (preferably in constant time)
 - Long-term scheduler keeps track of threads in ready queue

Operating Systems

Implementation Issues

- At least from a conceptual point of view, the three lists are implemented as such
 - Due to efficiency constraints this may be different in practice
- All required state information is stored in a thread-specific data structure (tcb = thread control block)



Operating Systems

Operating Systems

Threads
Monoprocessor Scheduling

System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

Motivation



- Scheduling strategy depends on operational field
 - Batch processing
 - Non interactive, long execution times
 - Real-time processing
 - Guarantees for execution deadlines
 - Interactive processing
 - Quick responses to users, short CPU bursts
- Hardware platform
 - Monoprocessor - Single CPU is a scarce resource
 - Multiprocessor
- Minimal number of context switches
 - Cheap but not for free (major problem earlier)
 - Beware of cold caches and TLBs

Structural Consequence

- Two-level approach
- Short-term scheduler
 - Called by dispatcher during state transition “assign”
 - Scheduling decision must be made very fast
- Long-term scheduler
 - Continuous update of state information relevant for all scheduling decisions
 - Preprocessing of expensive scheduling strategies
 - Information preparation for short-term scheduler

Operating Systems

Scheduling Criteria

- CPU utilization
 - Use maximum of CPU cycles for applications
- Throughput
 - Maximize the number of completed jobs per time interval
- Turnaround time
 - Minimize time between two consecutive job activations
- Waiting time
 - Minimize the time a thread stays in ready queue
- Response time
 - Favor interactive applications that react on user input
- Real time
 - Meet application-specific deadlines

Operating Systems

Preemptive and Non-preemptive

- Preemptive Scheduling
 - Thread will be forced to withdraw CPU
 - Reason may be exceeded time quantum or another more privileged thread must run
- Pro
 - Only way to prevent monopolization
 - Does improve response time
- Contra
 - Additional context switches
 - May add substantial overhead to total execution time of a thread
- Non-preemptive Scheduling
 - Thread yields CPU voluntarily
 - Implicitly in case of blocking operation
 - Explicitly
- Pro
 - Total execution time can be assessed (realtime systems)
- Contra
 - CPU monopolization possible
 - Fast response times can't be guaranteed in all cases

Operating Systems

First-Come, First-Served (FCFS)

- CPU assignment in the order of arrival
- Non-preemptive scheduling strategy
- Simple implementation via queue
 - Thread at queue head will be assigned next
 - Threads entering the ready queue will be appended to tail
- Average waiting time may vary



Operating Systems

The Convoi Effect of FCFS

- A single thread with a very long CPU burst may slow down the whole system
- Something you may recognize from crowded highways ☺



Operating Systems

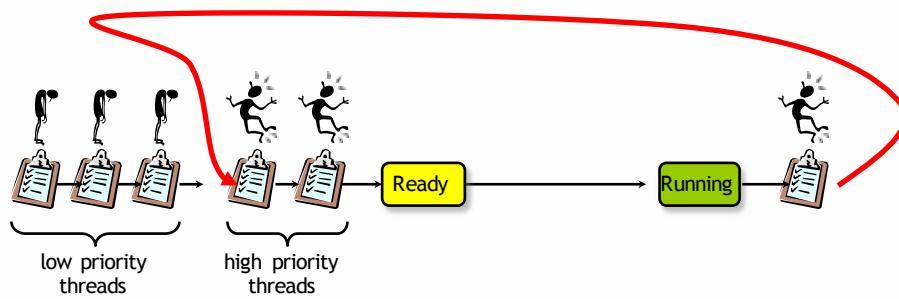
Priority-based Scheduling

- Priorities are assigned to threads
 - Explicitly: Defined by application itself, by user, or by some administrator
 - Implicitly: Defined by operating system
- Thread with highest priority will be assigned to next free CPU
- Preemptive version (common)
 - A running thread will be preempted if a thread with higher priority gets ready
 - New thread with higher priority created (add transition)
 - Previously blocked thread with higher priority is ready again (ready transition)
- Non-preemptive version
 - Thread with highest priority will be scheduled on next assign transition

Operating Systems

Starvation

- Priority-based scheduling doesn't guarantee fairness
- Example
 - Threads with high priorities are always ready to run
 - Threads with low priorities are always ready to run too



Operating Systems

Round-Robin (RR)

- RR = FCFS with preemption
 - Prevalent scheduling strategy for multitasking systems
 - Exceeding a given time quantum lead to preemption
 - Typical length of time quantum between 10 ms and 20 ms
- Size of time quantum critical for overall system performance
 - too large: FCFS (most threads block before their quantum exceeds)
 - too small: too many additional context switches
- Average waiting time
 - Depends on number of threads
 - n threads and time quantum tq

$$Waittime_{Average} \leq tq * (n - 1)$$

- Length of CPU burst has no influence



Operating Systems

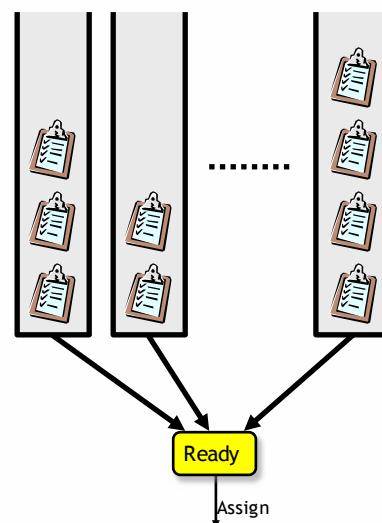
Threads Doing IO are Penalized

- RR favors threads with long CPU burst
 - exhaust their time quantum
- Threads doing IO yield CPU before time quantum exceeded
- Possible solution
 - Multi-level scheduling
 - Additional queue for threads not exhausting their quantum
 - Hierarchy of two scheduling strategies
 - 1st Level: Which queue to select next?
 - 2nd Level: Order within queue?

Operating Systems

Multi-Level Scheduling

- Ready queue partitioned into several sub-queues
 - Address different scheduling goals
- 1st level strategy
 - Which sub-queue to select?
 - Priorities are very common
- 2nd level strategy
 - Which thread within queue will get next CPU
 - Depending on scheduling goal, different strategies are used

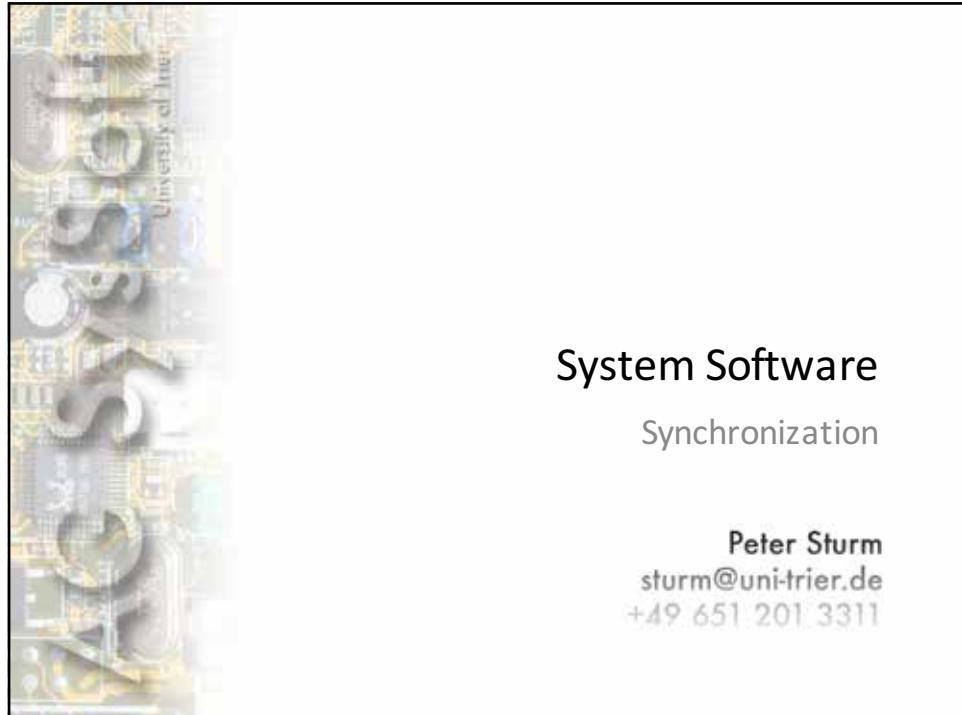


Operating Systems

Feedback Scheduling

- Many scheduling strategies take the past into account
 - SJF approximation makes scheduling decision on the length of previous CPU bursts
 - In interactive environments, threads exhausting their time quantum should be penalized
 - The application currently holding the UI focus should be favorized (priority boost done by Windows operating systems)
 - Aging of priorities in order to minimize risk of starvation
 - ...
- Requires additional tracking of all relevant state information
 - Stored within thread control block
- Feedback may lead to complex and unpredictable behavior

Operating Systems

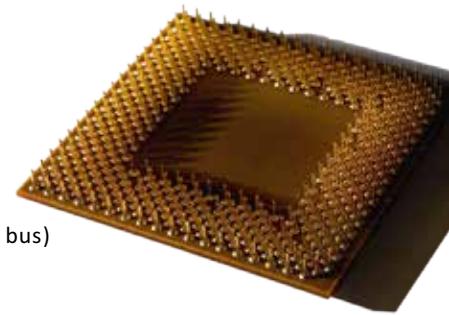


The slide is titled 'Das Kontinuum' (The Continuum). It features a horizontal orange bar divided into five horizontal segments. Two stylized black stick figures are running across the bar. The first figure is carrying a small object on its back. The text below the bar lists two main categories: 'Aktivität (Zeit)' and 'Raum', each with associated sub-points.

- Aktivität (Zeit)
 - Abfolge ausgeführter Instruktionen
 - Kontrollfluss oder Thread
 - Unterstützung für mehrere Threads
 - Concurrency Control
- Raum
 - Menge der adressierbaren Speicherzellen
 - Adressraum oder "Adress space"
 - Unterstützung für mehrere Adressräume

Adreßraum

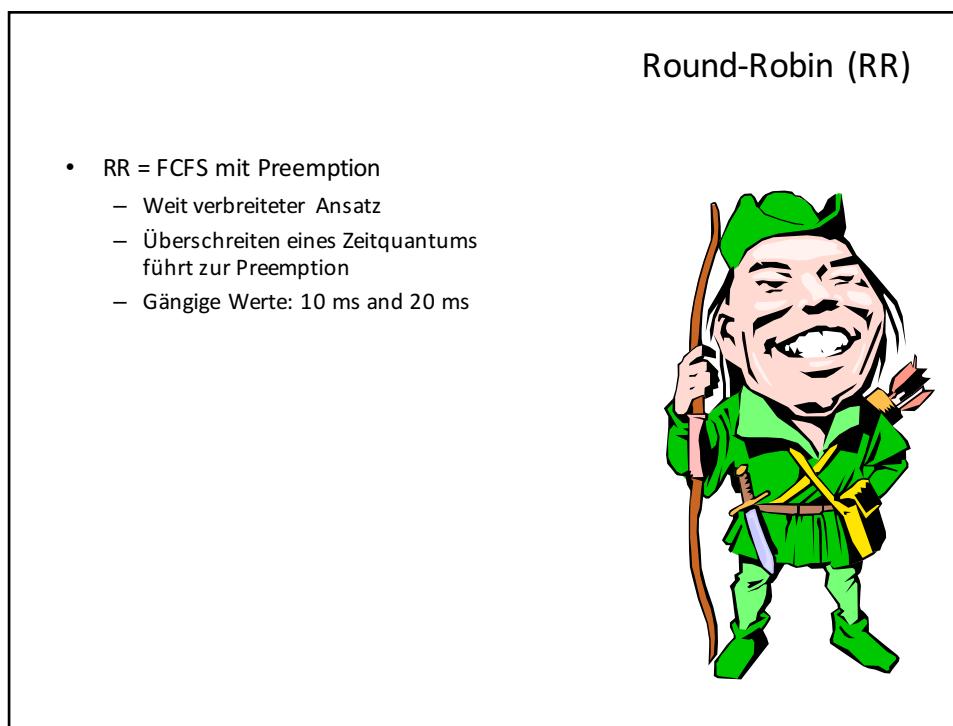
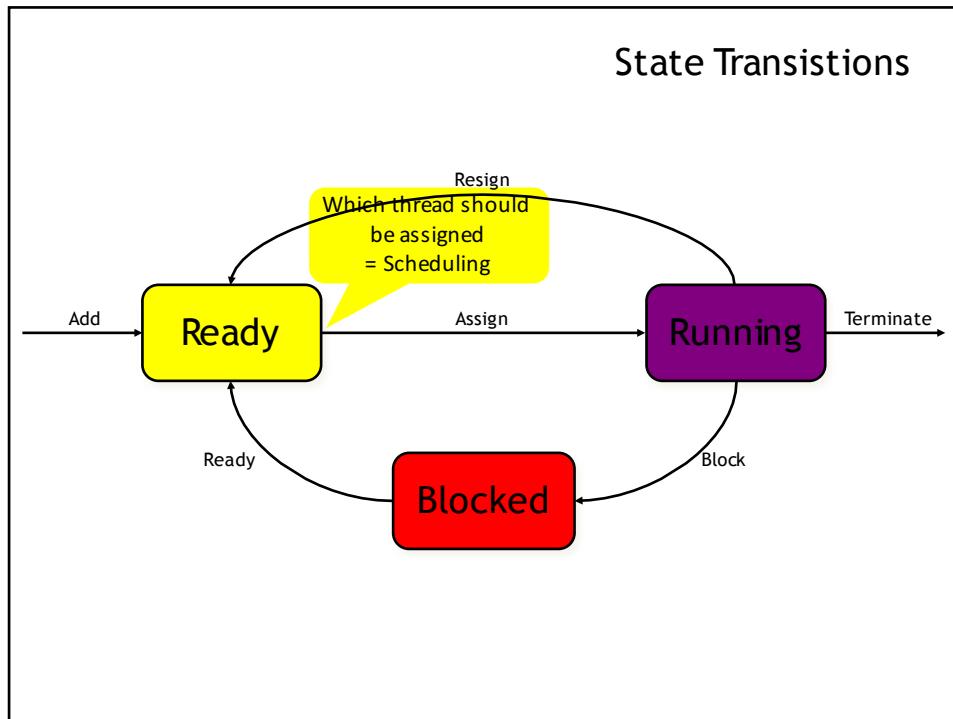
- CPU definiert ein theoretisches, oberes Limit
 - Anzahl der Adreß-Pins (address bus)
- Beispiele
 - 32 bit Architektur
(Intel Pentium, AMD Athlon XP, ...)
 - 32 Adressleitungen = 4 GByte Adreßraum
 - 64 bit Architektur
(AMD Athlon 64, Intel Itanium, ...)
 - 64 Leitungen = 4194304 TByte theoretischer Adreßraum
 - Real limitiert auf ~40 Pins (Server 48 Pins)

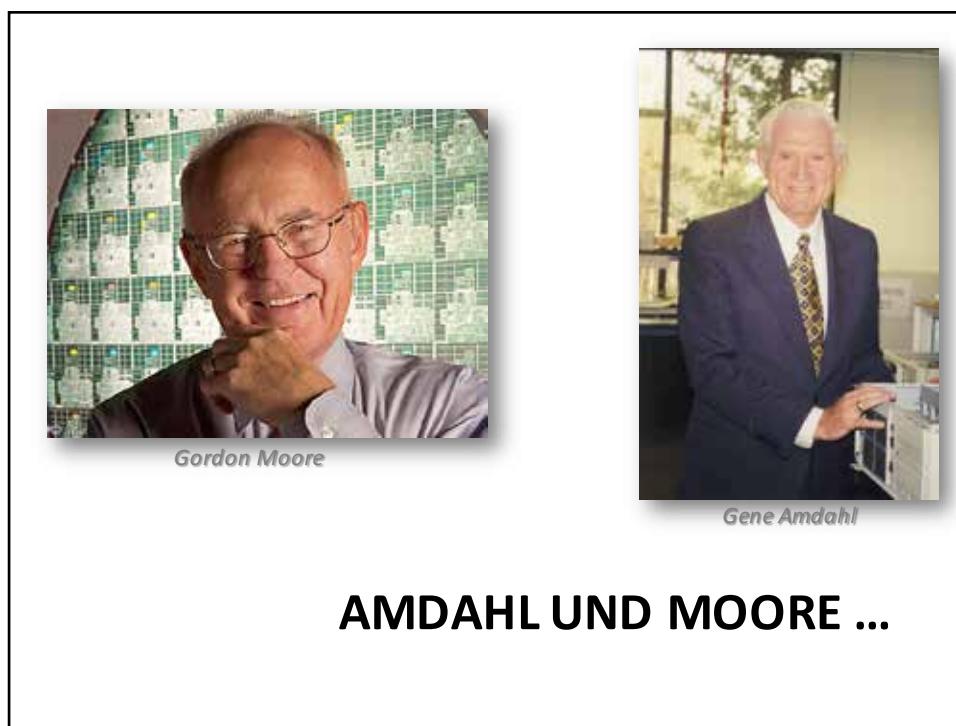
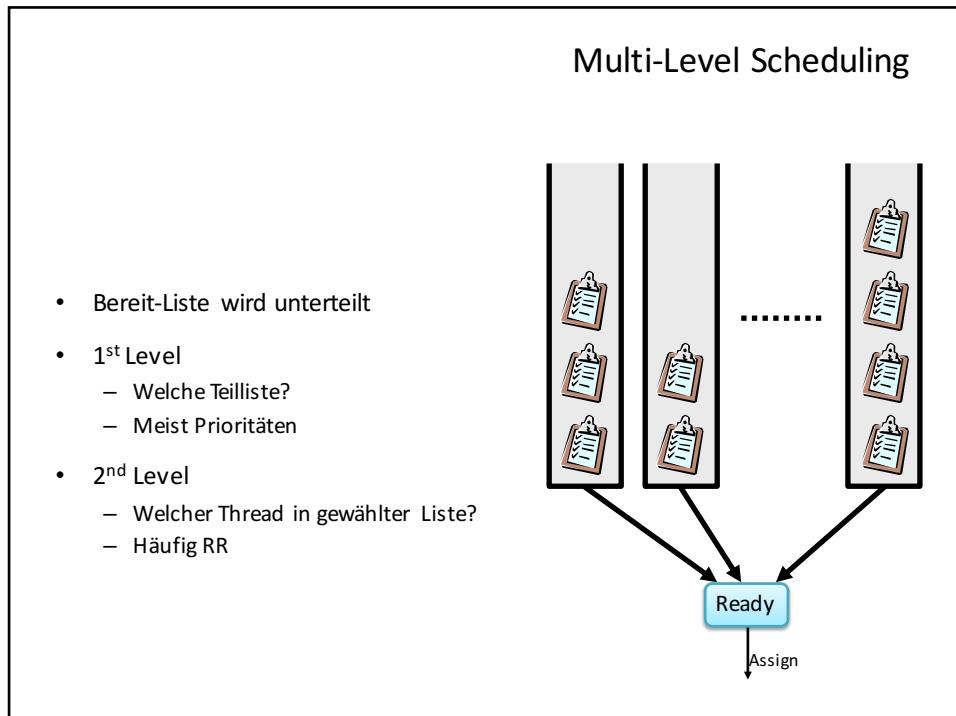


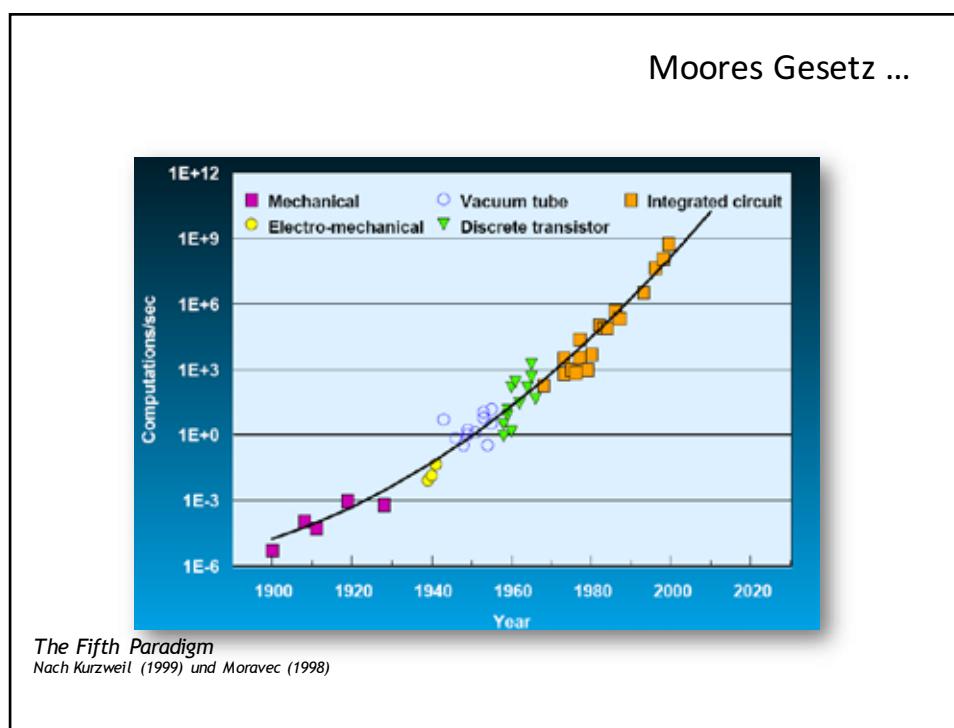
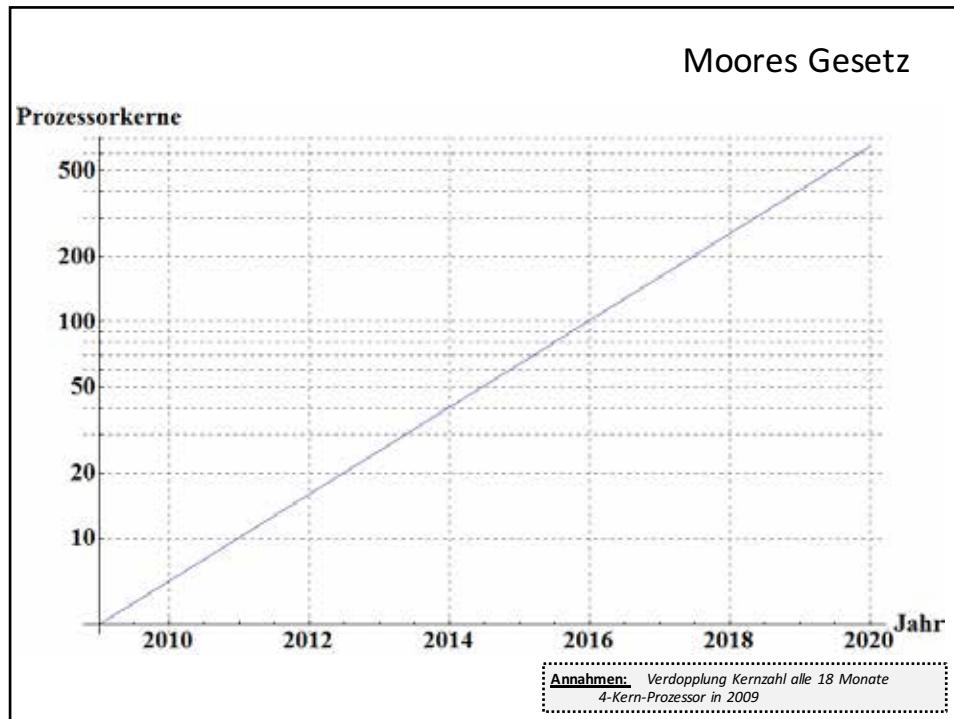
Virtuelle Prozessoren

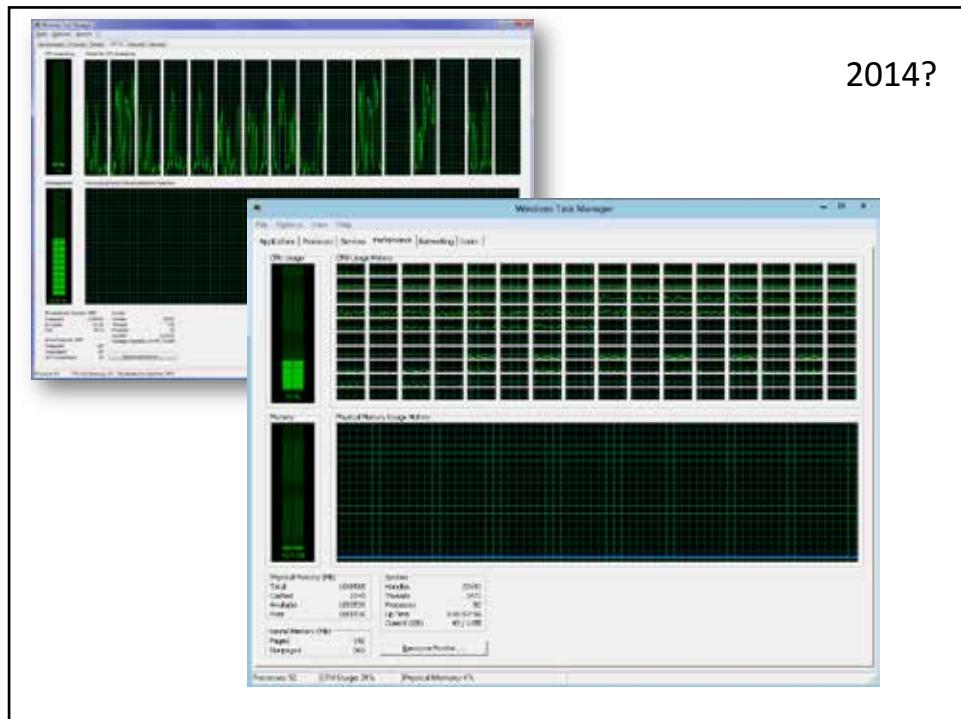
Threads







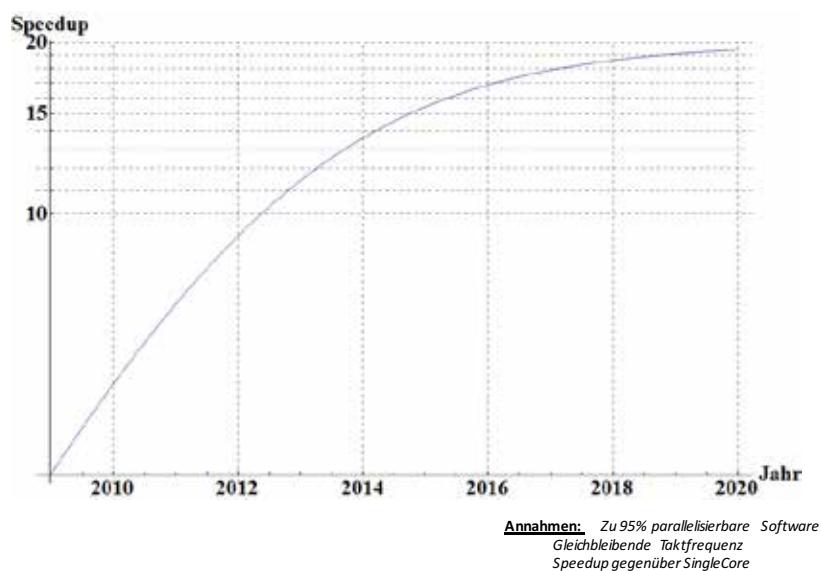


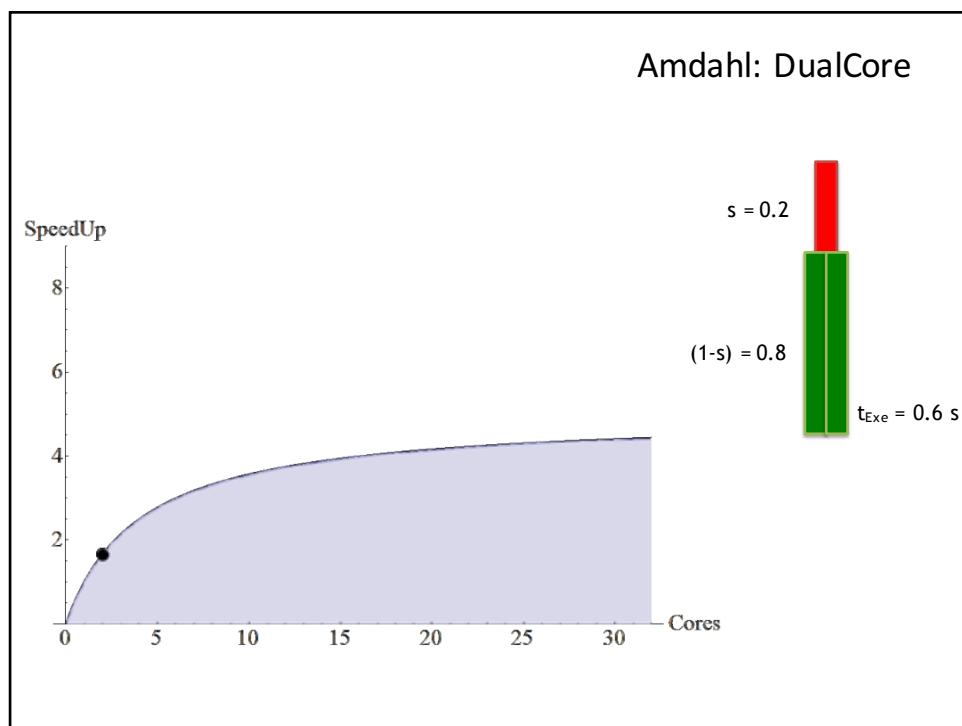
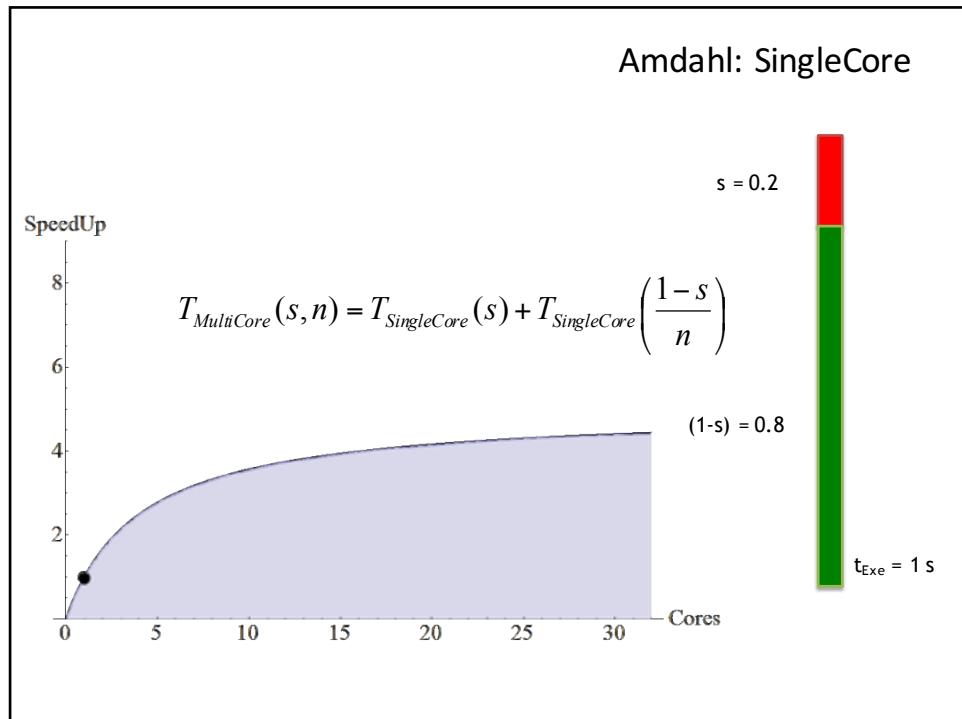


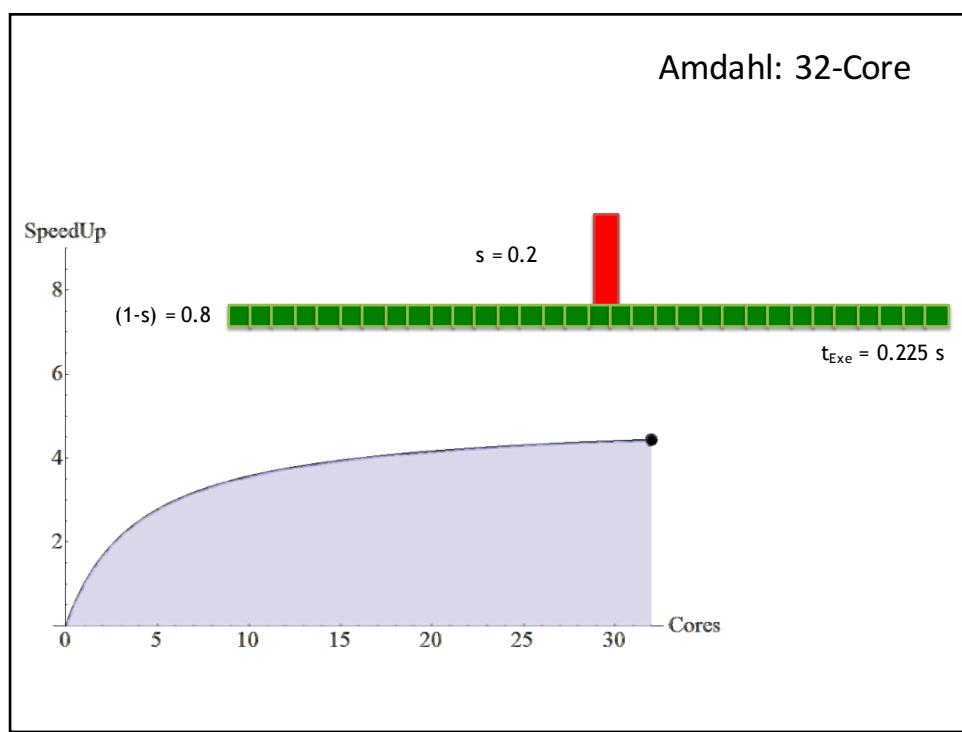
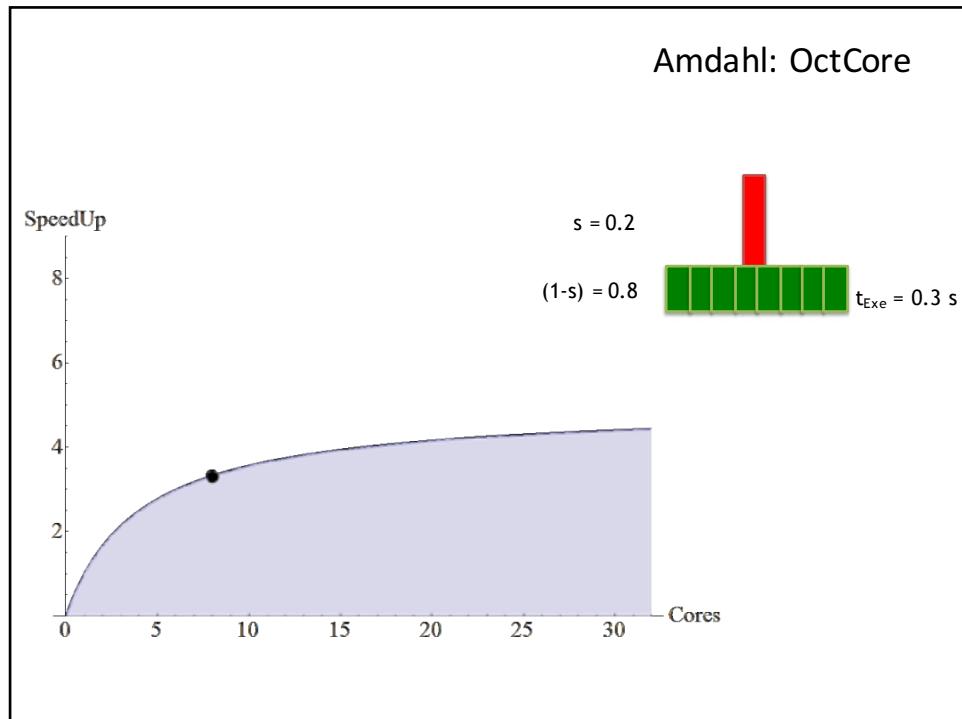
-
- Weiter?
- 2014
 - 64 Cores
 - 2018
 - 1024 Cores
 - 2028
 - 8192 Cores
 - Many Cores
 - ...
 - Enough Cores ☺

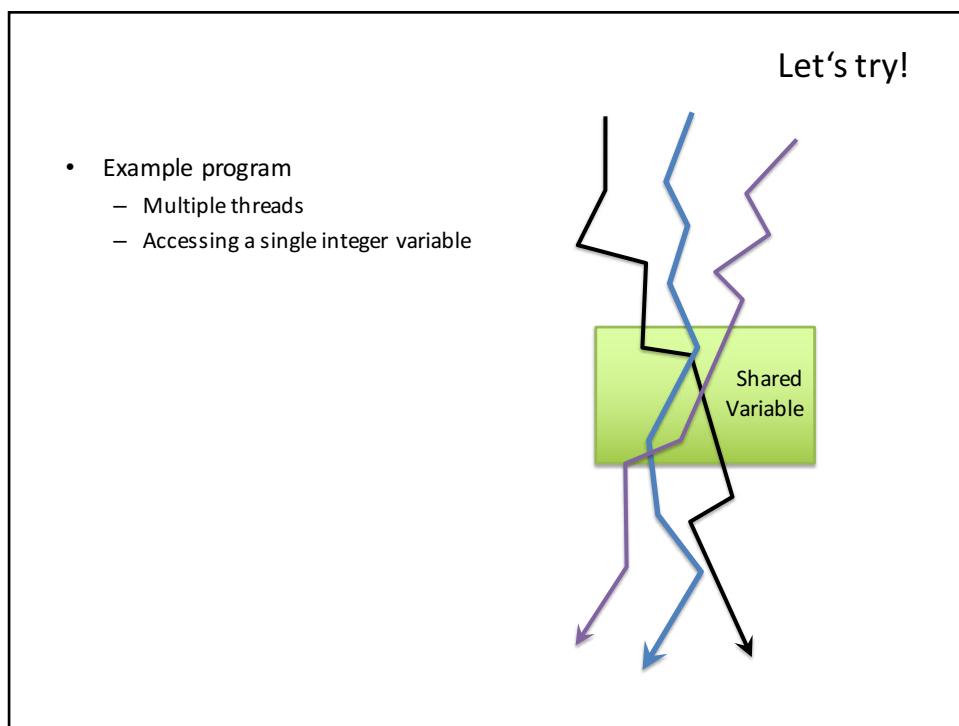
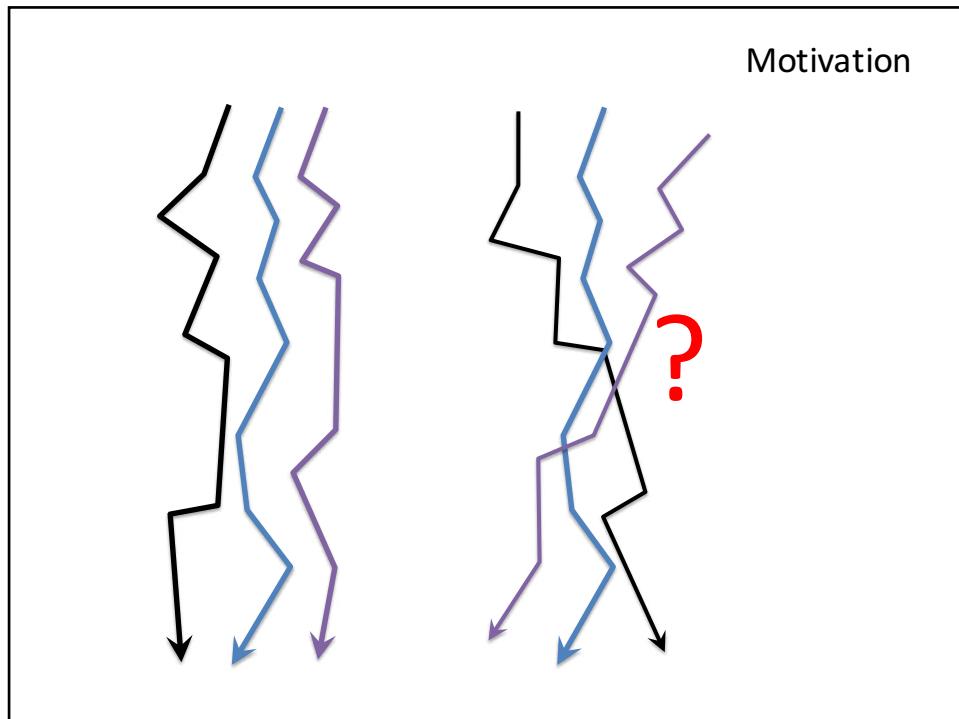
**Nicht-parallele
Software wird nie
mehr schneller
laufen als heute!**

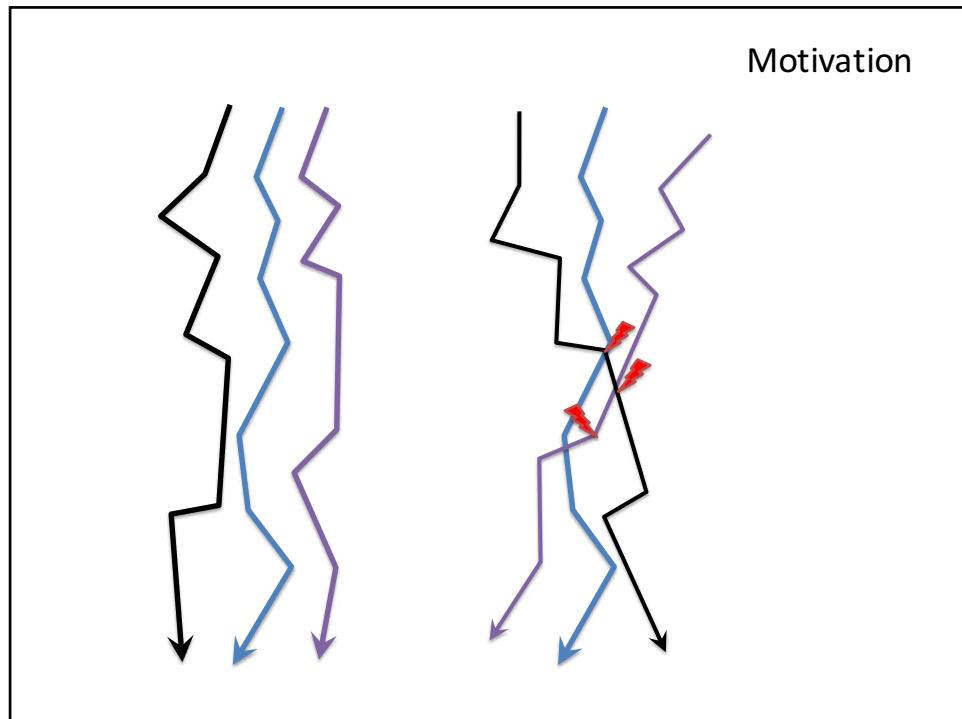
Amdahls Gesetz









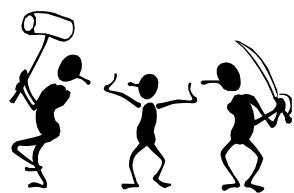




Konkurrenz und Kooperation

Konkurrenz

- Kampf um Ressourcen
 - Ziel: Alles für mich
- Implizite Synchronisation
 - Virtuelle Ressourcen
 - OS serialisiert
- Wechselseitiger Ausschluß



Kooperation

- Kooperierende Threads
- Teilt sich Ressourcen & Aufgaben
- Ziele
 - Speedup
 - Reaktionsgeschwindigkeit
- Beliebig komplexe Synchronisationsprobleme

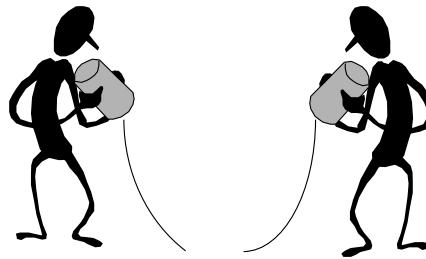


Mutual Exclusion

- Wechselseitiger Ausschluß
- Atomanität für eine Gruppe von Befehlen
- Höhere Synchronisationslösungen ableitbar
- Paßt zum Scheduling
 - Blockiert-Liste wenn belegt

IPC

- Interprozeßkommunikation
 - Synchronisation
 - A wartet auf Bedingung
 - B erfüllt Bedingung und informiert A
 - 1-Bit-Nachricht
 - Kommunikation
 - Austausch von Daten
- Vielfältige Verfahren

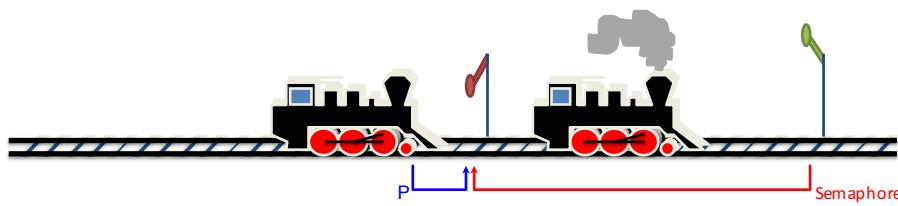


Abstand

- In einem Prozeß (Adreßraum)
 - Zugriff auf gemeinsamen Speicher
- Zwischen Prozessen auf einem Rechner
- Zwischen Prozessen auf verschiedenen Rechnern
 - Verteilte Systeme
 - Internet, IP, TCP, UDP

SEMAPHORE

Semaphore



- Fundamental synchronization primitive
- Two basic functions (on a semaphore s)
 - s.P()
 - May block in case semaphore already “occupied”
 - s.V()
 - Never blocks; releases semaphore and may free a blocked thread
- Dijkstra (1968), “THE” Multiprogramming System
 - P = passeren (request entry)
 - V = vrygeven (release)

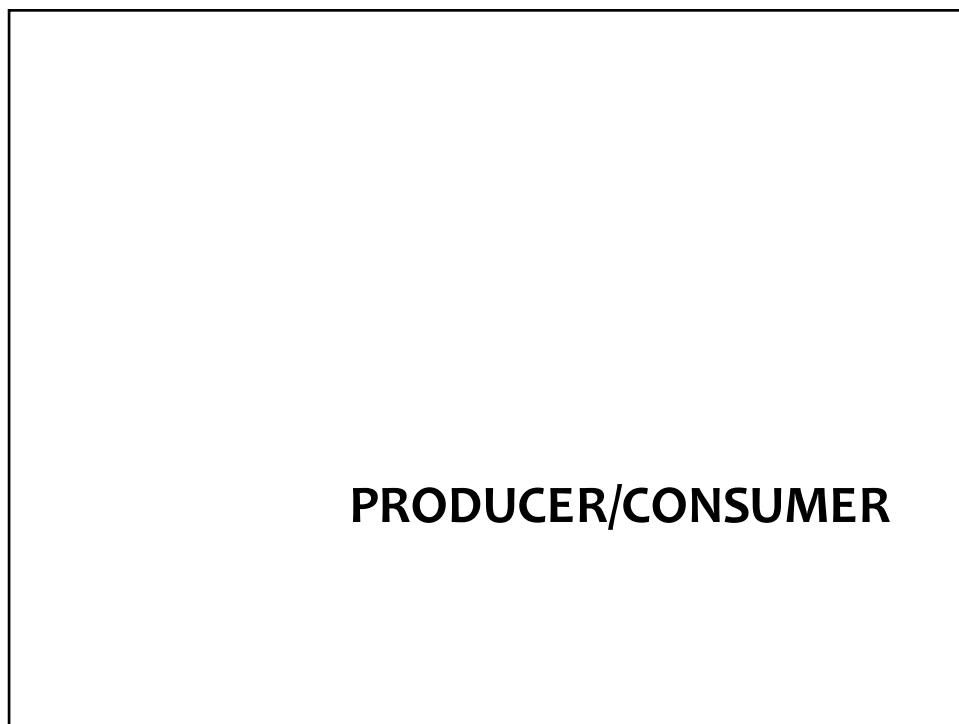
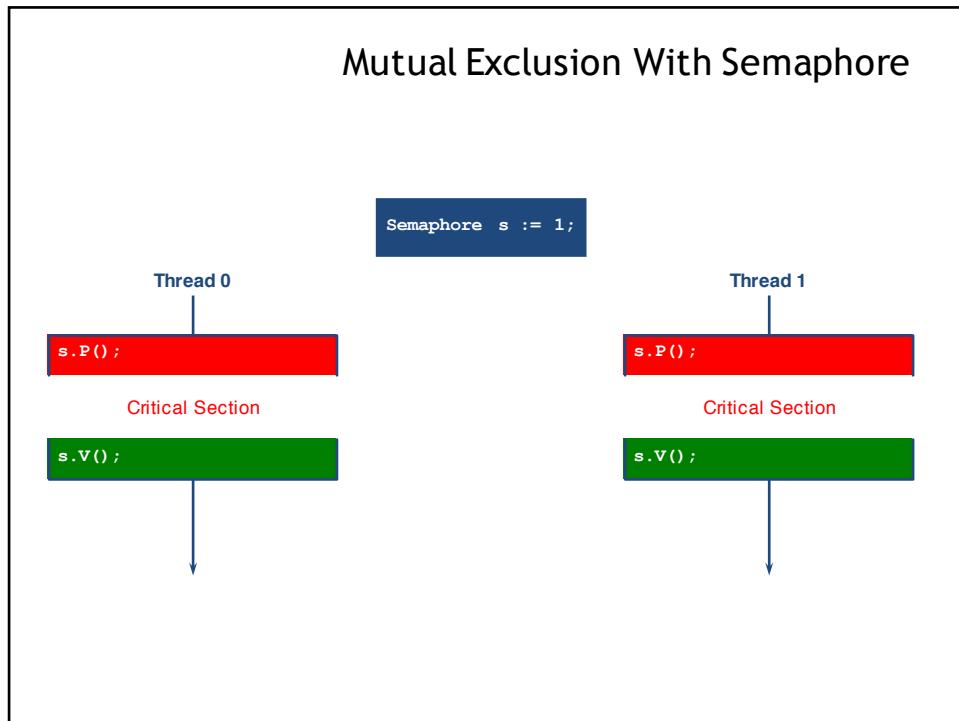
Semantic

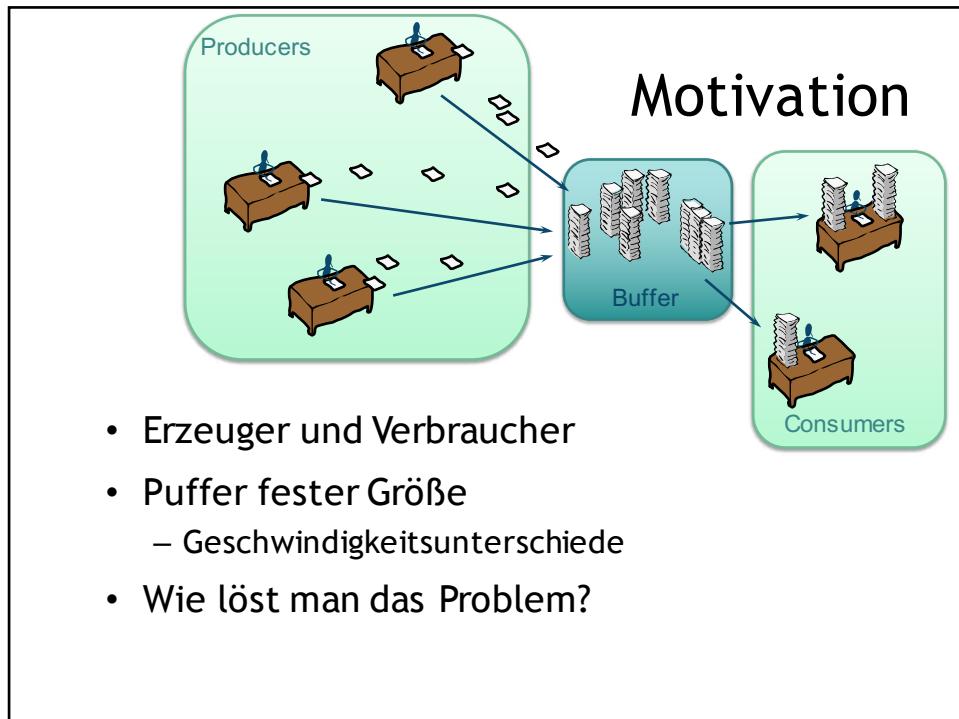
- Semaphores are counting signals
 - s.P(): Thread waits for some signal s
 - s.V(): Thread sends a signal s
- Semaphor = Integer (s.value)
 - Initialized with a s.value ≥ 0

```
void P () {
    s.value--;
    if (s.value < 0)
        Block calling thread;
}
```

Of course, P and V themselves are critical sections and must be protected!

```
void V () {
    s.value++;
    if (s.value < 1)
        Put a thread waiting in s into ready queue;
}
```





Outline of Buffer

```
public class Buffer
{
    public Buffer ( int n )
    {
    }

    /// Inserts another good into the buffer.
    /// May block until free space is available.
    public void Produce ( int good )
    {
    }

    /// Consumes a good stored inside the buffer.
    /// May signal blocked producer threads.
    public int Consume ()
    {
        return 42;
    }
}
```

Producer (C#)

```
public class Producer
{
    public Producer ( Buffer b )
    {
        buffer = b;
        my_id = this.GetHashCode();
        ThreadStart ts = new ThreadStart(Run);
        my_thread = new Thread(ts);
        my_thread.Start();
    }

    private void Run ()
    {
        Console.WriteLine("Producer {0}: started ...",my_id);
        int good = this.GetHashCode() * 1000000;
        while (true)
        {
            buffer.Produce(good);
            Console.WriteLine("Producer {0}: good {1} stored",my_id,good);
            good++;
        }
    }

    private Buffer buffer;
    private Thread my_thread;
    private int my_id;
}
```

Consumer (C#)

```
public class Consumer
{
    public Consumer ( Buffer b )
    {
        buffer = b;
        my_id = this.GetHashCode();
        ThreadStart ts = new ThreadStart(Run);
        my_thread = new Thread(ts);
        my_thread.Start();
    }

    private void Run ()
    {
        Console.WriteLine("Consumer {0}: started ...",my_id);
        while (true)
        {
            int good = buffer.Consume();
            Console.WriteLine("Consumer {0}: good {1} retrieved",my_id,good);
        }
    }

    private Buffer buffer;
    private Thread my_thread;
    private int my_id;
}
```

```

public class Buffer
{
    public Buffer ( int n )
    {
        this.n = n;
        slots = new int[n];
        mutex = new Semaphore (1);
        slots_available = new Semaphore (n) ;
        goods_available = new Semaphore (0) ;
    }

    ...

    private int n;
    private int [] slots;
    private int free = 0;
    private int used = 0;
    private Semaphore mutex;
    private Semaphore slots_available;
    private Semaphore goods_available;
}

```

The Buffer (C#)

```

public void Produce ( int good )
{
    slots_available.P();
    mutex.P();
    slots[free] = good;
    free = (free+1) % n;
    mutex.V();
    goods_available.V();
}

public int Consume ()
{
    goods_available.P();
    mutex.P();
    int good = slots[used];
    used = (used+1) % n;
    mutex.V();
    slots_available.V();
    return good;
}

```

```

public class Buffer
{
    public Buffer ( int n )
    {
        this.n = n;
        slots = new int[n];
        mutex_p = new Semaphore(1);
        mutex_c = new Semaphore(1);
        slots_available = new Semaphore (n) ;
        goods_available = new Semaphore (0) ;
    }

    ...

    private int n;
    private int [] slots;
    private int free = 0;
    private int used = 0;
    private Semaphore mutex_p, mutex_c;
    private Semaphore slots_available;
    private Semaphore goods_available;
}

```

The Optimal Buffer (C#)

```

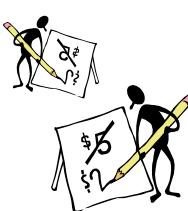
public void Produce ( int good )
{
    slots_available.P();
    mutex_p.P();
    slots[free] = good;
    free = (free+1) % n;
    mutex_p.V();
    goods_available.V();
}

public int Consume ()
{
    goods_available.P();
    mutex_c.P();
    int good = slots[used];
    used = (used+1) % n;
    mutex_c.V();
    slots_available.V();
    return good;
}

```

READER/WRITER

The Problem



- A shared data structure is used by multiple threads
- Writer threads
 - These threads modify the shared data and require mutual exclusion
- Reader threads
 - Since readers don't modify data, they can access the shared data structure simultaneously

1. Mutual Exclusion

- As a starting point, implement a simple solution based on mutual exclusion for readers and writers. Of course, this version doesn't exploit the potential for parallelism by allowing multiple readers to access the shared data simultaneously.
- [Start from scratch \(C#\)](#)
- Inspect the [complete example program \(C#\)](#)

Works as Expected

```
F:\OCP - Open Curriculum Project\OS Episodes\31 - IPC Reader Writer with Semaphores\r...
Episode_31.Reader average waiting time 0.0001592371599359
Episode_31.Reader average waiting time 0.00021244526155591
Episode_31.Reader average waiting time 0.00021244526155591
Episode_31.Writer average waiting time 0.8794661949925124
Press any key to continue
```

- But only one reader inside sanctum at any given time

Starting Point: Mutual Exclusion

```
Semaphore Sanctum = new Semaphore(1);
```

Shared Data

```
while (true) {  
    Sanctum.P();  
    // Change data  
    Sanctum.V();  
}
```

Writer

```
while (true) {  
    Sanctum.P();  
    // Read data  
    Sanctum.V();  
}
```

Reader

2. Reader Preference

- Improve the previous program to enable multiple readers to access shared data concurrently
- Hint
 - Only the first reader must acquire semaphore
 - Subsequent readers just enter
 - Only the last reader leaving must release semaphore
- We accept the risk of starvation for writers
- This is a still [faulty version of the program ☺](#)

Reader Preference (Faulty)

```
Semaphore Sanctum = new Semaphore(1);
int readers_inside = 0;
```

```
Shared Data

while (true) {
    Sanctum.P();
    // Change data
    Sanctum.V();
}
```

Writer

```
while (true) {
    if (readers_inside == 0)
        Sanctum.P();
    readers_inside++;
    // Read data
    readers_inside--;
    if (readers_inside == 0)
        Sanctum.V();
}
```

Reader

Why Does It Fail?

```
while (true) {
    if (readers_inside == 0)
        Sanctum.P();
    readers_inside++;
    // Read data
    readers_inside--;
    if (readers_inside == 0)
        Sanctum.V();
}
```

Reader

- Multiple readers may access the shared variable `readers_inside` concurrently
- Testing and setting variable must be atomic among readers
- Mutual exclusion required

2. Reader Preference (Correct)

```
Shared Data      Semaphore Sanctum = new Semaphore(1);  
                  Semaphore RMutex = new Semaphore(1);  
                  int readers_inside = 0;
```

```
while (true) {  
    Sanctum.P();  
    // Change data  
    Sanctum.V();  
}
```

Writer

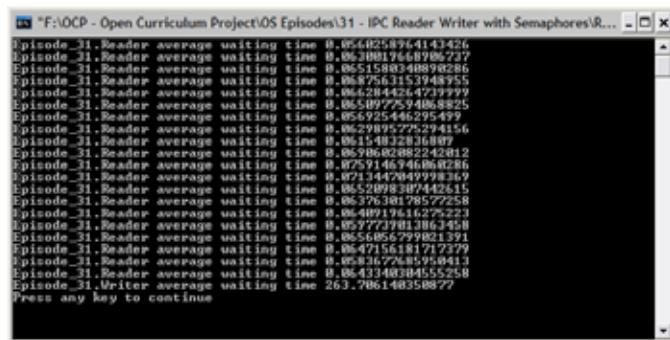
```
while (true) {  
    RMutex.P();  
    if (readers_inside == 0)  
        Sanctum.P();  
    readers_inside++;  
    RMutex.V();  
    // Read data  
    RMutex.P();  
    readers_inside--;  
    if (readers_inside == 0)  
        Sanctum.V();  
    RMutex.V();  
}
```

Reader

Obviously Reader Preference

```
F:\OCP - Open Curriculum Project\OS Episodes\31 - IPC Reader Writer with Semaphores\&...  
Episode_31.Reader average waiting time 0.0418106922040472  
Episode_31.Reader average waiting time 0.0523497953217256  
Episode_31.Reader average waiting time 0.0523497953217256  
Episode_31.Writer average waiting time 0.24223242386724  
Episode_31.Writer average waiting time 0.261152651181742  
Episode_31.Writer average waiting time 0.242290652086607  
Press any key to continue
```

But Writer May Starve



```
"F:\OCP - Open Curriculum Project\OS Episodes\31 - IPC Reader Writer with Semaphores\R... <input type="button" value="X">
Episode_31.Reader average waiting time 0.06642589764143426
Episode_31.Reader average waiting time 0.0630817668706237
Episode_31.Reader average waiting time 0.065158834889826
Episode_31.Reader average waiting time 0.0642844264733999
Episode_31.Reader average waiting time 0.0658972594068825
Episode_31.Reader average waiting time 0.066925446295499
Episode_31.Reader average waiting time 0.0629895775294156
Episode_31.Reader average waiting time 0.06154812836897
Episode_31.Reader average waiting time 0.065868289242812
Episode_31.Reader average waiting time 0.0657518646162616
Episode_31.Reader average waiting time 0.0713447949998369
Episode_31.Reader average waiting time 0.0652998307442615
Episode_31.Reader average waiting time 0.0637638178577258
Episode_31.Reader average waiting time 0.0649919616275223
Episode_31.Reader average waiting time 0.0577739813863450
Episode_31.Reader average waiting time 0.0656862791821391
Episode_31.Reader average waiting time 0.0656862791821391
Episode_31.Reader average waiting time 0.0583672685958013
Episode_31.Reader average waiting time 0.0641348034555258
Episode_31.Writer average waiting time 263.796148358877
Press any key to continue
```

- It takes a long time for a single writer, if there are 20 readers continuously accessing shared data
- Sometimes forever ☺

3. Writer Preference

- The moment, a single writer is requesting to enter the shared data section, no other reader might enter anymore
- The writer has to wait, until all readers inside left the sanctum
- Subsequent writers are preferred
- We accept a possible starvation among readers
- Hint
 - Similar to reader preference, we count the number of writers willing to enter
 - Another semaphore locks readers in case of writer interest

3. Writer Preference

```

Semaphore Sanctum = new Semaphore(1);
Semaphore RMutex = new Semaphore(1);
Semaphore WMutex = new Semaphore(1);
Semaphore PreferWriter = new Semaphore(1);
int readers_inside = 0;
int writers_interested = 0;

while (true) {
    WMutex.P();
    if (writers_interested == 0)
        PreferWriter.P();
    writers_interested++;
    WMutex.V();
    Sanctum.P();
    // Change data
    Sanctum.V();
    WMutex.P();
    writers_interested--;
    if (writers_interested == 0)
        PreferWriter.V();
    WMutex.V();
}

```

Shared Data

Writer

```

while (true) {
    PreferWriter.P();
    RMutex.P();
    if (readers_inside == 0)
        Sanctum.P();
    readers_inside++;
    RMutex.V();
    PreferWriter.V();
    // Read data
    RMutex.P();
    readers_inside--;
    if (readers_inside == 0)
        Sanctum.V();
    RMutex.V();
}

```

Reader

Works Great

... and Writer Doesn't Starve

```
F:\OCP - Open Curriculum Project\OS Episodes\31 - IPC Reader Writer with Semaphores\R...
Episode_31.Reader average waiting time 0.418289949079648
Episode_31.Reader average waiting time 0.355864852594586
Episode_31.Reader average waiting time 0.429141661921241
Episode_31.Reader average waiting time 0.351578142829192
Episode_31.Reader average waiting time 0.437671282615524
Episode_31.Reader average waiting time 0.351578142829192
Episode_31.Reader average waiting time 0.42874939632885
Episode_31.Reader average waiting time 0.373580548588364
Episode_31.Reader average waiting time 0.434237057817711
Episode_31.Reader average waiting time 0.356813837351911
Episode_31.Reader average waiting time 0.356813837351911
Episode_31.Reader average waiting time 0.369429810521338
Episode_31.Reader average waiting time 0.443978264317181
Episode_31.Reader average waiting time 0.349584185278389
Episode_31.Reader average waiting time 0.472281698239524
Episode_31.Reader average waiting time 0.369181676463732
Episode_31.Reader average waiting time 0.435798535877862
Episode_31.Reader average waiting time 0.35561823622847
Episode_31.Reader average waiting time 0.459728754862149
Episode_31.Reader average waiting time 0.35561823622847
Episode_31.Writer average waiting time 0.643866387968945
Press any key to continue...
```

4. Writer Preference (2nd Edition)

- Any ideas to improve the program with writer preference even further?
- Hint
 - A single writer may compete with multiple readers when calling `PreferWriter.P()`
 - Is it possible to limit that to at most one reader?

4. Writer Preference (Improved)

```

Shared Data
Semaphore Sanctum = new Semaphore(1);
Semaphore RMutex = new Semaphore(1);
Semaphore WMutex = new Semaphore(1);
Semaphore PreferWriter = new Semaphore(1);
Semaphore ReaderQueue = new Semaphore(1);
int readers_inside = 0;
int writers_interested = 0;

while (true) {
    WMutex.P();
    if (writers_interested == 0)
        PreferWriter.P();
    writers_interested++;
    WMutex.V();
    Sanctum.P();
    // Change data
    Sanctum.V();
    WMutex.P();
    writers_interested--;
    if (writers_interested == 0)
        PreferWriter.V();
    WMutex.V();
}

```

Writer

```

while (true) {
    ReaderQueue.P();
    PreferWriter.P();
    ReaderQueue.V();
    RMutex.P();
    if (readers_inside == 0)
        Sanctum.P();
    readers_inside++;
    RMutex.V();
    PreferWriter.V();

    // Read data
    RMutex.P();
    readers_inside--;
    if (readers_inside == 0)
        Sanctum.V();
    RMutex.V(); Reader
}

```

4. Writer Preference (Improved but faulty)

```

Shared Data
Semaphore Sanctum = new Semaphore(1);
Semaphore RMutex = new Semaphore(1);
Semaphore WMutex = new Semaphore(1);
Semaphore PreferWriter = new Semaphore(1);
Semaphore ReaderQueue = new Semaphore(1);
int readers_inside = 0;
int writers_interested = 0;

while (true) {
    WMutex.P();
    if (writers_interested == 0)
        PreferWriter.P();
    writers_interested++;
    WMutex.V();
    Sanctum.P();
    // Change data
    Sanctum.V();
    WMutex.P();
    writers_interested--;
    if (writers_interested == 0)
        PreferWriter.V();
    WMutex.V();
}

```

Writer

```

while (true) {
    ReaderQueue.P();
    PreferWriter.P();
    RMutex.P();
    if (readers_inside == 0)
        Sanctum.P();
    readers_inside++;
    RMutex.V();
    ReaderQueue.V();
    PreferWriter.V();

    // Read data
    RMutex.P();
    readers_inside--;
    if (readers_inside == 0)
        Sanctum.V();
    RMutex.V(); Reader
}

```

Works Even Better

```
F:\OCP - Open Curriculum Project\OS Episodes\31 - IPC Reader Writer with Semaphores\R... □ x
Episode_31.Reader average waiting time 2.55761086177943
Episode_31.Reader average waiting time 2.23687177330562
Episode_31.Reader average waiting time 3.2948696319469
Episode_31.Reader average waiting time 1.6325648492265
Episode_31.Writer average waiting time 0.8137261825693348
Episode_31.Writer average waiting time 0.82727693482285
Press any key to continue
```

... and Writer Still Doesn't Starve

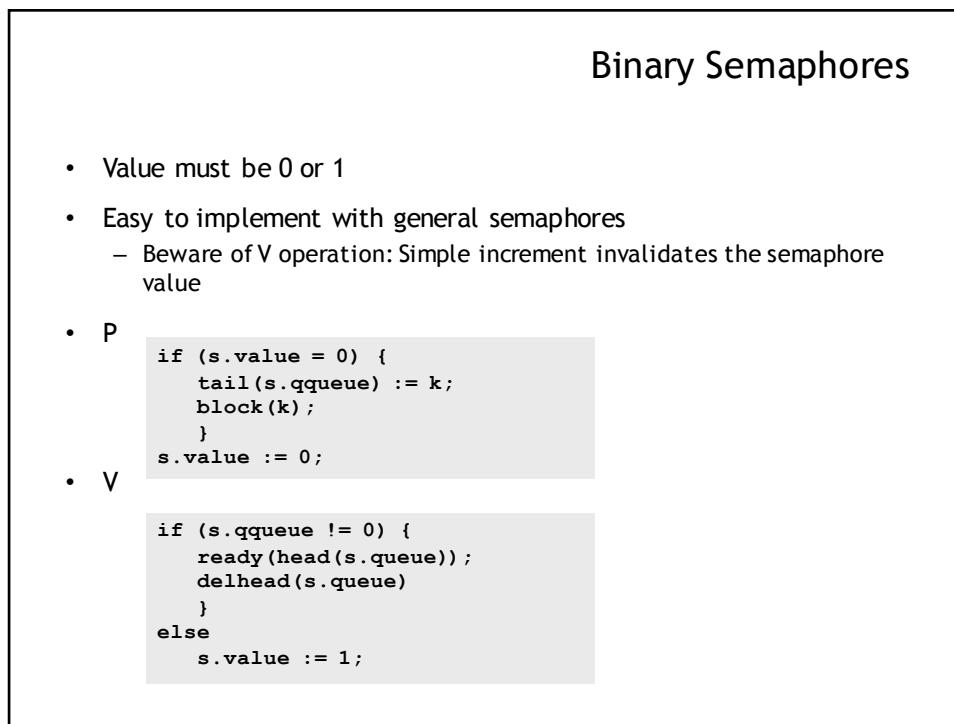
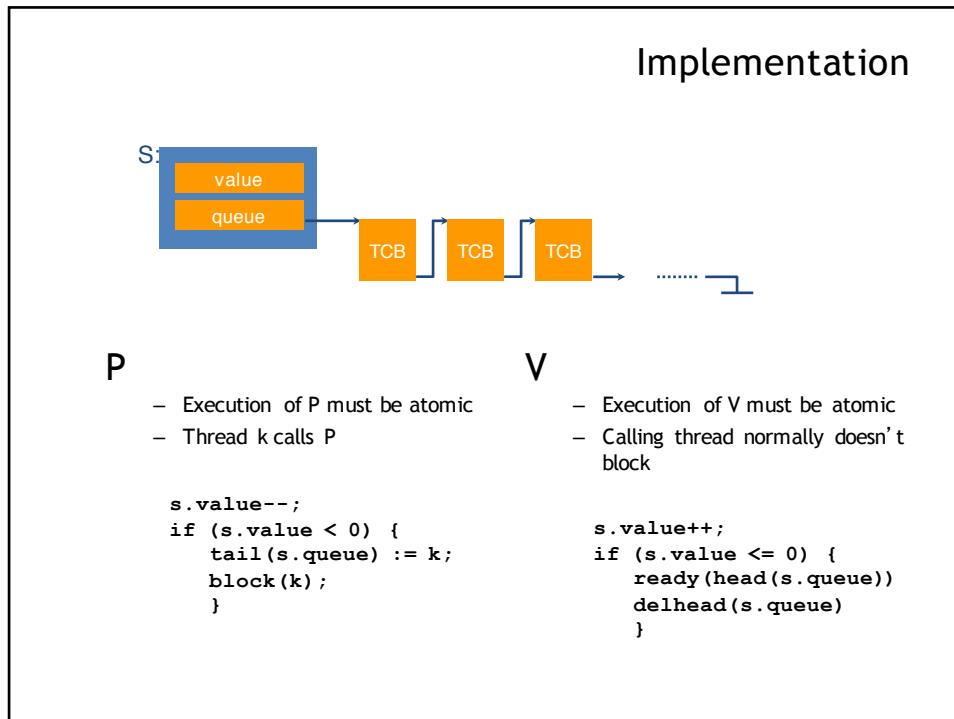
```
F:\OCP - Open Curriculum Project\OS Episodes\31 - IPC Reader Writer with Semaphores\R... □ x
Episode_31.Reader average waiting time 1.63256648492262
Episode_31.Reader average waiting time 1.221948490609363
Episode_31.Reader average waiting time 1.6325648492265
Episode_31.Reader average waiting time 1.9885896381579
Episode_31.Reader average waiting time 1.64492222721174
Episode_31.Reader average waiting time 1.25295528941582
Episode_31.Reader average waiting time 1.684905731357
Episode_31.Reader average waiting time 1.21814468370995
Episode_31.Reader average waiting time 1.6729485597717636
Episode_31.Reader average waiting time 1.17994254015565
Episode_31.Reader average waiting time 1.73775263771238
Episode_31.Reader average waiting time 1.19794254015565
Episode_31.Reader average waiting time 1.61228088751684
Episode_31.Reader average waiting time 1.21437142135399
Episode_31.Reader average waiting time 1.69873819348923
Episode_31.Reader average waiting time 1.221963649558612
Episode_31.Reader average waiting time 1.6729485597717636
Episode_31.Reader average waiting time 1.27062479227257
Episode_31.Reader average waiting time 1.618229319741822
Episode_31.Writer average waiting time 1.18408621118812
Episode_31.Writer average waiting time 0.8595989274451206
Press any key to continue
```

One Semaphore too much?

- Why not throw semaphore ReaderQueue away and move RMutex before PreferWriter? Wouldn't that serialize all readers too?

```
while (true) {
    ReaderQueue.P();
    PreferWriter.P();
    RMutex.P();
    if (readers_inside == 0)
        Sanctum.P();
    readers_inside++;
    RMutex.V();
    PreferWriter.V();
    ReaderQueue.V();
    // Read data
    RMutex.P();
    readers_inside--;
    if (readers_inside == 0)
        Sanctum.V();
    RMutex.V();
}
}                                Reader
```

IMPLEMENTATION DETAILS



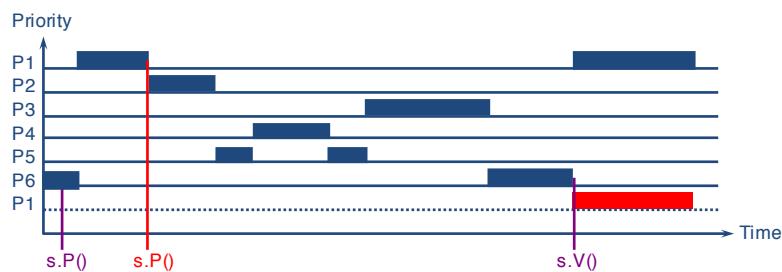
Nested (binary) Semaphores

```
Bsemaphor Mutex(1);

f1 (...) { ← f2 (...) {
    Mutex.P();      Mutex.P();
    ...
    Mutex.V();      f1();
    } }           .....           ...
                    Mutex.V();       ...
                                ...           Mutex.V();
                                } }
```

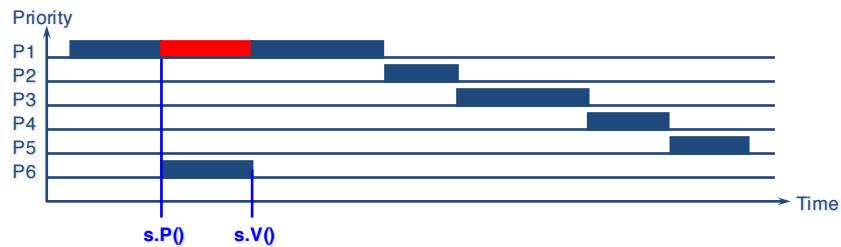
- Make multiple critical sections (procedures) thread-safe
 - Mutual exclusion among library functions
 - Reentrance
- Most semaphore implementations are smart or friendly

Semaphores and Scheduling Priorities



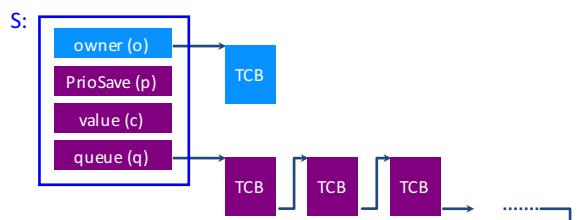
- P1 and P6 are synchronized using semaphore s
- Priority inversion
 - Apparent priority of P1 is less than priority of P6
- Dangerous for real-time systems

Priority Inheritance



- Low priority thread inherits higher priority until calling V
- Any modern OS supports inheritance
- Not necessarily dequeue head of blocked list (FCFS)

Implementing Priority Inheritance



- Limited to binary semaphores
- P operation
 - No blocking: Store ownership and current priority
 - Blocking: Transfer higher priority to owner thread if needed
- V operation
 - Restore original priority
 - Put blocked thread with highest priority in ready queue

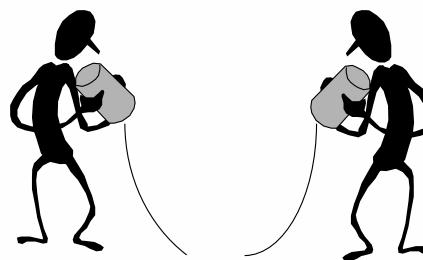
Operating Systems

Communication
Motivation

System Software Group
University of Trier
D-54286 Trier, Germany
Contact: Peter Sturm
sturm@uni-trier.de

Synchronization and Communication

- Interprocess Communication (IPC) typically comprises synchronization as well as communication issues
 - Synchronization
 - A waits until a specific condition holds
 - B fulfills this specific condition and “informs” A
 - Only a single bit is the information entity
 - Communication
 - True exchange of data
- Various communication concepts exist



Operating Systems

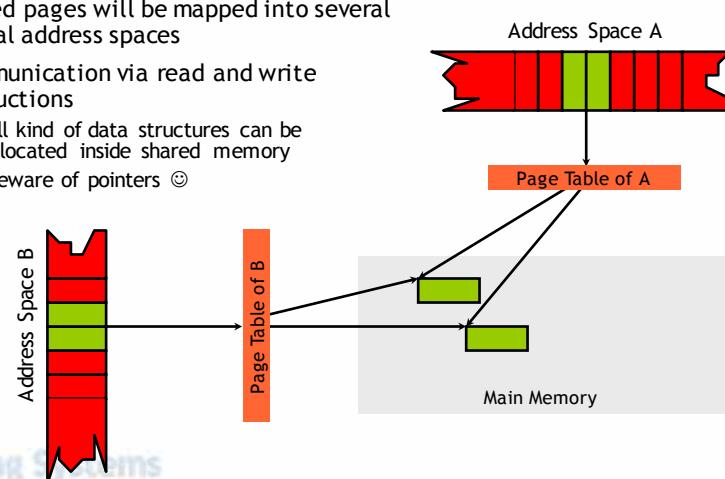
3 Classes of Communication

- Within a single process
 - All threads inside a process share the complete address space
 - Communication through any data structures provided by the programming language
- Between processes on a single machine
 - Explicit use of memory portions shared between different address spaces
 - Fast but complicated way to communicate
- Between processes on different machines
 - Distributed systems
 - Message passing is the primary form of communication
 - Some people prefer to use various kinds of distributed shared memory

Operating Systems

Sharing Memory Between Processes

- Multiple processes share a given number of pages explicitly
- Shared pages will be mapped into several virtual address spaces
- Communication via read and write instructions
 - All kinds of data structures can be allocated inside shared memory
 - Beware of pointers ☺



Operating Systems

Beware of Pointers!

The diagram shows two address spaces, A and B, represented as vertical rectangles divided into segments. Address Space A contains a yellow segment at the top labeled 'C' and a white segment below it labeled 'X'. Address Space B contains a yellow segment at the top labeled 'C' and a white segment below it labeled 'Y'. A green arrow points from 'X' in Address Space A to 'C' in Address Space B, indicating a pointer from one process to another. A red arrow points from 'Y' in Address Space B back to 'C' in Address Space A, indicating a pointer from one process to another.

- Dangerous to use pointers inside shared memory
- Pointers within shared memory itself are okay if mapped to the same address in all processes
 - Therefore, the address to be mapped to can be fixed
- Pointers stored inside but pointing outside of shared memory
 - are deadly since they point to something different in other address spaces

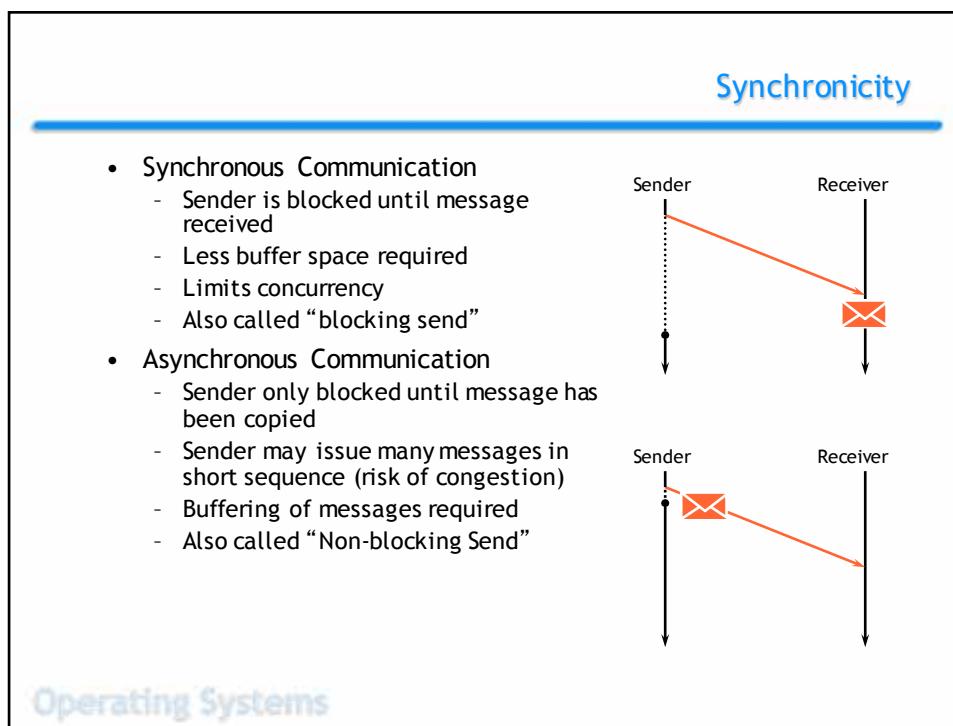
Message-based Communication

The diagram shows two computer systems. Each system has an application layer (green ovals) and a system layer (blue rectangle). The left system's application layer has a red envelope icon pointing to the system layer. The right system's application layer has a red envelope icon pointing away from the system layer. Below each system is a horizontal double-headed arrow labeled 'Rechner'.

- Message = Sequence of bytes
- 2 fundamental operations
 - send
 - receive
- For many message-based techniques, it is easy to cross computer boundaries

Operating Systems

| Comparison | |
|--|---|
| Memory-based Communication | Message-based Communication |
| <ul style="list-style-type: none"> Access shared memory via read and write operations | <ul style="list-style-type: none"> Send and receive messages |
| <ul style="list-style-type: none"> Pro <ul style="list-style-type: none"> Performance Transparency (for some) Contra <ul style="list-style-type: none"> Shared memory needed (requires monoprocessor or multiprocessor) Transparency (for others) You don't see when you communicate Need for synchronization to preserve data consistency | <ul style="list-style-type: none"> Pro <ul style="list-style-type: none"> More general concept (applicable to monoprocessors, multiprocessors, and even distributed systems) No transparency (for some) Contra <ul style="list-style-type: none"> Efficiency <ul style="list-style-type: none"> In case sender and receiver are on the same host Networks are slow compared to memory No transparency (for others) |
| Operating Systems | |



Communication Pattern

- Smallest unit of communication
 - Notification
 - Single message from sender to receiver
 - Service Invocation
 - Client request to server
 - Reply back to sender
- Additional patterns (mostly used in distributed systems)
 - Multicast
 - Send a message to a group of receivers
 - Broadcast
 - Send a message to any potential receiver on the network

Operating Systems

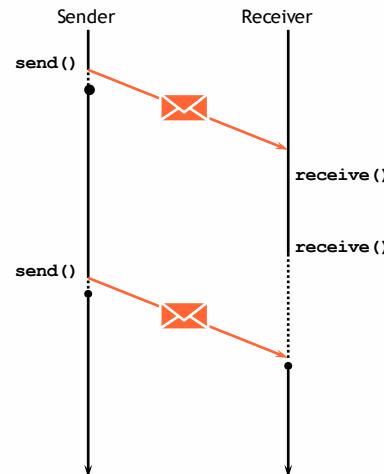
4 Basic Ways to Communicate

| | Asynchronous | Synchronous |
|--------------------|---------------------------------|--------------------------------|
| Notification | Datagram | Rendezvous |
| Service Invocation | Asynchronous Service Invocation | Synchronous Service Invocation |

Operating Systems

Datagram

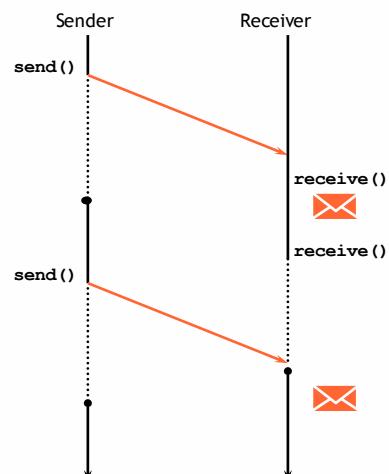
- Sender and receiver are decoupled
 - Concurrency
- Operating system has to buffer messages
 - Complex (overflow, ...)
 - Sometimes still need to block sender in case of resource limitations
- Examples
 - UDP (User Datagram Protocol)
 - Part of TCP/IP
 - Best Effort
 - max. 64 kbyte data
 - Signals
 - Software interrupts in UNIX
 - No data or a maximum of 4 Byte
- Sender knows ...
 - Nothing



Operating Systems

Rendezvous

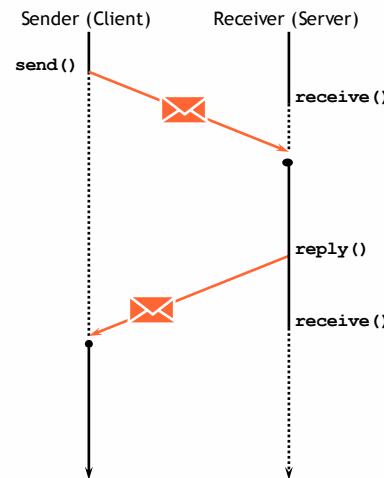
- Sender and receiver are coupled for the duration of message transfer
 - Limited concurrency
- No buffers required
 - Direct copy of data from sender to receiver
 - Very attractive when crossing address space boundaries
 - Simple to implement
- Examples
 - Ada, QNX, ...
- Sender knows ...
 - that message arrived safely



Operating Systems

Synchronous Service Invocation

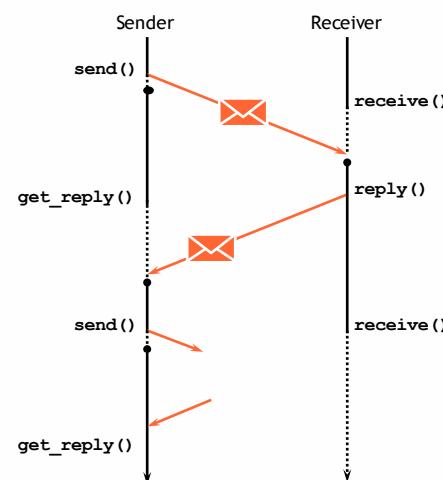
- Even less concurrency
- Bi-directional communication possible
- Example
 - Remote Procedure Call (RPC)
 - DCE, Corba, COM+, Java RMI, ...
 - QNX
- Sender knows ...
 - Message has been received
 - Request has been processed
 - The result



Operating Systems

Asynchronous Service Invocation

- More concurrency again
- Sometimes called “asynchronous RPC”
- Examples
 - V-Kernel
- Sender knows ...
 - Message has been received
 - Request has been processed
 - The result



Operating Systems

Direct and Indirect Communication

Direct

- The receiver is addressed, e.g. process id

Indirect

- The receiver is addressed via port or mailbox

Pro

- Efficiency

Contra

- Changing the receiver can be complicated

Pro

- Flexibility
 - Change of receiver
 - Multiple receivers

Contra

- Overhead induced by additional indirection

Operating Systems

Messages Across Address Spaces

Asynchronous

- 2 copies
 - From sender to buffer
 - From buffer to receiver
- Problems
 - Buffer management
 - Overflow handling

Synchronous

- Sender and receiver are blocked
 - Direct copy from sender to receiver possible
 - 1 copy
- Problems
 - Limited concurrency
 - Solved with multiple threads

Operating Systems

Beispiele

- Kommunikation zwischen zwei Adreßräumen
 - Signale
 - Asynchron, Mitteilung, Direkt, Paket
 - Pipes und Named Pipes
 - Synchron/Asynchron, Mitteilung, Direkt/Indirekt, Strom
 - Message Queues
 - Synchron/Asynchron, Mitteilung, Indirekt, Nachricht
- Client/Server-Modell
 - Remote Procedure Call (RPC)
 - Synchron/Asynchron, Auftrag, Direkt/Indirekt, Paket
 - RPC auf einem Rechner = Local Procedure Call (LPC)



Operating Systems

Signale

- Übermittlung asynchroner Ereignisse zwischen Prozessen
(keine zusätzliche Übertragung von Daten möglich)
- Ursprung
 - Software-Interrupts in Fehlerfällen und für "Job Control"
 - Arithmetikfehler
 - Adreß- und Busfehler
 - Unerlaubter Zugriff
 - Timer
 - ...
- Reaktion auf Signal
 - Blockieren: Zeitliche Empfangsbeschränkung (vgl. Disable Interrupts)
 - Ignorieren
 - Verarbeiten: Anmeldung eines Signal Handlers
 - Prozeß terminieren

Operating Systems

Beispiel

- Prozeß gibt Zustandsinformationen bei Erhalt eines Signals aus

```

int signal_arrived = 0;

void signal_handler ( int signo ) {
    signal_arrived = 1;
}

int main () {
    signal(SIGUSR1,signal_handler);
    while (1) {
        // Tue etwas sinnvolles ...
        if (signal_arrived) {
            // Reagiere auf Signal ...
            signal_arrived = 0;
        }
    }
}

```

- Absenden des Signals, z.B.

```
kill -USR1 pid
```

Operating Systems

Pipes

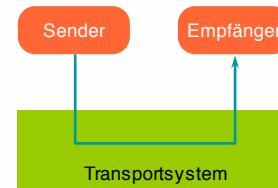
- “Pipes” und “Named Pipes” (auch FIFO genannt) für Prozeßkommunikation auf einem Rechner
- Pipes zwischen Vater- und Kindprozessen:
 - Verbindung über vererbte Datei-Deskriptoren
- Beispiel

```

main () {
    int pipe_ends[2];

    pipe(pipe_ends);
    if (fork() != 0) {
        // Vaterprozeß
        write(pipe_ends[1],Daten,...);
    }
    else {
        // Kindprozeß
        read(pipe_ends[0],Daten,...);
    }
}

```



Operating Systems

Exkurs: UNIX-Shell und Pipes

- Exemplarische Realisierung von: “a | b”

```

...
int pfd[2];
pipe(pfd);
if (fork()==0) {
    // Kind - Wird Prozeß a
    dup(pfd[1],1); /* stdout auf Pipe-Eingang */
    execve("a",...);
}

if (fork()==0) {
    // Kind - Wird Prozeß b
    dup(pfd[0],0), /* stdin auf Pipe/Ausgang */
    execve("b",...);

    if (Vordergrund-Bearbeitung) {
        wait(...);
    }
}

... zurück zur Eingabeschleife

```

Operating Systems

Realisierung



- Klassisches Erzeuger/Verbraucher-System
 - Puffer bestimmter Länge für jede Pipe
 - Empfänger wird blockiert, wenn keine Daten vorliegen
 - Sender wird blockiert, solange Puffer voll ist
- Synchrones Verfahren bei zu kleinem Puffer
 - Gefahr der Verklemmung?
- “Named Pipes”
 - Pipe-Deskriptoren nur über fork vererbbar
 - Pipe-Namen im Dateisystem aufbauen, damit unabhängige Prozesse miteinander kommunizieren können
 - Kommunizierte Daten werden nicht über Dateisystem geleitet!

Operating Systems

Client/Server Systems

- Based on synchronous service invocation
- Forms the platform for most distributed systems today
 - Worldwide availability of servers
 - Clients send requests
 - RPC dominant communication technique
 - Addresses open and heterogeneous systems
- But also a possible architecture for modern operating systems
 - Mikrokernel with minimal functionality
 - Lifting of operating system functionality into user-mode servers
 - Basis for distributed operating systems
 - Efficient communication among address spaces required

```

    graph LR
      Client[Client] -- "send()" --> Server[Server]
      Server -- "receive()" --> Client
      Client -- "reply()" --> Server
      Server -- "receive()" --> Client
  
```

The diagram illustrates the synchronous communication between a Client and a Server. It shows a vertical dotted line for the Client and a vertical solid line for the Server. A red arrow labeled "send()" points from the Client to the Server. A red arrow labeled "receive()" points from the Server back to the Client. A second set of red arrows labeled "reply()" and "receive()" indicates the return path from the Client to the Server.

Operating Systems

Client/Server Architecture

The diagram shows two separate computer systems, "Rechner 1" and "Rechner 2". Each system has a "Betriebssystem" (Operating System) layer and a "Rechner" (Computer) layer. On each system, there are multiple green "Client" boxes and one red "Server" box. Blue lines connect the clients on one system to the servers on the other, representing network communication.

- Lokal servers
 - File system
 - Input and Output
 - Graphics subsystem
 - Process management
 - Paging
 - ...
- Remote servers
 - Network file systems (e.g. shares or NFS)
 - Distributed file systems
 - Rlogin, Telnet, ...
 - WWW server (IIS, apache, ...)
 - Clock synchronization (ntp)

Operating Systems

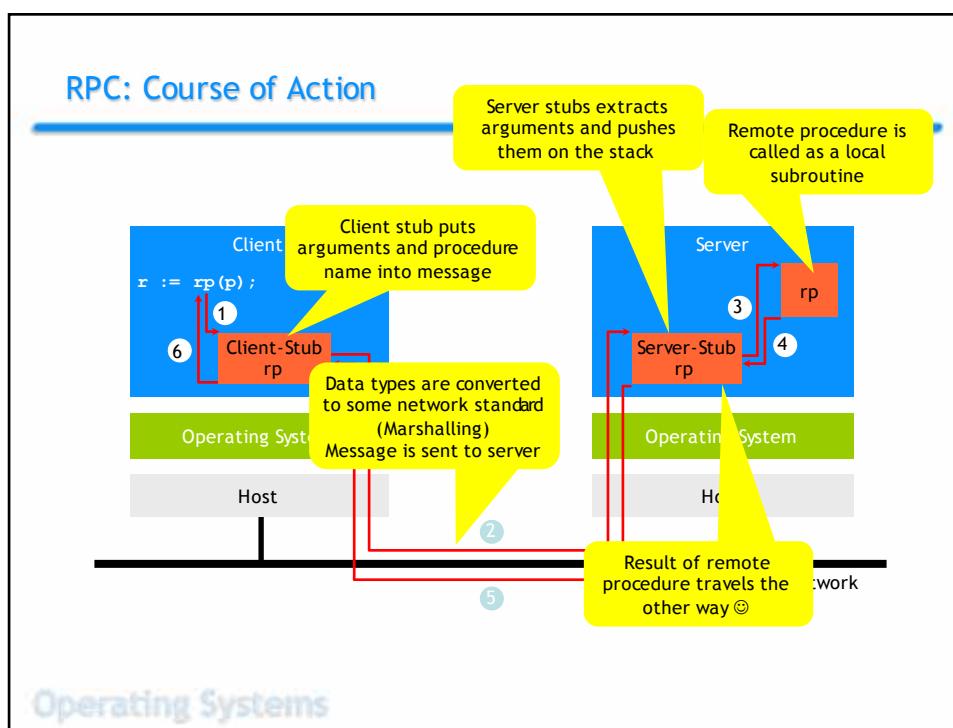
Remote Procedure Call

- Tries to resemble semantics of local procedure call
 - Push arguments on stack
 - Jump subroutine
 - Push return address on stack
 - Execute procedure
 - Store result in well-known register
 - Return from procedure
- Local procedure are part of address space of caller
- Now, procedure may be remote
 - In some other address space
 - Maybe on another machine

Local Procedure Call

Application
 $x := f(p);$
↓
f

Operating Systems



Marshalling

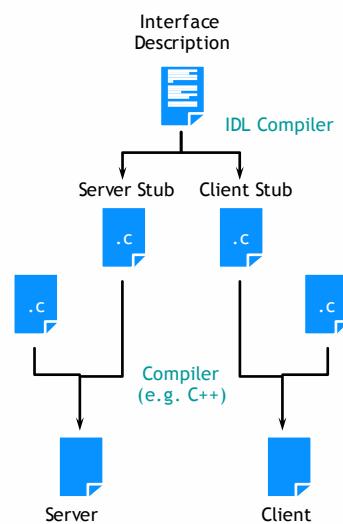
- Imported for distributed systems
 - Encoding of basic data types may differ
 - Little endian vs. big endian
 - ...
- Define a generic network standard
 - Sender converts data from local representation to network standard
 - Receiver converts data from network standard to local representation
- Not needed if we just want to cross address spaces not hosts



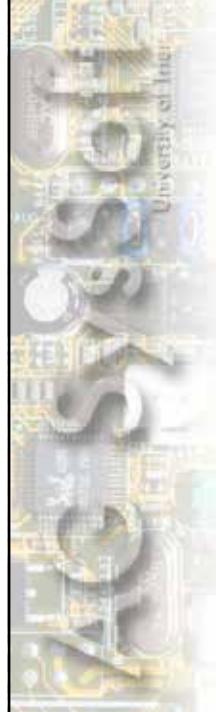
Operating Systems

RPC Compiler

- Server
 - Defines a set of remote procedures
- Service or interface description
 - Define signature for each method
 - Type and sequence of arguments
 - Name of procedure
 - Return type
- Interface definition language (IDL)
- IDL compiler automates most functionality to call a remote procedure
 - Client stubs
 - Server stubs
 - Additional server functionality



Operating Systems

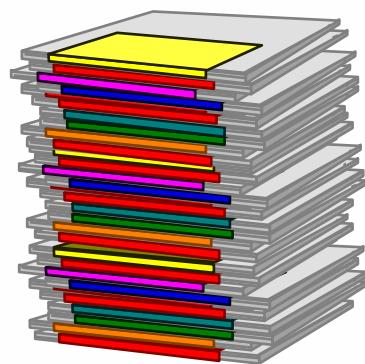


Filesystems

Peter Sturm
sturm@uni-trier.de
+49 651 201 3311

File Systems

- Most functionality of OS invisible for the user
 - Virtual address spaces
 - Threads
 - Synchronization
 - Communication
 - ...
- File system is the most visible OS component
 - Stores ALL persistent data
 - Data loss can be bothersome or disastrous



Peter Sturm, University of Trier

Basics

- Persistent storage of programs and data
- Basic abstraction = File
 - Sequence of data with a name
- File attributes (meta data)
 - Name
 - Type (in some file systems)
 - Physical position
 - Size
 - Access permissions
 - Time stamps
 - owner
- Directories
 - Hierarchical organization of files

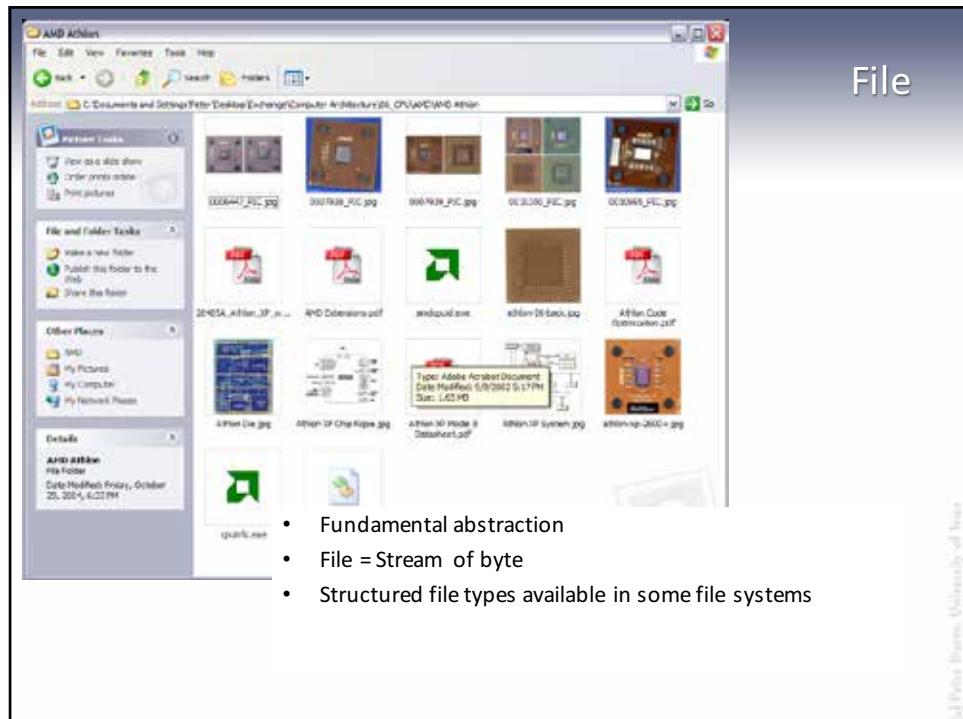


Alfred Sturm, University of Trier

Typical Usage of File

- Characteristics (identified empirically)
 - Most files are small (some KB) or very large (audio, video)
 - Reads are frequent, writes are less frequent
 - (Files are deleted rarely or never)
 - Sequential file access is dominant
 - Some files are shared by many users
 - Most files are not shared at all
- Characteristics determine functionality and structure of an efficient file system
 - Make the common case fast
 - Sequential read access for (small) files
 - Utilize caching
 - Frequent reading, 1 user, no sharing

Alfred Sturm, University of Trier



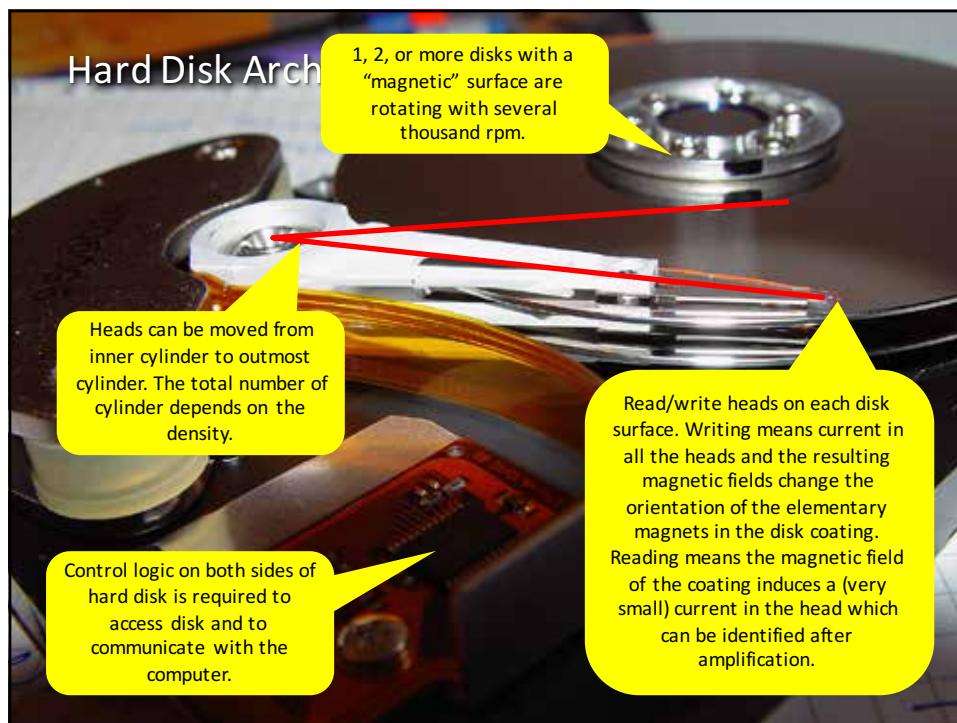
- Fundamental abstraction
- File = Stream of byte
- Structured file types available in some file systems

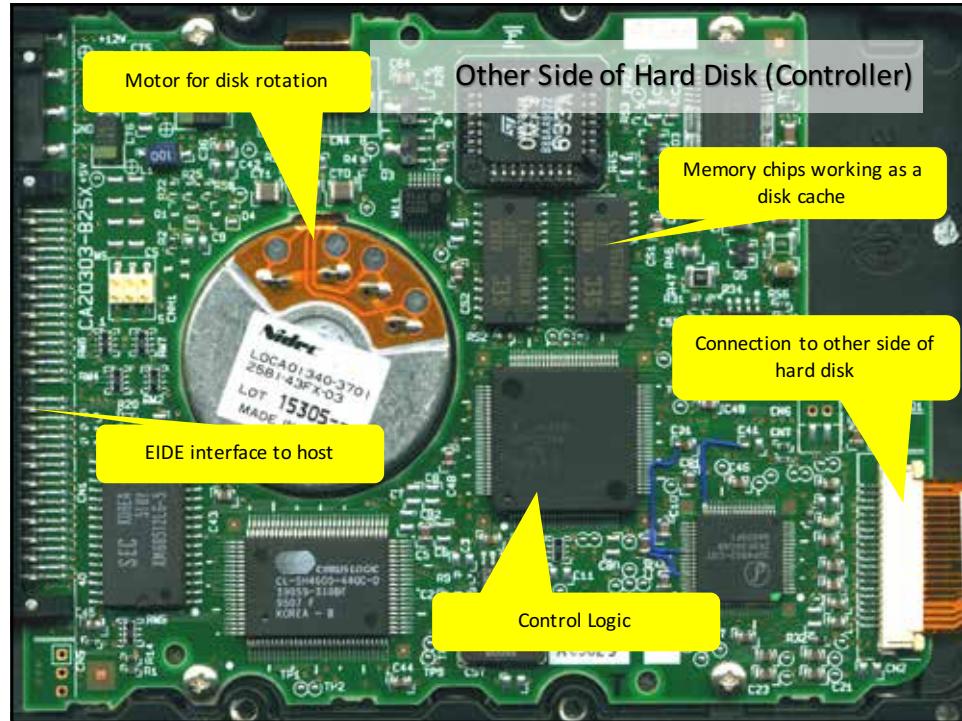


Directories

- Contains
 - Files
 - Directories
- Hierarchical structure

Peter Sturm, University of Trier





| Comparison: Main Memory and Disk Memory (of July 2005) | |
|--|--|
| Internal Memory <ul style="list-style-type: none"> – Volatile – Byte and word based access – Quite fast ($\leq 10\text{ns}$) – Several GB/s throughput • Expensive <ul style="list-style-type: none"> – 1 GB = 110 Euro – 1 MB = 11 Cent • Common capacities <ul style="list-style-type: none"> – 512 MB .. 4 GB on PCs | External Memory <ul style="list-style-type: none"> – Persistent – Block based access (512 Byte ... 8 KByte) – Very slow (4-8 ms) – 50 MB/s and more • Cheap <ul style="list-style-type: none"> – 200 GB = 100 Dollar – 1 MB ca. 0.05 Cent • Common capacities <ul style="list-style-type: none"> – 160 GB .. 400 GB and more |
| <p>These values change very fast, so they may differ substantially from your findings.</p> | |

Comparison: Main Memory and Disk Memory (of January 2009)

Internal Memory

- Volatile
- Byte and word based access
- Quite fast ($\leq 10\text{ns}$)
- Many GB/s throughput
- **Expensive**
 - 2 GB = 30 Euro
 - 1 MB = 1.46 Cent
- **Common capacities**
 - 1 GB .. 16 GB on PCs

External Memory

- Persistent
- Block based access
(512 Byte ... 8 KByte)
- Very slow (4-8 ms)
- 100 MB/s and more
- **Cheap**
 - 1000 GB = 90 Euro
 - 1 MB ca. 0.009 Cent
- **Common capacities**
 - 200 GB .. 1500 GB and more

Al Peter Sturm, University of Trier

Comparison: Main Memory and Disk Memory (of January 2012)

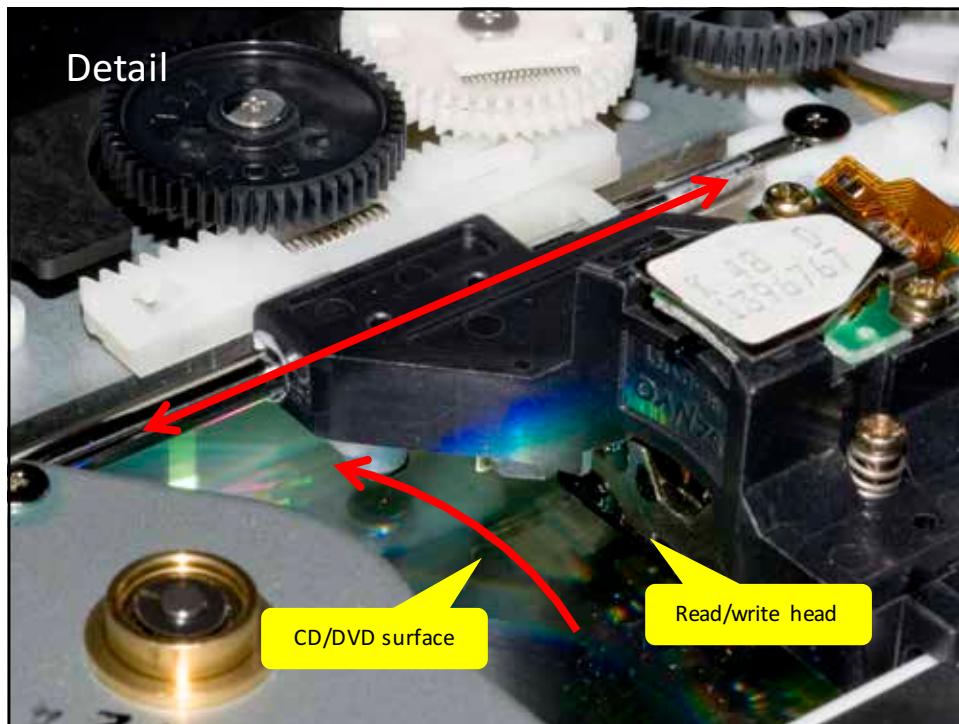
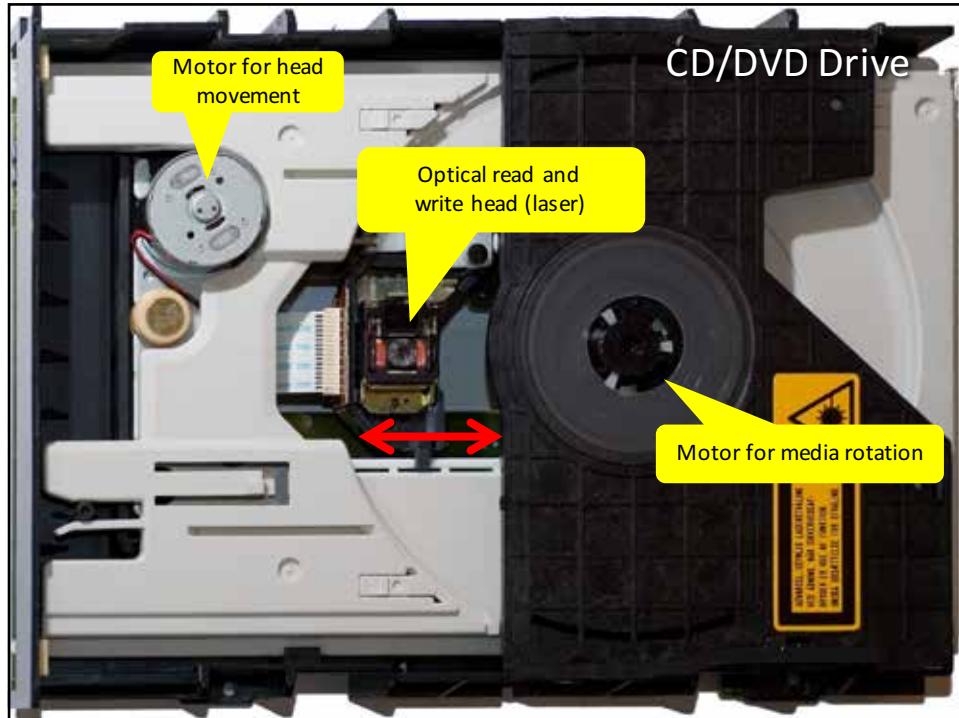
Internal Memory

- Volatile
- Byte and word based access
- Quite fast ($\leq 10\text{ns}$)
- Many GB/s throughput
- **Expensive**
 - 4 GB = 20 Euro
 - 1 MB = 0.48 Cent
- **Common capacities**
 - 4 GB .. 16 GB on PCs

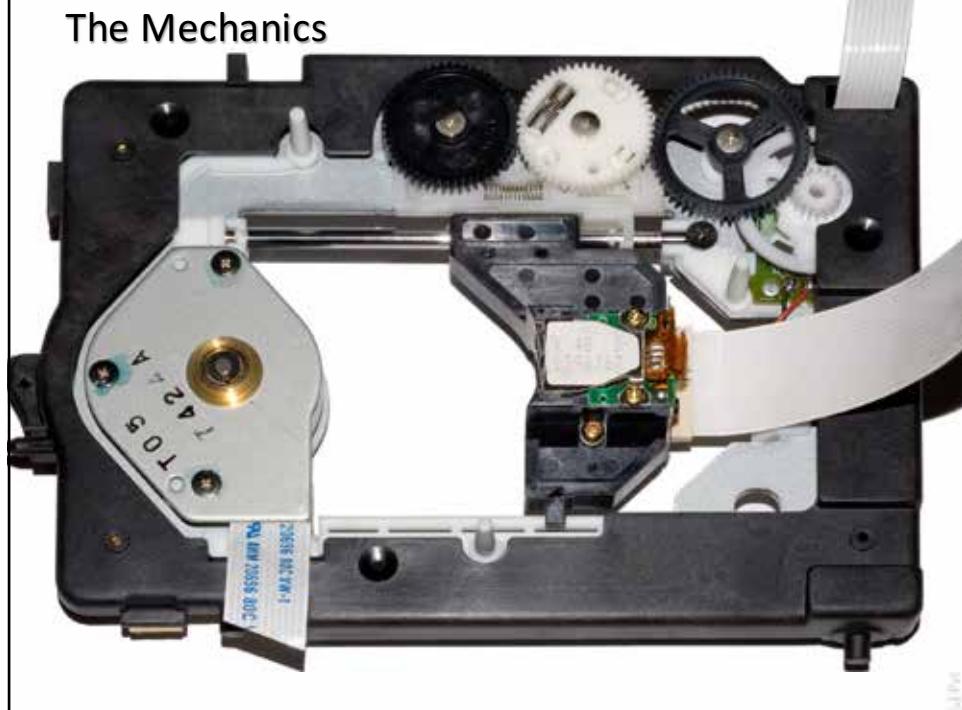
External Memory

- Persistent
- Block based access
(512 Byte ... 8 KByte)
- Very slow (4-8 ms)
- 100 MB/s and more
- **Cheap**
 - 2000 GB = 120 Euro
 - 1 MB ca. 0.005 Cent
- **Common capacities**
 - 500 GB .. 3000 GB and more

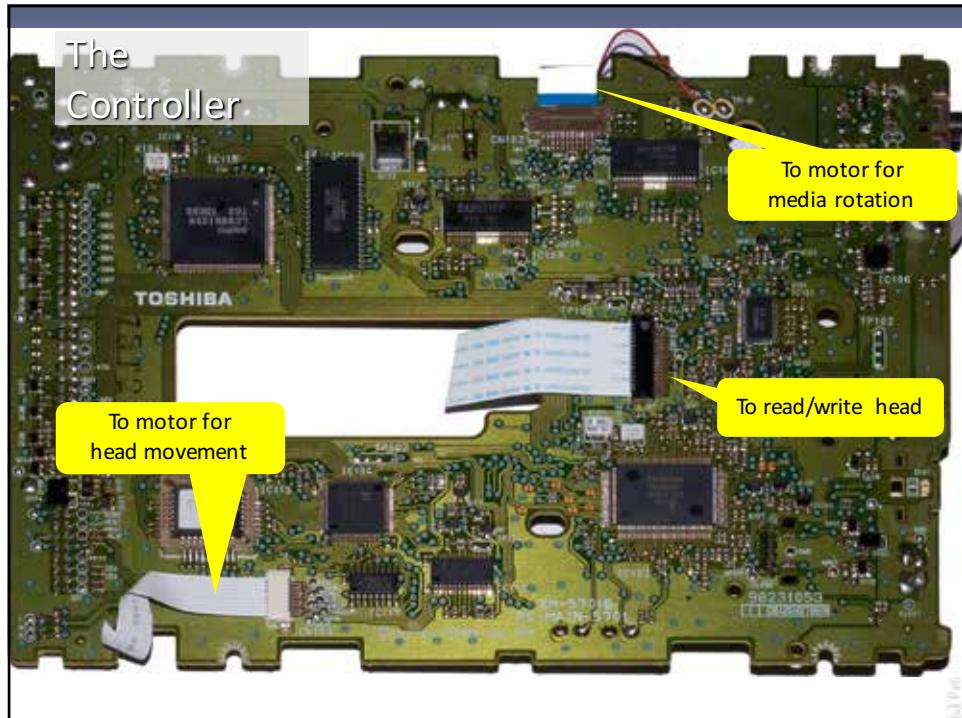
Al Peter Sturm, University of Trier

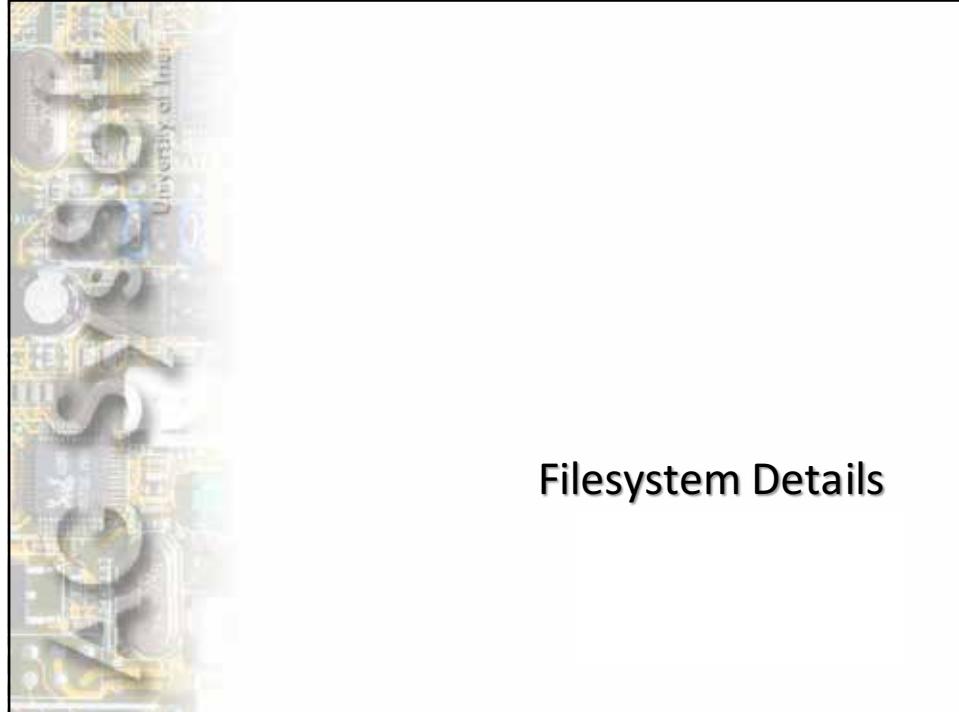


The Mechanics



The Controller





Introduction

- Persistent storage of applications, user files, configuration data, operating system
- Basic abstraction
 - File = Sequence of bytes
- Characteristics of hard drives require block-based access
 - Head movement limiting factor
 - Typical size is 4 kbytes
- Make the common case fast
 - Reading files sequentially
 - Read ahead, Caching

Requirements for Modern Filesystems

- Minimize time-consuming head movement
 - Store file data and all accompanying meta-data in close neighborhood
 - Collect multiple disk requests and optimize movement
 - Cluster blocks consecutive blocks of a file physically
- Improve read/write performance
 - Utilize caching aggressively
(within hard disk itself and operating system)
 - Exploit striping and mirroring to improve performance and/or fault tolerance
- Keep pace with increasing disk capacity
 - “Since disk usage is like an ideal gas, which takes up every available space”, file systems must deal with more and more files
 - Recovery in case of system failure must be fast

Functions

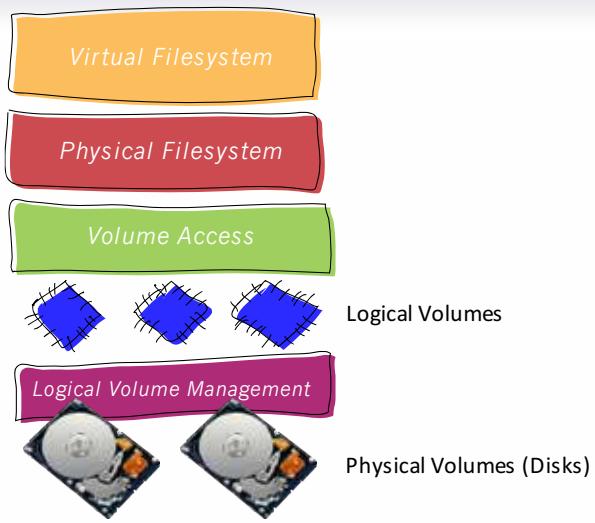
- Formatting and initializing media
 - Create initial content to access media
- Keeping track of bad blocks
 - Recognize bad blocks and exclude them from future use
- Managing of meta-data
 - Small set of data but high update frequency
- Storing and retrieval of files
 - Sequential reading and writing
 - Positioning of file pointer
 - Mapping files into virtual address space
- Security
 - Who's allowed to access a file?

Meta Data

- File information (inode)
 - Size, owner, access rights, timestamps
 - Where are the blocks on disk
- Directory information
 - Content, owner, access rights, timestamps
- Device information
 - root directory
 - Free blocks
 - Bad blocks

Al Peter Sturm, University of Trier

Hierarchical Structure



Al Peter Sturm, University of Trier

Classes of filesystems

- Media characteristics
 - Read/Write
 - Hard disk filesystems (most common)
 - FATx, NTFS, HFS, ext2, ext3, reiserfs, XFS, ZFS, ...
 - Limited write
 - Filesystems for flash media
 - Readonly
 - CD/DVD
 - ISO9660, UDF, ...
- Networking capabilities

Al Peter Sturm, University of Trier

In Place vs. Log

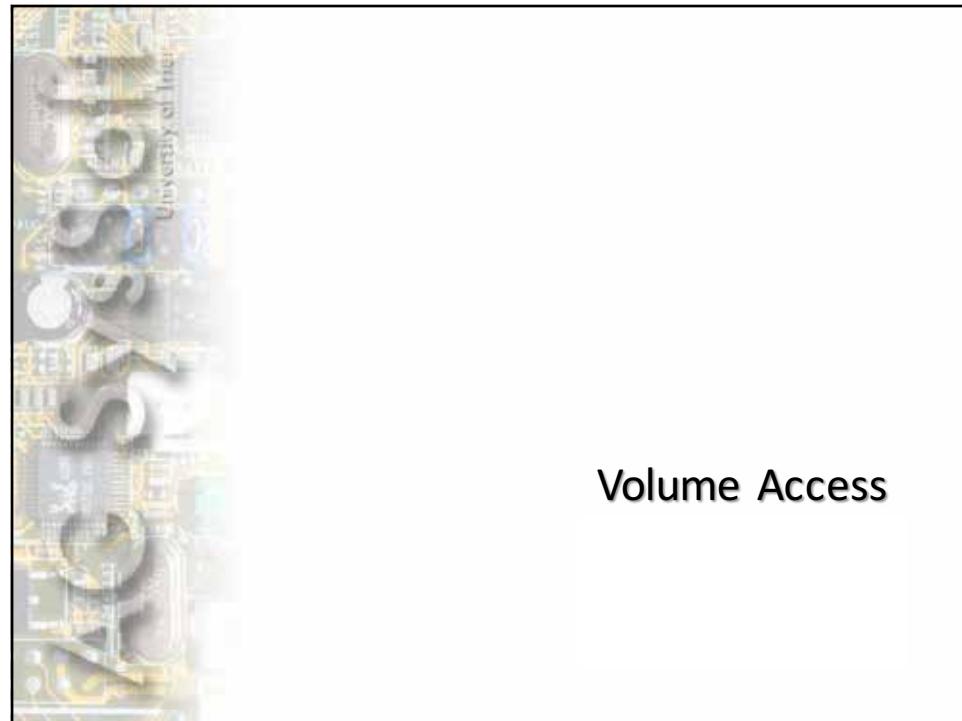
- Sometimes high update frequency
 - e.g. every access changes some timestamps
- Small updates
 - 4 bytes timestamp
- Cache changes might lead to inconsistencies
- How to update
 - In place
 - Log-based



Al Peter Sturm, University of Trier

| | | Classification | |
|-----------|-----------------|---|---|
| | | Data | |
| | | Update in Place | Log |
| Meta Data | Update in Place | Traditional filesystems ext2, FAT, ... | unknown |
| | Log | Journaling filesystems ext4, NTFS, HFS, reiserfs, XFS, ... | Log-based filesystems ZFS, btrfs, ... |

Algoritmen, Datenstrukturen und Systeme



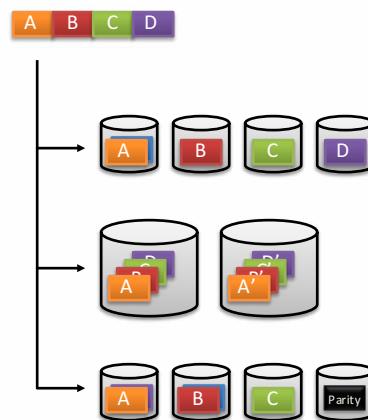
Basic access

- Defines unique address for each block
 - Outdated is Cylinder/Head/Sector (CHS)
 - Logical block address (LBA)
- Size of address determines maximum size of filesystem
 - Huge numbers for today's disks
- Basic functionality
 - Access block by address
- Maintaining bad block and free block lists
 - Typically stored in bitlists
- Implemented by BIOS and hardware

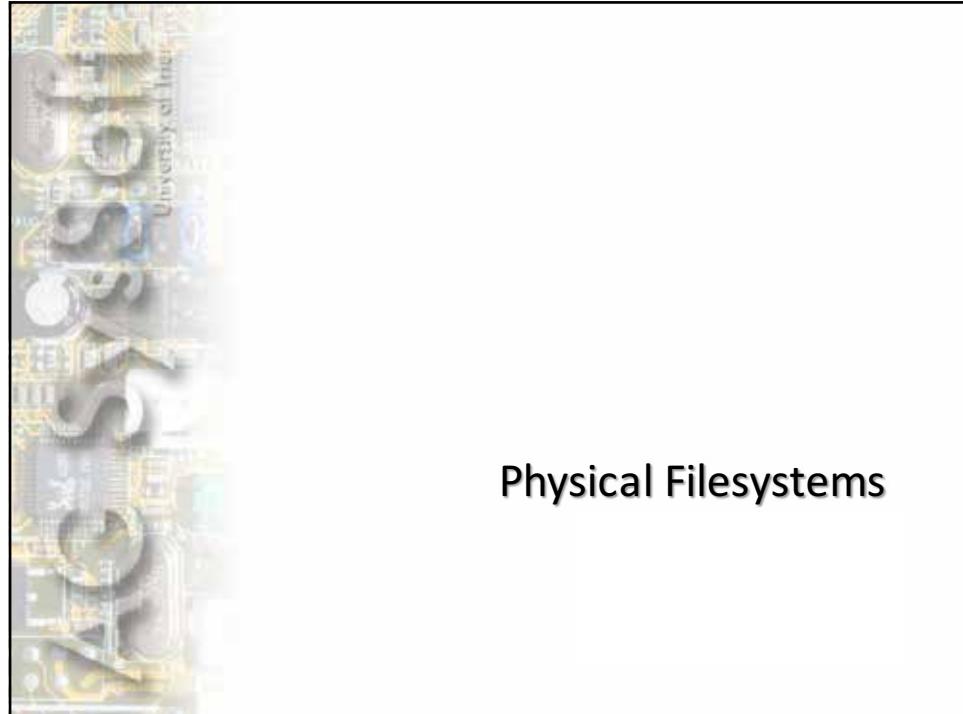
Al Peter Sturm, University of Trier

Redundant Arrays of Independent Disks (RAID)

- 6 Levels
 - RAID 0: Striped
 - RAID 1: Mirrored
 - RAID 2: Striped+ ECC
 - RAID 3: Striped+ Parity
 - RAID 4: Huge stripes + Parity
 - RAID 5: Huge stripes + Rotating Parity
 - RAID 6: RAID 5 + additional parity
- Most common are RAID 0,1,5 and combinations (RAID 10 = 1 + 0)
- Favor huge writes
 - Overhead $1/(n-1)$ for n disks
- Small writes are expensive
 - Read data
 - Write data
 - Read parity
 - Write parity



Al Peter Sturm, University of Trier



Physical Filesystems

The slide features a title "File Organization" at the top center. Below the title are two diagrams illustrating different file allocation methods. The left diagram shows contiguous allocation with a sequence of blocks numbered 0, 1, 2, followed by four empty blue blocks. The right diagram shows distributed allocation with a sequence of blocks numbered 2, 1, 0, followed by three empty blue blocks. Both diagrams consist of alternating orange and blue rectangular boxes.

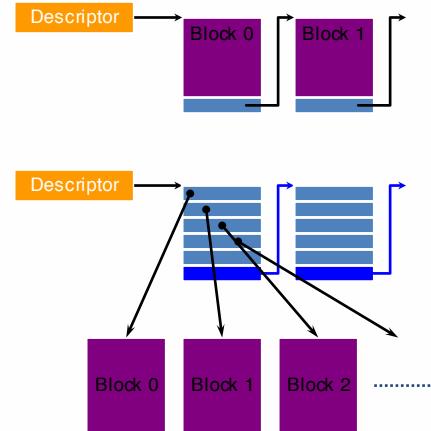
File Organization

| Contiguous Allocation (Clustering) | Distributed Allocation |
|------------------------------------|------------------------|
| | |

- Contiguous allocation (clustering)
- Advantages
 - Fast sequential and random access
 - Easy to manage
- Disadvantages
 - External Fragmentation
 - Shrinking of files creates holes
 - File growth can be expensive
 - Defragmentation (space)
- Distributed allocation
- Advantages
 - Internal fragmentation only
 - Size change of file simple
- Disadvantages
 - No uniform access time
 - Complex to manage
- Hybrid approaches are favored
 - Clustering as far as possible
 - But distributed allocation always possible
 - Defragmentation (performance)

Distributed Allocation

- Dedicated index structure to keep track of block order
- Internal
 - Index part of data blocks
 - Exotic block sizes (no multiple of 2)
 - Random access expensive
- External
 - Dedicated index blocks
 - Single-linked and double-linked lists
 - Trees, balanced trees
 - Hybrid approaches



File Allocation Table (FAT)

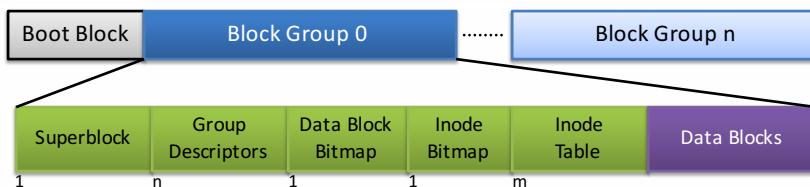
- Developed by Microsoft for MS-DOS (1977)
- Best supported filesystem
- Disk space is allocated to files in contiguous groups (clusters)
 - FAT (FAT12) = 2^{12} clusters \Rightarrow max. 32 MBytes disk size
 - FAT16 \Rightarrow 2 GByte
 - FAT32 \Rightarrow 8 TByte
- Linked list of clusters for bigger files
- Newly created files tend to be fragmented
 - Defragmentation needed

ext2

- The original Linux filesystem
 - Origin was Minix filesystem (just as in case of Linux itself)
 - Extended filesystem (inefficient extension)
 - 2nd Extended filesystem (since 1994)
- Properties
 - Selectable block size between 1024 and 4096 bytes
 - Selectable number of inodes per partition
 - Clustering of blocks
 - Preallocation of file blocks (enlargement)

Alfred Sturm, University of Trier

Structure



- Reason for block groups?
- Maximum size of block group:
 - Data Block Bitmap = 1 Block
 - 8 MB (1 KB blocks) up to 128 MB (4 KB blocks)
- Superblock and group descriptors are replicated in each group
 - Kernel will access block group 0 only

Alfred Sturm, University of Trier

Superblock

- Total number of inodes
- Block size
- Size of overall filesystem in blocks
- Number of reserved blocks
 - Only accessible by root processes
 - Accessing an otherwise full filesystem
- Number of free blocks and free inodes
- Characteristics of block groups
- Statistics
 - Timestamps (Last mount operation, Last write, Last filesystem check)
 - Mount statistics (Number of mount operations, Maximum mounts before check)
- Status flag
 - 0 mounted or missed unmount, 1 cleanly unmounted, 2 Filesystem has errors



Alfred Sturm, University of Trier

Trends

- Reinvent the file system
 - Apply the state of the art
 - Copy-on-Write, Transactional, Checksums, Ease of use
- Precursor
 - ZFS (Sun Microsystems)
 - Available in Solaris, FreeBSD, (Linux), (Mac OS)
 - Pooling, Software RAID
- Imitators ☺
 - btrfs (GPL, Oracle)
 - ReFS (Microsoft, Windows 8 Server)

Alfred Sturm, University of Trier