

# STUDENT'S HANDBOOK GUIDE TO CMPUT 379

## OPERATING SYSTEM CONCEPTS



A STUDENT'S GUIDE TO THE COMPLEX NATURE OF  
OPERATING SYSTEMS

Learn about the definition of a process; process states and state transitions; process control block; operations on processes; interrupt processing; parallel processing; resource allocation; shared and unshared allocation; critical sections; semaphores; deadlock; deadlock prevention, avoidance, detection, and recovery; memory management; memory allocation schemes; virtual memory; paging and segmentation; page replacement strategies; working sets; demand paging; job and processor scheduling; scheduling levels, objectives, and criteria; various scheduling algorithms; multi-processor considerations; file system functions; file organization; tree structured file systems; space allocation; file catalogs; file access control mechanisms; operating systems security.

By mastering the complex sub-field of Operating Systems will lead you to a heroic victory against the machination of the professor. By slaying his programming minions, defeating his midterm champion and ending the final themselves. The heroic legacy of your achievement will be remembered on your academic transcript.

---

Disclaimer: With the sole exception of the cover this book will crush your hopes and dreams as it does not contain prehistoric or high fantasy content like elves, dwarfs, dragons, or aarakocras.

# CONTENTS

<b>1: Introduction and Overview</b>	<b>1</b>	CPU State . . . . .	19
Abstraction . . . . .	1	Resource Management . . . . .	20
Example . . . . .	1	Computer Organization . . . . .	20
Compiling, Linking, Loading . . . . .	1	Activities . . . . .	20
Compiling . . . . .	1	OS Operating Cycle . . . . .	21
Linking . . . . .	1	Process Context . . . . .	21
Loading . . . . .	2	Suspended State . . . . .	21
Intro to x86 Memory Layout . . . . .	2	Threads vs Processes . . . . .	22
Basic Storage classes in C . . . . .	2		
Memory layout in the User Space . . . . .	3		
<b>2: ELF: Executable and Linking Format</b>	<b>4</b>	<b>5: Synchronization</b>	<b>24</b>
Introduction . . . . .	4	Example of Thread Concurrency . . . . .	24
File Format . . . . .	4	Mutual exclusion . . . . .	24
EFL Header . . . . .	5	Example of mutual exclusion attempts .	25
Program Header Table . . . . .	5	A Better Way . . . . .	27
Section Header Table . . . . .	5	Taxonomy of Interrupts . . . . .	28
Example of an ELF . . . . .	5	Intel x86 interrupts . . . . .	28
<b>3: DLL</b>	<b>11</b>	Another view of mutual exclusion . . .	29
How are DLLs loaded . . . . .	13	Interrupt masking . . . . .	30
<b>4: Processes</b>	<b>16</b>	<b>6: Synchronization Abstraction</b>	<b>31</b>
Conceptual Framework . . . . .	16	Semaphores . . . . .	31
Process Control Block . . . . .	17	Counting Semaphores . . . . .	31
Process Queues . . . . .	17	Binary semaphores . . . . .	33
Components of a Procedure . . . . .	17	Producer-Consumer semaphore . . . .	34
Context Switching . . . . .	17	Readers and Writer Problem . . . . .	35
		Limitation of classic semaphores . . .	35
		Monitors . . . . .	36
		Monitor Semantics . . . . .	36

Monitor Implementation of a Binary Semaphore . . . . .	36	Round Robin . . . . .	50
Semaphore implementation of monitors (Hoare) . . . . .	37	FPPS: Fixed Priority Preemptive Scheduling . . . . .	50
Java Example . . . . .	38	Hard deadlines . . . . .	51
Limitation of Monitors . . . . .	39	Multi-mode systems/Multiple policies . . . . .	52
Atomic Transaction . . . . .	39	Multi-level queues . . . . .	52
Serializability . . . . .	39	Feedback queues . . . . .	53
Conflicting operations . . . . .	40	Real-time vs Standard System . . . . .	53
Attaining concurrency and serializability . . . . .	40	Primary sources of problems . . . . .	53
General Problem of Synchronization . . . . .	40	Examples of Non Real-time paradigms . . . . .	54
<b>7: Deadlocks</b>	<b>41</b>	Sync Issues . . . . .	54
Resource Allocation . . . . .	41	Priority Inversion . . . . .	54
How to avoid deadlock? . . . . .	41		
Resource Types . . . . .	41		
Necessary conditions for a deadlock . . . . .	42		
Ordering Resources . . . . .	42		
Banker's Algorithm . . . . .	43		
Safe vs Unsafe State . . . . .	43		
General Case and Algorithm . . . . .	43		
Deadlock Detection . . . . .	46		
Is deadlock detection important? . . . . .	46		
Another useful operation . . . . .	46		
<b>8: Scheduling</b>	<b>47</b>		
Short Term Scheduling . . . . .	47		
Guessing the CPU burst . . . . .	47		
SJF: Shortest Job First . . . . .	48		
Non-Preemptive vs Preemptive . . . . .	48		
Priority . . . . .	49		
Starvation . . . . .	49		
<b>9: Memory Management</b>	<b>56</b>		
Memory Allocation (management) . . . . .	56		
Overlays . . . . .	56		
Memory layouts of early programs . . . . .	57		
Why separate data and code? . . . . .	58		
Trampolines . . . . .	58		
Signal Handlers . . . . .	59		
The Goal of Memory Allocation schemes	60		
Fixed Partition . . . . .	61		
Fragmentation . . . . .	62		
Variable partition . . . . .	62		
Gap filling strategies . . . . .	63		
Compaction . . . . .	63		
Memory Sharing . . . . .	64		
Swapping . . . . .	64		
When to swap out? . . . . .	64		
Paging . . . . .	65		
Caching . . . . .	66		
TLB . . . . .	66		

Address Translation . . . . .	66	Optimum page fault frequency . . . . .	86
Multi-level Tables . . . . .	66	So how much memory? . . . . .	86
A Few notes . . . . .	68	Example . . . . .	87
Caching . . . . .	68	A simple solution . . . . .	88
Locking . . . . .	69	The triple point goal . . . . .	88
Associativity types . . . . .	69	How it all looks like in real life . . . . .	89
IBM 370 TLB . . . . .	70	Example . . . . .	89
Principle . . . . .	72	Victimizer . . . . .	90
Common Segment . . . . .	73	Dividing memory among programs . . . . .	91
Dynamics of address spaces . . . . .	73	Keeping track of pages (page directory) . . . . .	91
A simple trick . . . . .	74	Modern trends . . . . .	92
Copy on Write . . . . .	75	Going beyond 32-bits . . . . .	93
More tricks . . . . .	75	Inverted page tables . . . . .	93
Prerequisites for CRAPE . . . . .	76	What about having no tables at all . . . . .	93
<b>10: Virtual Memory</b>	<b>78</b>	Hardware assist tools . . . . .	93
Swapping out . . . . .	78	paradigm shift . . . . .	94
Swapping at the page-level . . . . .	79	Power PC . . . . .	94
How much memory does does a process need . . . . .	79	MIPS . . . . .	95
Page replacement policies . . . . .	79	IA-64 . . . . .	95
Some ideas . . . . .	79	Superpages . . . . .	96
FIFO . . . . .	80	<b>11: File Systems</b>	<b>97</b>
LRU (Least Recently Used) . . . . .	80	Tape-based file system . . . . .	97
Performance . . . . .	81	Problems . . . . .	97
Monotonic strategies . . . . .	82	Blocking . . . . .	97
What can a page fault mean? . . . . .	83	Disks . . . . .	97
How much memory? . . . . .	83	Surface mapping . . . . .	97
Optimization Problem . . . . .	84	Hopscoch . . . . .	97
Queuing Convention . . . . .	85	Index sequential files . . . . .	98
Global vs Local thrashing . . . . .	86	Direct access files . . . . .	98
		Representing files . . . . .	99

Continuous allocation . . . . .	99	RAID 1 . . . . .	113
Disk Partitions . . . . .	99	RAID 2 . . . . .	113
Linked allocation . . . . .	99	RAID 5 . . . . .	114
Separate links from blocks . . . . .	99	Reading ahead (for file operations) . . . . .	114
Bitmap . . . . .	100	Plugging . . . . .	115
FAT . . . . .	100	How it works . . . . .	115
UNIX file system . . . . .	101		
File identification . . . . .	101		
Directory structure . . . . .	102		
Links . . . . .	102		
Free space management . . . . .	103		
Modern Enhancements . . . . .	104		
Extents . . . . .	108		
B-Trees . . . . .	108		
WAFL . . . . .	108		
Snapshots . . . . .	109		
<b>12: Disk scheduling</b>	<b>111</b>		
SSF: Shortest Seek First . . . . .	111		
IS SSF optimal . . . . .	111	Credits . . . . .	121
LOOK strategy . . . . .	112	General . . . . .	121
Multiple Disks . . . . .	112	Art . . . . .	121
RAID . . . . .	113	Disclaimer . . . . .	121
<b>13: Virtual Machines</b>	<b>116</b>		
Soft Virtualization . . . . .	116		
Cygwin . . . . .	116		
OS Virtual Machines . . . . .	116		
Example: System call . . . . .	117		
I/O request . . . . .	117		
Memory mapping . . . . .	117		
The goal of virtualizing OSes . . . . .	118		
Fucking Intel . . . . .	118		
Binary rewriting . . . . .	119		
Hardware assist (Intel's VT-x) . . . . .	119		
<b>Professor's Villains and Credits</b>	<b>120</b>		

# 1: INTRODUCTION AND OVERVIEW

An operating system is a bunch of software layers separating your programs from the computer hardware.

## ABSTRACTION

The idea that some details are removed from the user. It removes details that the user does not need to worry about. This however can lead to performance loss as the user loses some of the control over the hardware.

## EXAMPLE

---

Files are an abstraction of information that are stored on persistence storage, i.e. non-volatile.

To maintain or improve performance, the operating system needs to control the hardware optimally. The operating system is responsible to make the user experience fast and seamless even with limiting hardware. Moreover, the user doesn't know this process is happening.

## COMPILING, LINKING, LOADING

### COMPILING

---

When compiling the source code the GCC compiler is converting the human-readable code into a **object file**. This object file or **modules** is a tentative executable, it is missing certain things like functions from external files, etc. To represent these missing information, the object file contains symbolic information in the form of a symbol table. It also contains procedures to fix the object file into a fully functioning executable.  
The object file is in a specific format called ELF or Executable and Linking Format.

### LINKING

---

Linking takes the objects file(s), and using the information from the object file(s) fix and link the files together to create a fully functioning executable that can be loaded into memory.

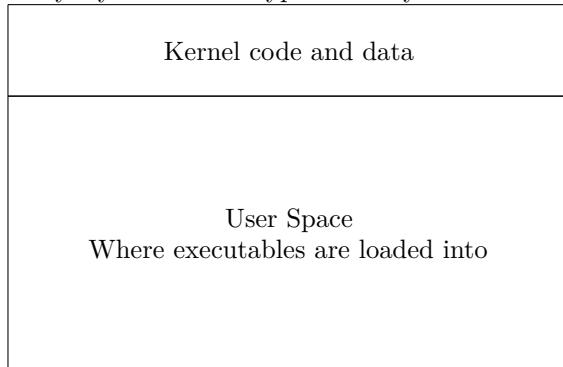
## LOADING

---

Loading takes the executable and loads the information into memory. The actual location of the executable is flexible but not arbitrary. Each platform/OS use specific conventions and format where an executable will go. The loader uses information from the previous steps to indicate where and how to place (or initialize) the various memory areas of the process-to-be. Some libraries are not loaded until needed by the program, these are dynamically loaded libraries.

## INTRO TO X86 MEMORY LAYOUT

In the figure below shows the memory layout of the a typical x86 system.



The Kernel space is protected and have privileged access to certain services while the user space does not have privileged access or is un-privileged.

To access the kernel code and use the kernel's services a system call is called by the program where it runs a subroutine in the kernel space protected from the user. A system call is not the same as a function call as a system call calls a subroutine in the operating system while a function call will call a subroutine in the program. Moreover, a system call will interrupt the current process for the system call process making the program a bit slower then its program's equivalent.

## BASIC STORAGE CLASSES IN C

Here is a sample code describing the storage classes in C

```
1 int read_something (void);
2 int do_something (int);
3 void write_something (const char *);
4
5 int some_global_variable;
6 static int some_local_variable;
7
8 main () {
9     int some_stack_variable;
10    some_stack_variable = read_something ();
11    some_global_variable = do_something (some_stack_variable);
12    write_something ("I am done");
13 }
```

```
int some_global_variable;
```

This variable is globally visible by all modules (common, bss)

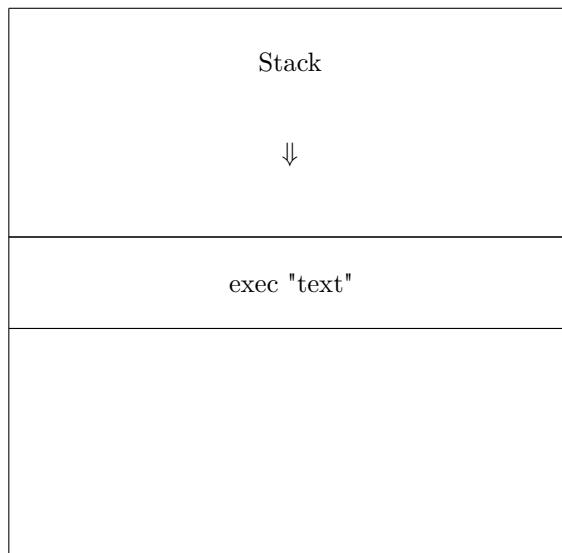
```
static int some_local_variable;
```

This variable is globally visible by all function within the current module, but not outside the module (data)

```
int some_stack_variable;
```

This variable is allocated on stack, visible only within the block (called auto for automatic)

## MEMORY LAYOUT IN THE USER SPACE



The stack grows downward for this system.

# 2: ELF: EXECUTABLE AND LINKING FORMAT

## INTRODUCTION

This describes the format for object files compiled in GCC.

### RELATIONSHIP WITH OTHER ELVES

Note: These ELF's are not related to the Dungeons and Dragon's variety of elf, nor are they related to Game Workshops' varieties of Aeldari from Warhammer 40k, or the Asur, Druchii, or Asrai from Warhammer Fantasy.

### DWARF

Is a common debugger format used to debug ELF.

Like the relationship status of the ELF these DWARFs are not related to the Dungeon and Dragon's variety of dwarf, nor are they relate to Game Workshops' varieties of Squats from Warhammer 40k, or the many dwarf species from Warhammer Fantasy.

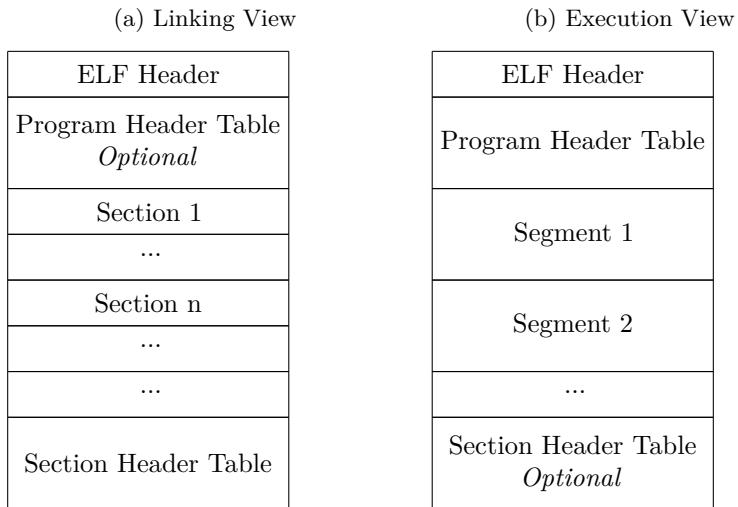
### THE THREE MAIN TYPES OF AN OBJECT FILE

1. **Relocatable file:** holds code and data suitable for linking with other object files to create an executable or a shared object file.
2. **Executable file:** holds a program suitable for execution.
3. **Shared object file:** holds code and data suitable for linking in two contexts. First, the link editor may process it with other relocatable and shared object files to create another object file. Second, the dynamic linker combines it with an executable file and other shared objects to create a process image.

Created by the assembler and link editor, object files are binary representations of programs intended to execute directly on a processor. Programs that require other abstract machines are excluded.

## FILE FORMAT

Because this format is used in the linking (building a program) and execution (running a program). The object file format provides parallel views of a file's contents, reflecting the different needs of these activities.



## ELF HEADER

---

An *ELF header* resides at the beginning and holds a "road map" describing the file's organization, it also contains the number of sections in the file. Sections hold the bulk of object file information for the linking view: instructions, data, symbol table, relocation information, etc.

## PROGRAM HEADER TABLE

---

A *Program header table*, if present, tells the system how to create a process image (execute a program). Files used to build a process image must have a program header table; relocatable files do not need one.

## SECTION HEADER TABLE

---

A *section header table* contains information describing the file's sections. Every section has an entry in the table, each entry gives information such as the section name, location, size, and so on. Files used during linking must have a section header table, other object files may or may not have one.

## EXAMPLE OF AN ELF

---

Here is an example of an ELF format object file for this source code below:

```

1 int read_something (void);
2 int do_something (int);
3 void write_something (const char*);
4
5 int some_global_variable;
6 static int some_local_variable;
7
8 main () {
9     int some_stack_variable;
10    some_stack_variable = read_something ();
11    some_global_variable = do_something (some_stack_variable);
12    write_something ("I am done");
13 }
```

# ELF format

## SIGNATURE

```
000 7f 45 4c 46 01 01 01 00 00 00 00 00 00 00 00 00 00  
-- E L F 32 LE FV -----
```

## ELF HEADER

```
010 01 00 03 00 01 00 00 00 00 00 00 00 00 00 00 00 00  
0001 0003 00000001 00000000 00000000  
↑ target architecture ↑ program header (executable module only)  
relocatable starting address for execution (executable module only)  
module  
020 04 01 00 00 00 00 00 00 34 00 00 00 00 00 28 00  
00000104 00000000 0034 0000 0000 0028  
SHT offset CPU flags hdr no PHT SHT entry size  
length  
030 0b 00 08 00  
000b 0008  
↑ SHT = Section Header Table  
index (into SHT) of the string section containing section names  
PHT = Program Header Table  
number of entries in SHT
```

## WHAT DOES THESE MEAN?

- **Signature:** Tells the linker that this is a ELF formatted object file and defines whether to use 32 or 64 bit addresses. In this case it uses a 32-bit address.
- **SHT offset:** Is the relative offset from the start of the file to the start of the Section Header Table.
- **SHT entry size:** Tells the size of each entry in the Section Header Table. To deal with variable length sections, each entry will have a pointer to the section as well as its size. There are other information as well but that will be described later.
- **# of entries in SHT:** Tells the number of entries in the Section Header Table.
- **index (into SHT) of the string section containing section names:** This is the index of the section in the Section Header Table that contains the location, size, and other information regarding the string section containing section names. This section is called **.shstrtab** or Section Header STRing TABLE.

# ELF contents: SHT (0-1)

## Section 0 (NULL)

```
104          00 00 00 00 00 00 00 00 00 00 00 00  
110 00 00 00 00 00 00 00 00 00 00 00 00  
120 00 00 00 00 00 00 00 00 00 00 00 00
```

Entry size = 28;  $104 + 28 = 12c$   
28 bytes = a words

## Section 1 (.text)

12c 1f 00 00 00 = 0000001f	name pointer (index into a string section)
01 00 00 00 = 00000001	section type (PROGBITS = belongs to the program)
06 00 00 00 = 00000006	flags: in program's memory + executable
00 00 00 00 = 00000000	address in memory (unknown because module is relocatable)
34 00 00 00 = 00000034	offset to section in ELF file
38 00 00 00 = 00000038	size $34 + 38 = 6c$
00 00 00 00 = 00000000	extra link (none - typically points to a related section)
00 00 00 00 = 00000000	extra info (none)
04 00 00 00 = 00000004	address alignment (word boundary)
00 00 00 00 = 00000000	entry size (if the section consists of fixed-size entries)

## WHAT DOES THESE MEAN?

- **name pointer**: Also called sh\_name, is an offset to a null-terminated character sequence in the section name (.shstrtab) section that provides the name for this section.
- **section type**: Also called sh\_type, describes the type of the section, in this case 0x01 means SH\_PROGBITS or Program data.
- **flags**: Also called sh\_flags, identifies the attributes of the section, in this case 06 means in program's memory and executable.
- **address in memory**: Also called sh\_addr, is the virtual address of the section in memory, for sections that are loaded. In this case it is unknown (NULL) because the module is relocatable.
- **offset to section in ELF file**: Also called sh\_offset, is the offset of the section in the ELF file. The offset is relative from the head of the ELF file.
- **size**: Also called sh\_size, is the size in bytes of the section in the ELF file. May be 0 due to compiler optimization, and other things.
- **extra link**: Also called sh\_link, contains the section index of an associated section. This field is used for several purposes, depending on the type of the section.
- **extra info**: Also called sh\_info, contains extra information about the section. This field is used for several purposes, depending on the type of the section.
- **address alignment**: Also called sh\_addralign, contains the required alignment of the section. This field must be a power of 2.
- **entry size**: Also called sh\_entsize, contains the size, in bytes, of each entry, for sections that contain fixed-sized entries. Otherwise, this field is zero.

# ELF contents: section body 8, 5, 1

## Section 8 (STRTAB .shstrtab)

```
0b1    00 2e 73 79 6d 74 61 62 00 2e 73 74 72 74 61  
      -- . s y m t a b -- . s t r t a  
0c0 62 00 2e 73 68 73 74 72 74 61 62 00 2e 72 65 6c  
      b -- . s h s t r t a b -- . r e l  
0d0 2e 74 65 78 74 00 2e 64 61 74 61 00 2e 62 73 73  
      . t e x t -- . d a t a -- . b s s  
0e0 00 2e 72 6f 64 61 74 61 00 2e 6e 6f 74 65 2e 47  
      -- . r o d a t a -- . n o t e . G  
0f0 4e 55 2d 73 74 61 63 6b 00 2e 63 6f 6d 6d 65 6e  
      N U - s t a c k -- . c o m m e n  
100 74 00  
      t --
```

## Section 5 (PROGBITS .rodata)

```
074 49 20 61 6d 20 64 6f 6e 65 00  
      I a m d o n e --
```

## Section 1 (PROGBITS .text)

```
034 55 89 e5 83 ec 08 e8 fc ff ff ff 89  
040 c0 89 45 fc 83 ec 0c ff 75 fc e8 fc ff ff ff 83  
050 c4 10 89 c0 a3 00 00 00 00 83 ec 0c 68 00 00 00  
060 00 e8 fc ff ff ff 83 c4 10 c9 c3 90
```

- As you can see being 0x1f from Section 8 points to the beginning of .text\0. The section name for Section 1.
- In section 5, it's only "I am done\0", .rodata means read-only data as the string that is being used is in a function where its only parameter takes in a constant char array. Therefore the string is immutable.
- Section 1 contains the code segment or the executable code for the program.
- Section 8 is a string table (SH\_STRTAB) that contains the name for each section.

## DATA SEGMENT

### Components in the data segments

**Text** Is where a portion of an object file or the corresponding section of the program's virtual address space contains the executable instructions is stored and is generally read-only and fixed in size.

**Data** contains any global or static variables which have a pre-defined or initialized value and can be modified.

**BSS** Uninitialized data.

**Heap** Heap space.

**Stack** Call stack.

## ELF contents: sections a and 7

```
Section a (STRTAB .strtab)
39c                                     00 70 72 6f
                                         p r o
3a0 67 72 61 6d 2e 63 00 73 6f 6d 65 5f 6c 6f 63 61
   g r a m . c -- s o m e _ l o c a
3b0 6c 5f 76 61 72 69 61 62 6c 65 00 6d 61 69 6e 00
   l _ v a r i a b l e -- m a i n --
3c0 72 65 61 64 5f 73 6f 6d 65 74 68 69 6e 67 00 64
   r e a d _ s o m e t h i n g _ d
3d0 6f 5f 73 6f 6d 65 74 68 69 6e 67 00 73 6f 6d 65
   o _ s o m e t h i n g _ s o m e
3e0 5f 67 6c 6f 62 61 6c 5f 76 61 72 69 61 62 6c 65
   g l o b a l _ v a r i a b l e
3f0 00 77 72 69 74 65 5f 73 6f 6d 65 74 68 69 6e 67
   w r i t e _ s o m e t h i n g
00
```

```

Section 7 (PROGBITS .comment)
07e                                     00 47
                                         G
080 43 43 3a 20 28 47 4e 55 29 20 33 2e 33 2e 33 20
    C   C   (   G   N   U   )   3   .   3   .   3
090 32 30 30 34 30 34 31 32 20 28 52 65 64 20 48 61
    2   0   0   4   0   4   1   2   (   R   e   d   )   H   a
0a0 74 20 4c 69 6e 75 78 20 33 2e 33 2e 33 2d 37 29
    t   L   i   n   u   x   3   .   3   .   3   -   7
0b0 00

```

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## ELF contents: sections 9 and 2

Section 9 (.symtab)									
2bc	00	00	00	00	00	00	00	00	00
=====	=====	=====	=====	=====	=====	=====	=====	=====	=====
01	00	00	00	00	00	00	00	04	00 f1 ff
=====	=====	=====	=====	=====	=====	=====	=====	=====	=====
name	value	size	info	SHT	IDX				
00	00	00	00	00	00	00	03	00	01 00
=====	=====	=====	=====	=====	=====	=====	=====	=====	=====
00	00	00	00	00	00	00	03	00	03 00
00	00	00	00	00	00	00	03	00	04 00
00	00	00	00	00	00	00	03	00	05 00
0b	00	00	00	00	00	04	00	00	01 00 04 00
00	00	00	00	00	00	00	03	00	06 00
00	00	00	00	00	00	00	03	00	07 00

Each entry consists of:  
3 words  
2 bytes  
1 half-word

These correspond to sections - those that describe the contents of the program's memory.

Section 2 (.rel.text)  
404 07 00 00 00 02 0a 00 00  
===== =====  
offset relocation type  
17 00 00 00 02 0b 00 00  
===== =====  
21 00 00 00 01 0c 00 00  
29 00 00 00 01 05 00 00  
2e 00 00 00 02 0d 00 00

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski 25

- Section 1 is a string table that contains null-terminated character sequences that contain the names for the symbol table.
  - Section 7 is part of the program and contains comment about the compiler used. - Section 9 is the symbol table for the EFL file.
  - Section 2 is .rel.text or relocatable .text

## SYMBOL TABLE

Each entry in the symbol table is made up with 3 words (32-bits in this case), 2 bytes, and 1 half-word (16-bits in this case). These represent the name, value, size, info and the Section Header Table index respectively.

- name is an offset from the head of the .strtab section that provides the name for each entry, NULL values are possible.
  - value is the value for that symbol if its a variable.
  - size is the size in bytes for that symbol if its a piece of code.
  - info gives information about the symbol, is it a global function, global variable, etc.
  - SHTIDX is the index in the section header table to the section that the symbol is related to.

## AN EXAMPLE

The symbol for main is called main as the offset of 0x1f in the .strtab is .main, it has no value as its not a variable, it has a size of 61 byte, its a global function, and its related to the .text section.

```

0: 55          push    %ebp
1: 89 e5       mov     %esp,%ebp
3: 83 ec 08    sub    $0x8,%esp
6: e8 fc ff ff ff  call   7 <main+0x7>
b: 89 c0       mov     %eax,%eax
d: 89 45 fc    mov     %eax,0xfffffff(%ebp)
10: 83 ec 0c   sub    $0xc,%esp
13: ff 75 fc   pushl   0xfffffff(%ebp)
16: e8 fc ff ff ff  call   17 <main+0x17>
1b: 83 c4 10   add    $0x10,%esp
1e: 89 c0       mov     %eax,%eax
20: a3 00 00 00 00  mov    %eax,0x0
25: 83 ec 0c   sub    $0xc,%esp
28: 68 00 00 00 00  push   $0x0
2d: e8 fc ff ff ff  call   2e <main+0x2e>
32: 83 c4 10   add    $0x10,%esp
35: c9         leave
36: c3         ret

```

This is the .text section disassembled back into x86 assembly. An address doesn't have to be located outside the current module to be essentially unknown. The exact location of the module in memory is only known when the **complete** program is loaded for execution. The benefits of DLL are many:

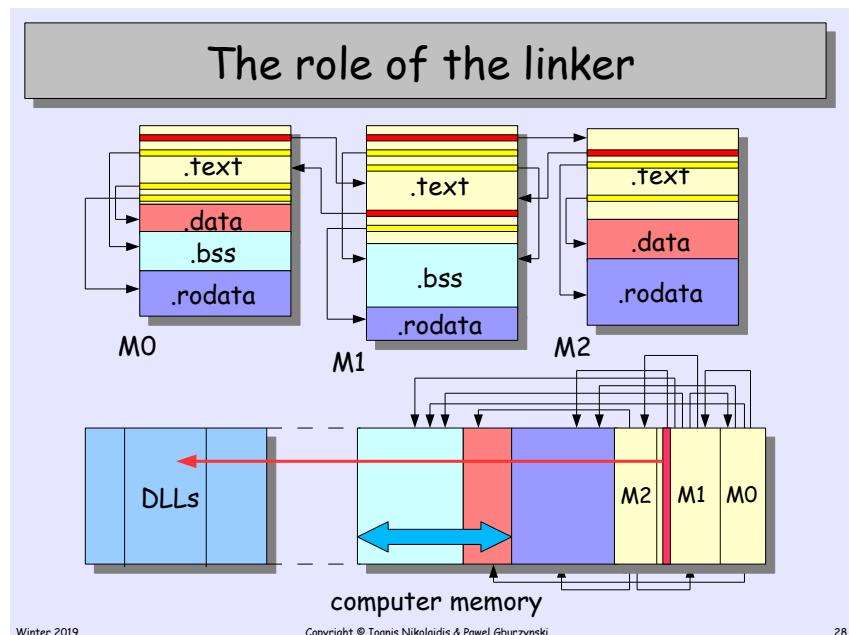


Figure 2: The role of the linker is to combine the modules and consolidate each data segment from the modules together, so for example all .data sections from the modules that have one will be grouped together and have the pointers in the sections updated to those new locations.

### THE BENEFITS OF DLL

- **Speed of Loading/executable size:** Not everything needs to be loaded during execution.
- **Incremental Loading (as needed):** Things are loaded during run time, therefore increasing the speed of loading
- **Debugging:** You don't have to recompile everything, just the modules that are being tested/debugged.



## The .text section

```
Section b (.text) = 2ac8 + c * 28
2ca8 80 00 00 00 = 00000080 name pointer
01 00 00 00 = 00000001 section type (PROGBITS)
06 00 00 00 = 00000006 flags: alloc + exec
90 83 04 08 = 08048390 address in memory
90 03 00 00 = 00000390 offset in ELF file
d0 01 00 00 = 000001d0 size
00 00 00 00 = 00000000 extra link (none)
00 00 00 00 = 00000000 extra info (none)
01 00 00 00 = 00000001 address alignment (byte)
00 00 00 00 = 00000000 entry size
```

Sections representing memory chunks, like `.text`, `.data`, `.rodata`, are simply pieces to be written into the respective locations in memory to build the program's executable image. Some allocatable sections, like `.bss`, have no image in the file: `size` and `location` are the only attributes that matter.

Winter 2019

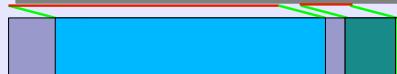
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Things are a bit more complicated:

Program Header: (this is what in fact describes how the program is loaded)

```
PHDR off 0x000000034 vaddr 0x08048034 paddr: 0x08048034 align 2**2
    filesz 0x000000c0 memsz 0x000000c0 flags r-x
INTERP off 0x0000000f4 vaddr 0x080480f4 paddr: 0x080480f4 align 2**0
    filesz 0x00000013 memsz 0x00000013 flags r-
LOAD off 0x000000003 memsz 0x000000003 flags r-x
LOAD off 0x00000059d vaddr 0x08048000 paddr: 0x08048000 align 2**12
    filesz 0x00000059d memsz 0x0000059d flags r-x
LOAD off 0x0000005a0 vaddr 0x080495a0 paddr: 0x080495a0 align 2**12
    filesz 0x00000014 memsz 0x00000014 flags rw-
DYNAMIC off 0x000000568 vaddr 0x080495ec paddr: 0x080495ec align 2**2
    filesz 0x000000c8 memsz 0x000000c8 flags rw-
NOTE off 0x000000108 vaddr 0x08048108 paddr: 0x08048108 align 2**2
    filesz 0x00000020 memsz 0x00000020 flags r--
```

ELF HDR PHDR INTERP NOTE ... .text .data



When the loader is called to load the executable into memory, it will load the correct files into the logical memory in the most appropriate way, for example DLLs may not be loaded if they are not being used. Each file that are being loaded can have read, write, and execute permissions set depending on it's content. Certain files have to be aligned to the memory's page size.

## How do we interface with the OS?

### 80484d0 <read\_something>:

```
int read_something (void) {
    int res;
    scanf ("%d", &res);
    return res;
}
80484d0:
    55
    89 e5
    83 ec 08
    83 ec 08
    8d 45 fc
80484dc:
    50
80484dd:
    68 92 85 04 08
80484e2:
    e8 59 fe ff ff
80484e7:
    83 c4 10
80484ea:
    8b 45 fc
80484ed:
    c9
80484ee:
    c3
push    %ebp
mov     %esp, %ebp
sub    $0x8, %esp
sub    $0x8, %esp
lea     0xfffffffffc(%ebp), %eax
push    %eax
push    $0x8048592
call    8048340 <_init+0x38>
add    $0x10, %esp
mov     0xfffffffffc(%ebp), %eax
leave
ret
```

### 08048320 <.plt>: (stands for Procedure Linkage Table)

```
08048320: ff 35 c8 95 04 08    pushl  0x80495c8
08048326: ff 25 cc 95 04 08    jmp    *0x80495cc
0804832c: 00 00
0804832e: 00 00
08048330: ff 25 d0 95 04 08    jmp    *0x80495d0
08048336: 68 00 00 00 00
0804833b: e9 e0 ff ff ff
08048340: ff 25 d4 95 04 08    jmp    *0x80495d4
08048346: 68 08 00 00 00
0804834b: e9 d0 ff ff ff
.....
push    $0x0
jmp    8048320 <_init+0x18>
push    $0x8
jmp    8048320 <_init+0x18>
```

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

38

When scanf is called in the program, the PC does not directly go to the actual function. Because this function is at an unknown location (can be a DLL). The PC first goes to the procedure linkage table as a level of redirection where it then jumps to the entry point of scanf.

# HOW ARE DLLS LOADED

Libraries that are dynamically linked and loaded often have their own methods of being loaded into memory. One way is by lazy linkage.

### Lazy linkage

```
08048320 <.plt>: (stands for Procedure Linkage Table)
08048320: ff 35 c8 95 04 08 pushl 0x80495c8
08048326: ff 25 cc 95 04 08 jmp *0x80495cc
0804832c: 00 00 add %al, (%eax)
0804832e: 00 00 add %al, (%eax)
08048330: ff 25 d0 95 04 08 jmp *0x80495d0
08048336: 68 00 00 00 00 push $0x0
0804833b: e9 e0 ff ff ff jmp 8048320 <_init+0x18>
08048346: 68 08 00 00 00 jmp *0xb0495d4
0804834b: e9 d0 ff ff ff jmp 8048320 <_init+0x18>
...

```

Contents of section .got: (global offset table)

```
080495c4 080495ec 00000000 00000000 08048336
080495d4 08048346 08048356 08048366 08048376
080495e4 08048386 00000000
```

This address will be replaced on the first call by the actual address of the called function. Just by looking at the executable, we have no clue where it will be located. These days, it is a standard practice to put all system functions in DLLs.

When the program is loaded, the loader identifies:  
 the "interpreter" responsible for dynamic linkage  
 the "DLLs" needed by the program  
 this information is available in the ELF headers

All calls to DLL functions are turned into lazy linkage pointers involving a call through the `.plt` section, and an offset in the `.got` section.

The initial setting of the offset triggers a jump to the interpreter with parameters that allow it to identify the DLL and the called function.

Before calling the function, the interpreter replaces the offset in `.got` with the direct pointer to the function, so that the next call will be straightforward (or as straightforward as it can get under the circumstances 😊). `.got` it?

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

39

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

40

### Figures can be more illuminating

(a) calling a subroutine is straightforward in a statically linked executable

(b) the layout is drastically different in the case of dynamically linked and loaded libraries (we will trace the first invocation to `disubr`)

### Invoking the dynamic linker/loader

(c) the indirect jump passes control over to the dynamic linker/loader

(d) we skip the details and go directly to the end result of invoking the dynamic linker/loader

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

41

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

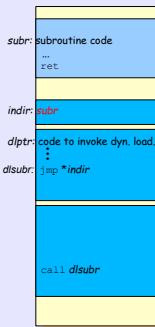
42

- b) During the first function call of the procedure in the executable. The linker first adds a level of indirection to the procedure linkage table. The **Procedure Linkage Table or .plt** section contains code that will dynamically load the function into memory.
- c) Because the new function can be placed really anywhere in memory the **Global Offset Table or .got** is used as another layer of indirection to point the jump call from before to the entry point of the function that will load the requested function into memory.
- d) Once the function has been loaded into memory the `.got` is updated so that it points to the entry point of the function. Note: the return address has not been touch, it still points back to the main executable.

3: DLL

13

## Following through the 1<sup>st</sup> invocation



(e) for convenience let's assume that the execution after invoking the loader continues to the same indirect jump we had encountered earlier

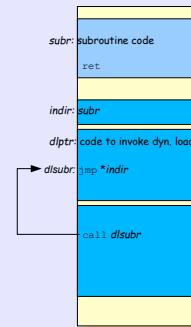


(f) upon (re-executing) the jump we branch to the address of the function that has been loaded

## The final touches



(g) this step looks obvious but it implies that the state of the stack when the `jmp` (figure (f)) was executed is as it was after the `call` (figure (b))



(h) the second and subsequent calls (from the same point or elsewhere in the code) to `subr` pay only the cost of the one additional indirect jump instruction

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

43

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

44

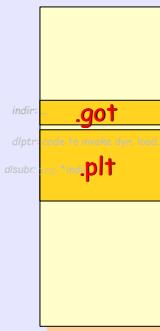
e) Right after the code that loaded the function into memory is the same jump command that points to the same entry in the .got table, only this time it had been updated to point to the entry point of the function.

f) When we recall the same jump command instead of reloading the function again which is idiotic it redirect the PC to the entry point of the loaded function.

g-h) As this loading process only used pure jump commands the return address register was not effected.

Therefore, calling return in the subroutine updates the PC back to the main executable. All subsequent calls only needs to pay one additional level of redirection to execute the subroutine.

## The names for the two special areas



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Now for the easy way

gcc -o executable -static program.o extras.o

Executable size:	dynamically linked	= 14221 bytes
	statically linked	= 1691082 bytes

```

8048218 <read_something>:
8048218: 55          push  %ebp
8048219: 89 e5        mov    %esp,%ebp
804821b: 83 ec 08      sub    $0x8,%esp
804821e: 83 ec 08      sub    $0x8,%esp
8048221: 8d 45 fc      lea    0xfffffff(%ebp),%eax
8048224: 50          push  %eax
8048225: 68 f2 e4 08 08  push  $0x808e4f2
804822a: e8 cd 04 00 00  call   80486fc <scanf>
804822f: 83 c4 10      add    $0x10,%esp
8048232: 8b 45 fc      mov    0xfffffff(%ebp),%eax
8048235: c9          leave 
8048236: c3          ret    
8048237:                   . . .
80486fc <scanf>:
80486fc: 55          push  %ebp
80486fc: 89 e5        mov    %esp,%ebp
80486fc: 83 ec 08      sub    $0x8,%esp
80486fc: 83 ec 08      sub    $0x8,%esp
80486fc: 8d 45 fc      lea    0xfffffff(%ebp),%eax
80486fc: 50          push  %eax
80486fc: 68 f2 e4 08 08  push  $0x808e4f2
80486fc: e8 cd 04 00 00  call   80486fc <scanf>
80486fc: 83 c4 10      add    $0x10,%esp
80486fc: 8b 45 fc      mov    0xfffffff(%ebp),%eax
80486fc: c9          leave 
80486fc: c3          ret    
80486fc:                   . . .
8048711: e8 ba 34 01 00  call   805bbd0 <_IO_vfscanf>
8048711:                   . . .

```

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

46

You can statically link all the libraries together to get a massive executable if you don't want to deal with this mess. It also increases portability (slightly) at the cost of size.

This is the function call for scanf. It looks complicated but the thing you need to know is the int \$0x80 system call. The %eax identifies which system call is invoked

```
080651c0 <__libc_read>:  
080651c0: 53          push  %ebx  
080651c1: 8b 54 24 10  mov   0x10(%esp,1),%edx  
080651c5: 8b 4c 24 0c  mov   0xc(%esp,1),%ecx  
080651c9: 8b 5c 24 08  mov   0x8(%esp,1),%ebx  
080651cd: b8 03 00 00  mov   $0x3,%eax  
080651d2: cd 80        int   $0x80  
080651d4: 5b          pop   %ebx  
080651d5: 3d 01 f0 ff ff  cmp   $0xfffff001,%eax  
080651da: 0f 83 70 da fe ff  jae   8052c50 <__syscall_error>  
080651e0: c3          ret
```

When we want input from the user the function scanf uses a system call that will transfer control from the program or the user in user space to the kernel that will retrieve the information on your behalf. This is to improve security, and reliability as the kernel developer is responsible for this function. When the system call returns control is given back to the program.

Hardware can also call kernel functions by the means of an external software interrupt. This is usually handled by an interrupt handler.

When you create a new process, you fork the parent process so that the running process which can be completely different is a child of the parent process. The only exception to this idea is the first process which is created differently.

## Blocking system calls

Because of the inability to proceed with the execution of the system call (see the last `read` call on the previous slide), the process that issued the system call needs to be (temporarily) blocked from continuing execution. Therefore we need:

- a mechanism to **switch** from the current process to another one that is not blocked (for productive and efficient use of the CPU)
- a mechanism to **unblock** the blocked processes when the reason (internal or external event) they were waiting for is satisfied (and to continue its execution in the CPU)
- genuinely **external events** are asynchronously delivered to the CPU via interrupts (rarely, via polling); **internal events** originate from the execution of other code (typically, other processes); inter-process communication (IPC) mechanisms provide some of those internal events

\*\* The term "event" is overloaded. Here we use it informally.

# 4: PROCESSES

A process and its children are formed into a tree.

In modern OSes Processes are running concurrently with a "seamless" experience.

## DEFINITIONS

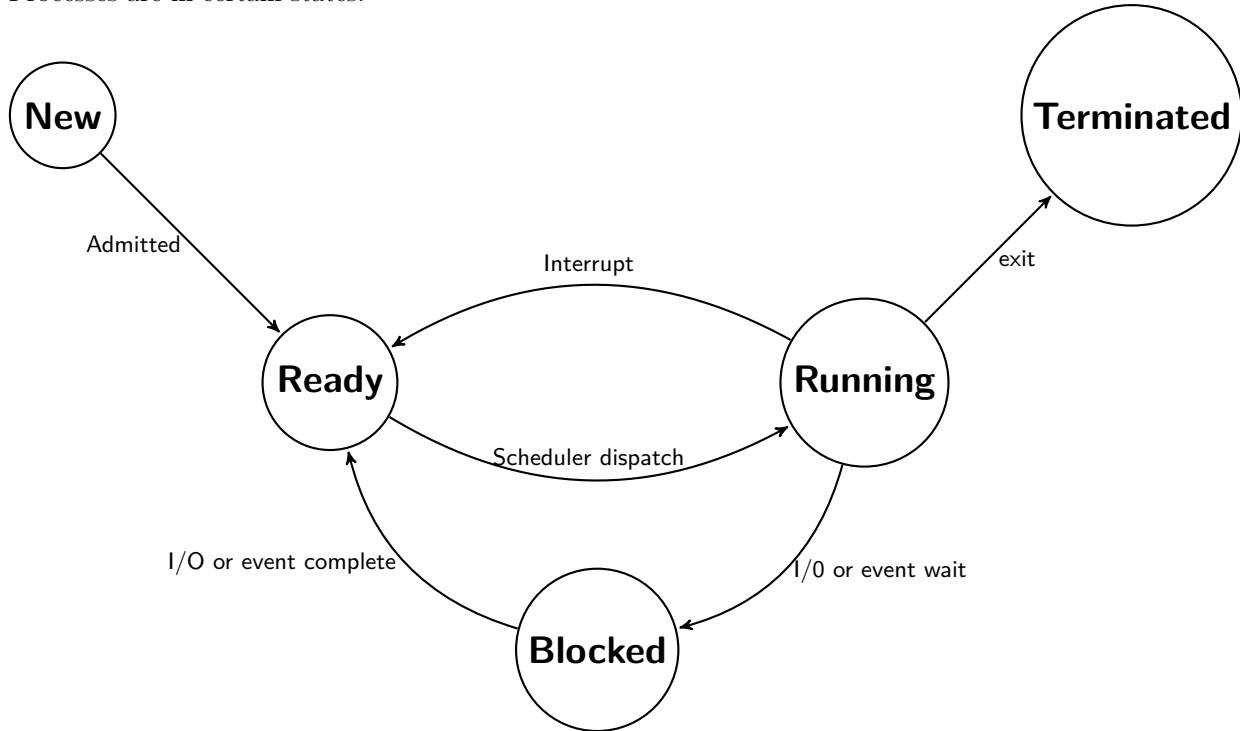
**Process:** A program executing in some environment.

**Program Executing:** text, bss, heap, and stack segments in memory.

**Context:** Shell environment variables, user privilege access, etc.

## CONCEPTUAL FRAMEWORK

Processes are in certain states.



When a new process is created the OS will try to allocate the resources to that process.

Once the resources are allocated the process is ready to be ran by the scheduler. When the process is running its in the running state, when the process is waiting for I/O or events it goes to the Blocked state where it then goes to the ready state once the event or I/O is completed. When OS does not require the process it is terminated.

# PROCESS CONTROL BLOCK

Contains information required by the scheduler to work with multiple processes. Each block contains information about the process's state, pid, PC, register being used, memory limit, etc.

## PROCESS QUEUES

### DIFFERENT TYPES OF QUEUES

**Ready Queue:** Are a queue used for processes in the ready state.

**I/O queue:** Are a queue used for processes in the waiting state.

## COMPONENTS OF A PROCEDURE

There are many components involved within a process. A process contains a **Memory image**, **The state of execution**, **Resources Allocated**, etc.

## CONTEXT SWITCHING

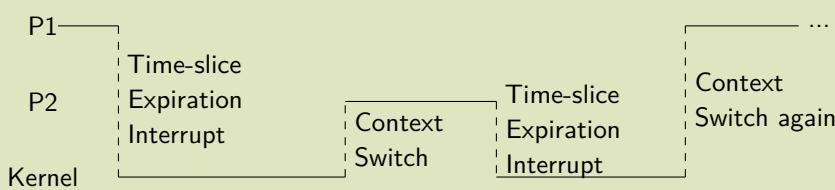
### CONTEXT SWITCHING

**Context Switching:** Is the transition from one process to another. The current process is **blocked** or paused while switching to the other process.

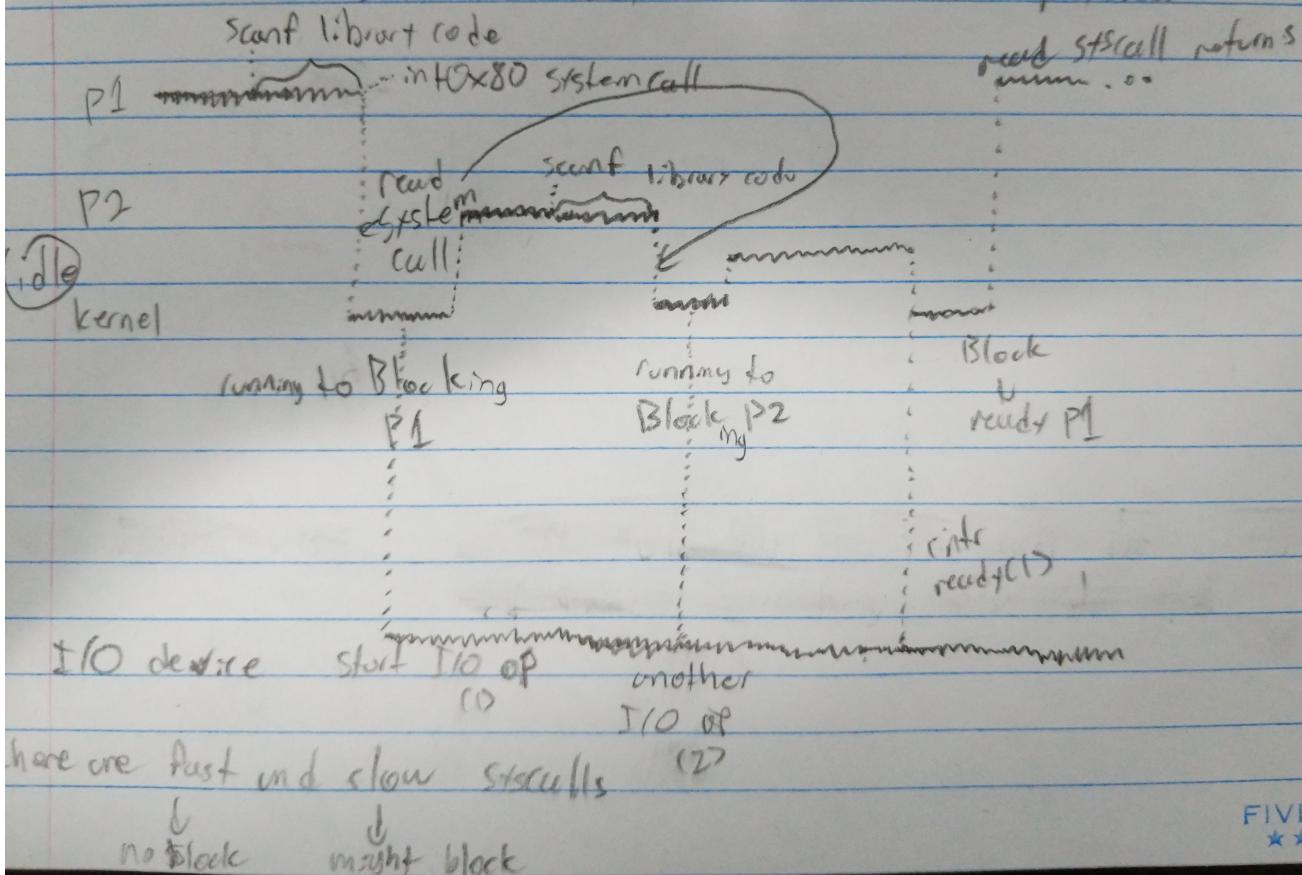
The timer used to tell the scheduler when to switch the process is done using hardware built into the processor. Once the timer has elapsed a certain amount of time it triggers a hardware trap that indicate to the scheduler that a new process should be processed.

### EXAMPLE

Here is an example of context switching between process P1, P2 and the kernel.



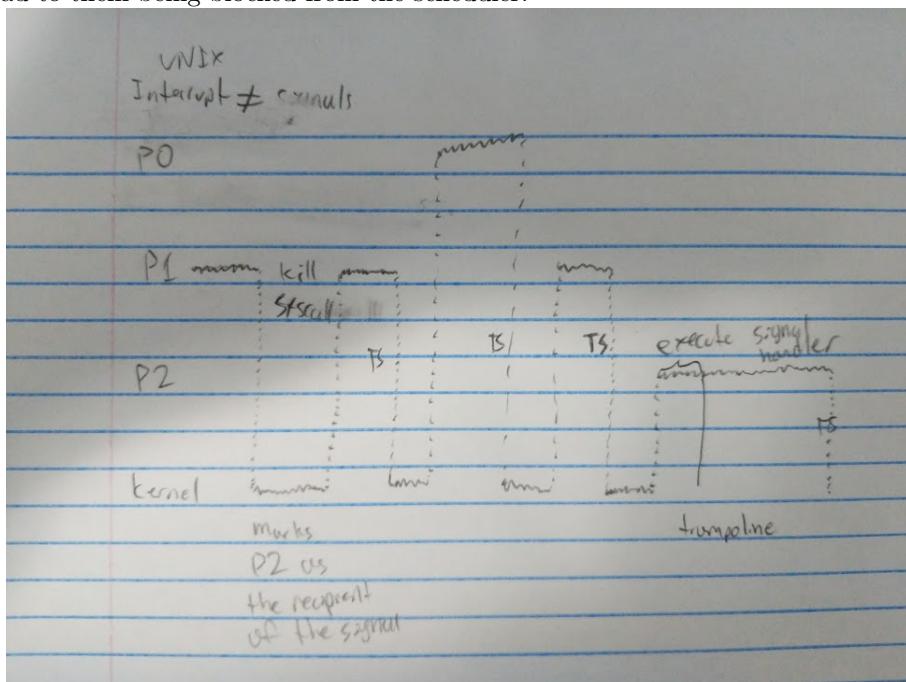
not to scare the countdown timer is handled by the kernel  
but the actual timer is physical.

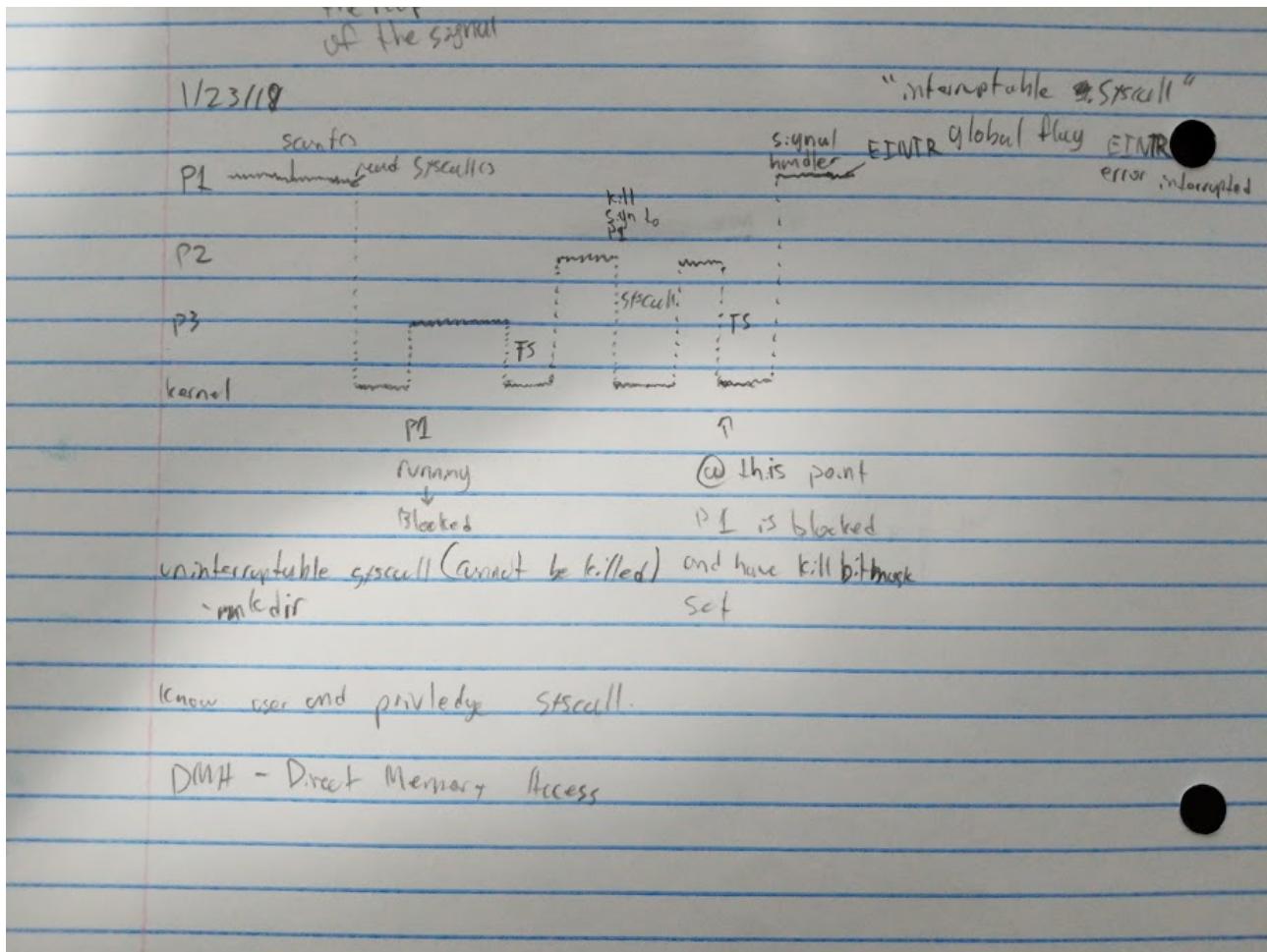


FIVE STARS

In the olden days there would be an idle process that the scheduler will jump into doing nothing if there were no processes in the ready state.

There are system calls that are fast and system calls that are slow. System calls that involve I/O operations are slow and might lead to them being blocked from the scheduler.





The diagram shown above shows a case where a process (P1) is blocked during which a kill signal is sent by process (P2). The default behaviour for most cases is that the signal handler in P1 will handle the signal and interrupt the `read()` system call but will have the `EINTR` (Error interrupt) flag set as the process is being interrupted.

Most system calls are interruptible, meaning the system call can be blocked if a signal is received. Some signals are not interruptible and therefore cannot be killed by standard means. For example `mkdir` is an uninterruptible system call.

## CPU STATE

There are some instructions that are illegal for the user in **User Mode** to execute. These instructions are called privileged instructions and the CPU must switch to the **kernel, or supervisor mode**. In this mode the CPU can execute these instructions that have direct access to fundamental resources like peripheral devices, memory mapping, CPU state, etc.

To gain access to this mode the user must invoke a special machine instruction (system call) that simultaneously calls a special "function" and switch the CPU state into supervisory mode. The actual flag that set this mode is stored in the **PSW (Program Status Word)**.

We generally have two modes to improve security and reliability as different users will have different processes running and you don't want one program to kill the rest of programs.

# RESOURCE MANAGEMENT

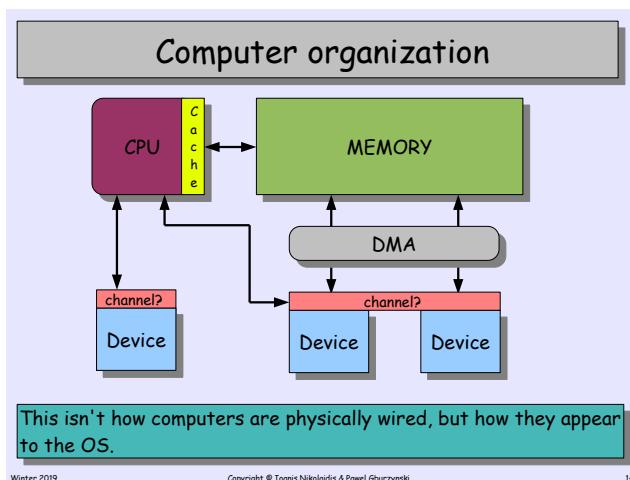
According to some point of view, operating systems are all about resource management. In the early days of computing, humans were the slowest links in the chain of computing as they took valuable time with loading and unloading of programs, debugging and so on.

## TWO IDEAS CHANGED THE WORLD OF COMPUTING

**Batch processing:** remove the clumsy human being from the computer room and have a queue of programs to run.

**Multi-programming:** find something to do for as many hardware components as possible. This only works if the components can operate **autonomously** and does not need to wait for other components to be finished.

# COMPUTER ORGANIZATION



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Górzynski

14

## COMPONENT DEFINITIONS

**Channels** are special "processors" associated with an independent peripheral device. They connect the device to other components in the computer. Some channels can have multiple devices connected to a bus.

**DMA or Direct Memory Access** is a component that will load information between the computer's memory and its devices. This is beneficial as there could be a lot of information from the devices that will bog down the CPU if the CPU has to do loading and storing of everything.

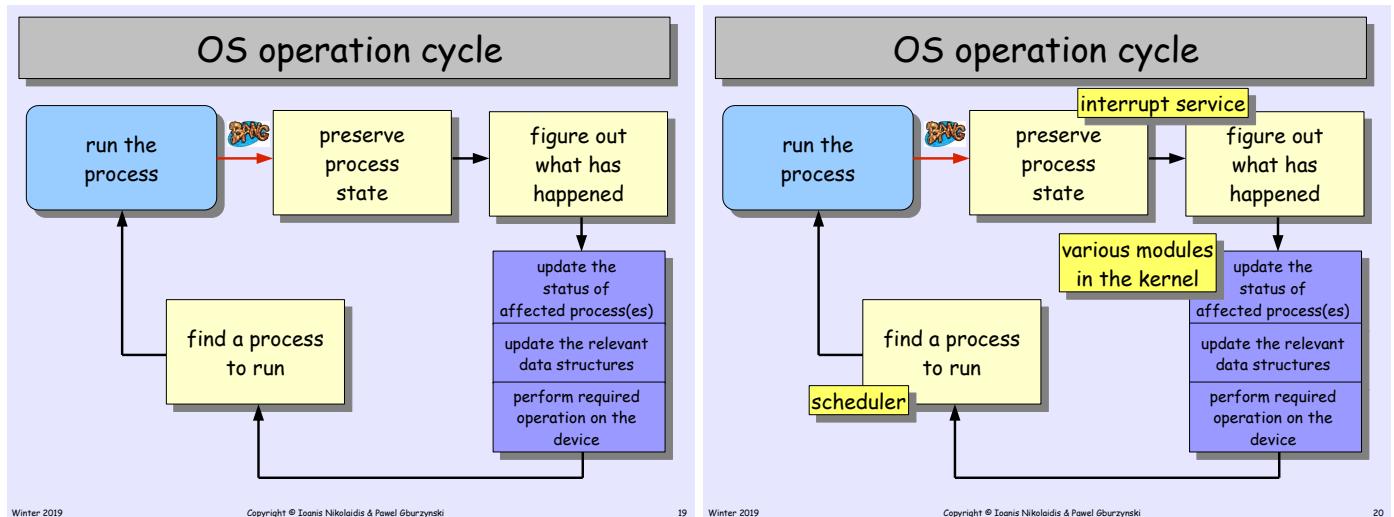
# ACTIVITIES

Typically, during its lifetime, your program can be **preempted** and then **resumed**, possibly many times. The system can handle other activities while running your program. There are computers today with multiple cores that can handle multiple processes, but in general we can consider the CPU as a "single" resource that the processors are competing for. On general OSes are designed such that they provide **fair** and **efficient** strategies for handling competing for resources. These competitors are called **processes**.

## TERMS

- **Processes:** Identifiable activities that can be stopped and resumed. They are described by special data structures and characterized by well-known sets of attributes. One such attribute is the process status
- **Interrupt service routines + syscall:** Activities to handle events in the system. Everything that happens in the system that may affect the status of a process boils down to the execution of an interrupt service routine
- **Others?:** Nope! The entire kernel is in fact a bunch of interrupt service routines, sometimes in minor disguise.

# OS OPERATING CYCLE



## PROCESS CONTEXT

Informally, (process) context is the answer the following question:

What needs to be saved such that a process can be suspended/frozen at an arbitrary point of its execution, such that it can be resumed at the exact same point later "as if nothing happened". Other than side-effects due to the passing of natural time, e.g., the current time advancing, the process is unable to determine on its own that this suspension/resumption cycle ever happened? (An "amnesia" of sorts.)

The answer must include:

- The state of the address space (including heap, stack, etc.)
- The state of all the general purpose registers used
- The state of special purpose registers (PC, SP, etc.)
- The state of the process used resources (file ptr, etc.)

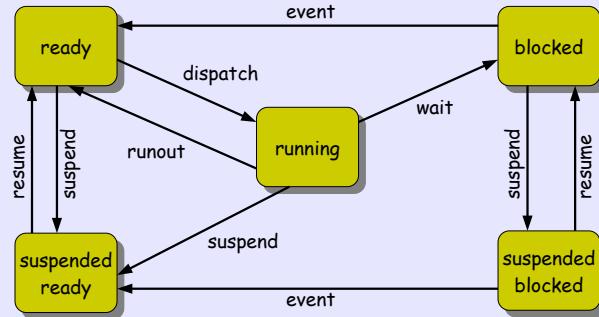
The transition between one context and another is called **Context Switch** and it generally considered an "expensive" operation. When terminating a process someone (usually the parent) needs to capture its exit code.

If the exit code is not caught by anyone then the PCB is not removed from the process queue and is considered a zombie process as it should be dead but the process is still alive by the OS

## SUSPENDED STATE

Sometimes we want to remove a process from the scheduler entirely but still want to keep its instance or state. Processes in suspended states are not considered by the scheduler to being operations so it is separated into its own separate pool. The process is then brought back to the scheduler pool once it resume back from suspension.

## A more complete transition diagram



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

27

# THREADS VS PROCESSES

### DIFFERENCE BETWEEN THREADS AND PROCESSES

A **process** expresses a single execution path over the address space. **Threads** are (multiple) independent execution paths that share the same address space.

Each thread has, at a minimum:

- its own stack,
- its own register (incl. Special purpose one) values

Of course there are other luxuries like thread-private data and resources, but strictly speaking, they are not needed to have a multi-threaded system.

There are different types of threads. One performed by user code alone are called **user space threads**, threads performed by the OS are called **OS threads**, or you can have a hybrid type of (OS/User thread support).

### Advantages of Threads

- Assuming you needed many instances of the same process to run (each with its own address space), a single multi-threaded process can **save on required memory** since all threads share the same memory image. (This is no longer a notable advantage as we will see in virtual memory systems)
- The most useful advantage of threads (within the same process) is that there is **no need for the entire process to block** if one thread of execution blocks (waiting for something to happen) because other threads may be able to continue.
- A context switch from a (thread) context to another thread's context but within the same process (i.e., within the same address space) is called a "**light-weight**" operation, because it involves only saving/reinstating the register values.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

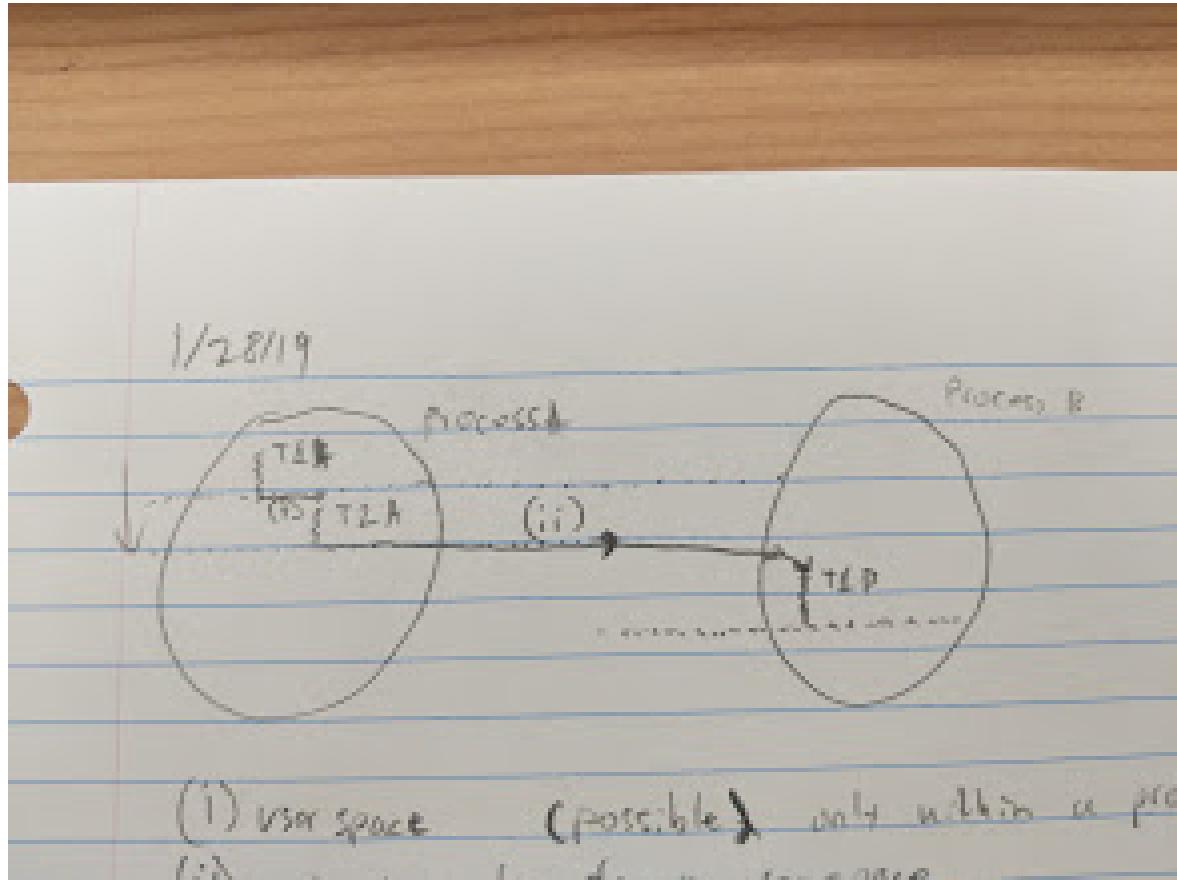
29

Threads can be switched in user space if the threads are within the same process. Note the kernel can also switch thread within the same process if no code is provided for the switch.

Threads in different processes however, will need to call the kernel version of the thread switching instruction.

Generally the kernel operation of thread switching is done in two ways:

1. Time-slice Expiration (if no kernel thread support) OR
2. Kernel is aware of the threads within the process and it automatically.



## Linux

Linux has, essentially, only threads, and hence each "PCB" corresponds to, and describes, a thread.

But Linux threads can share resources among them.

Hence, multiple threads sharing the same address space can be thought conceptually as being threads of a single process.

The PCB in Linux is `task_struct`

(in `include/linux/sched.h`)

The context switch is performed by `__switch_to_asm`

(in `arch/x86/entry/entry_64.s` for x86\_64)

## Digression: the term "Context"

- It means, in general, any state information that has to be maintained and managed carefully in order for correct, meaningful, computation (and not just by processes) to take place (or to save it and continue with it later).
- We will encounter use of the same term for interrupt handler execution (*interrupt context*) or execution within the kernel (*kernel context*). In each case there is a corresponding technicality of what is the scope of the context (boils down to at least the registers and stack involved). It is recommended to mention the type (process, thread, interrupt, kernel) corresponding to the context concept you are trying to convey. Do not let the meaning of context depend on context :-)

# 5: SYNCHRONIZATION

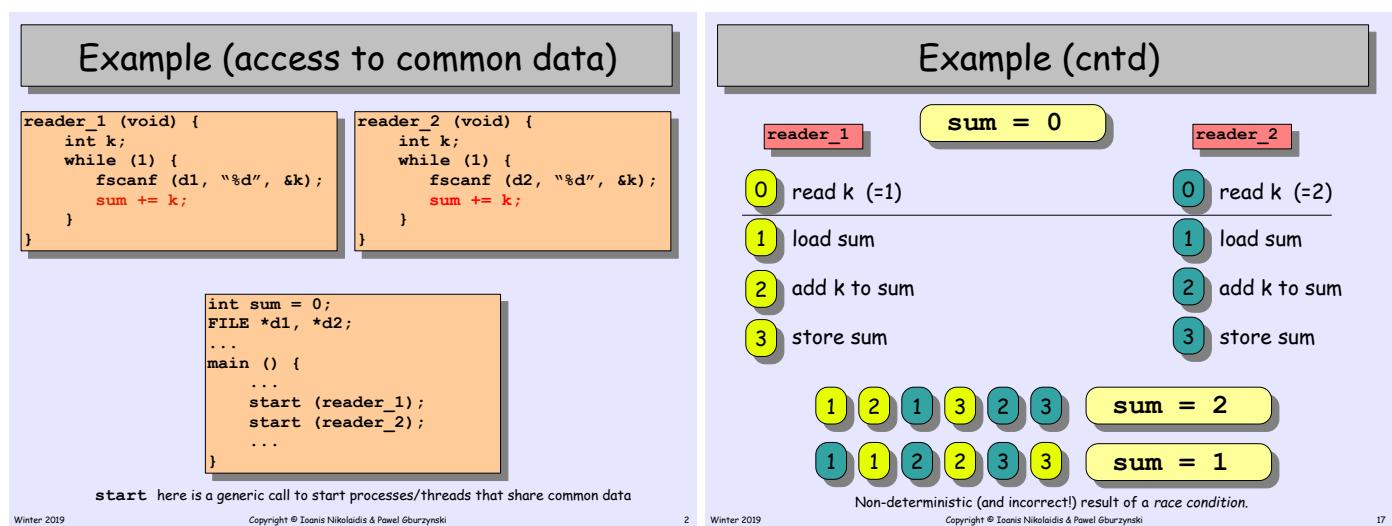
The synchronization of processes and threads are important to maintain the semantics of the program.

Suppose that we have a number of concurrent processes/threads

They could be interleaving where the CPU (single core) switches from one process to another over time. They could overlap where (in a multicore CPU) each core works on a different process.

Generally, unless we employ some mechanism, we have no way of telling in which way the instructions of the two processes will be interwoven in time. Sometimes it makes a difference.

## EXAMPLE OF THREAD CONCURRENCY



The correct behaviour of this program is to add up all the integers between two files concurrently. Bad ordering between the processes will lead to a race condition and incorrect behaviour.

We can use mutual exclusion to solve this issue.

## MUTUAL EXCLUSION

Where one thread of execution never enters its **critical section** at the same time that another concurrent thread of execution enter its own critical section. These critical sections may have operations on global variable where improper order of operations can lead to incorrect behaviour.

The critical sections in the code can be different from each other. The only reason why it's a critical section is because they have access to shared data and the concurrency of that data must be stable or semantically correct.

Note: There could be many independent critical sections so it's best to parameterize the operations.

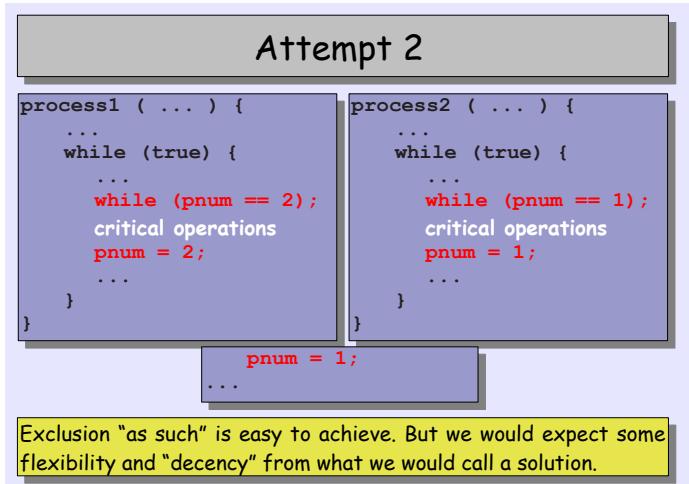
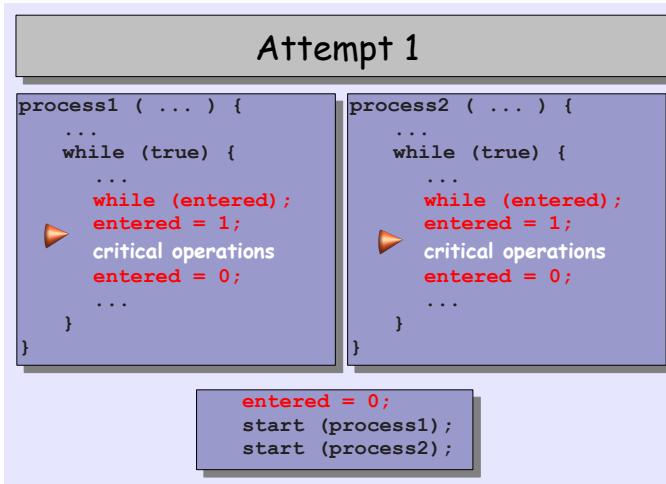
### WHAT IS REQUIRED FOR CONCURRENCY

1. Processes execute concurrently (it makes no difference whether they are interleaving or overlapping). We never know which one will go next and for how long.
2. Consistency of basic operations is guaranteed in hardware. i.e., if two processes try to store something in the same location at the same time, the location is going to end up with one thing or the other thing but not with garbage.

### POSTULATES

- **No "almost working" solutions:** If you think that some interleaving scenarios are not very likely, think again.
- **Simplicity is an asset:** Most critical sections are short, and the operations guarding them should not overshadow their cost.
- **No locked-in patterns of behavior:** Do not assume a particular pattern, e.g., that B will always follow A or anything like that
- **Determinism:** The amount of waiting time should not depend on circumstances other than for how long the section is going to be actually busy (no indefinite postponement).
- **Generalizability:** N processes should not be much more difficult to handle than 2.

## EXAMPLE OF MUTUAL EXCLUSION ATTEMPTS



Attempt 1 is incorrect as the context switch can happen in between the while instruction and the entered = 1 instruction. Thus leading to the same problem as before. Attempt 2 is correct, but now we are setting the order and its basically the same as a sequential instruction. If one of the process takes too long then the other process will have to wait until that process is finished, wasting time and resources.

## Attempt 3

```

process1 ( ... ) {
    ...
    while (true) {
        ...
        while (p2inside);
        plinside = true;
        critical operations
        plinside = false;
    }
}

process2 ( ... ) {
    ...
    while (true) {
        ...
        while (plinside);
        p2inside = true;
        critical operations
        p2inside = false;
    }
}

plinside = p2inside = false;
...

```

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Attempt 4

```

process1 ( ... ) {
    ...
    while (true) {
        ...
        pltrying = true;
        while (p2trying);
        critical operations
        pltrying = false;
    }
}

process2 ( ... ) {
    ...
    while (true) {
        ...
        p2trying = true;
        while (pltrying);
        critical operations
        p2trying = false;
    }
}

pltrying = p2trying = false;
...

```

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

35

**Attempt 3:** Same problem as attempt number 1, you can break this mutual exclusion.

**Attempt 4:** Deadlock if you context switch in between the intent to enter and the while loop.

## Attempt 5

```

process1 ( ... ) {
    ...
    while (true) {
        ...
        pltrying = true;
        while (p2trying) {
            pltrying = false;
            wait_a_little ();
            pltrying = true;
        }
        indefinite postponement
        critical operations
        pltrying = false;
    }
}

process2 ( ... ) {
    ...
    while (true) {
        ...
        p2trying = true;
        while (pltrying) {
            p2trying = false;
            wait_a_little ();
            p2trying = true;
        }
        critical operations
        p2trying = false;
    }
}

pltrying = p2trying = false;
...

```

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Attempt 6: Dekker's Algorithm

```

process1 ( ... ) {
    ...
    while (true) {
        ...
        pltrying = true;
        while (p2trying) {
            if (turn == 2) {
                pltrying = false;
                while (turn == 2);
                pltrying = true;
            }
        }
        critical operations
        turn = 2;
        pltrying = false;
    }
}

process2 ( ... ) {
    ...
    while (true) {
        ...
        p2trying = true;
        while (pltrying) {
            if (turn == 1) {
                p2trying = false;
                while (turn == 1);
                p2trying = true;
            }
        }
        critical operations
        turn = 1;
        p2trying = false;
    }
}

pltrying = p2trying = false;
turn = 1;

```

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

38

**Attempt 5:** Indefinite postponement, meaning each process can be waiting too long for the other to do something.

**Attempt 6:** Something that works finally, this type of mutual exclusion prevents the other process of entering the critical section based on its intent to enter the critical section.

## Attempt 7: Peterson's Algorithm

```

process1 ( ... ) {
    ...
    while (true) {
        ...
        pltrying = true;
        cookie = 2;
        while (p2trying &&
               cookie == 2);
        critical operations
        pltrying = false;
        ...
    }
    pltrying = p2trying = false;
    cookie = ...;
    ...
}

```

## Attempt 8: Hyman's Algorithm

```

process1 ( ... ) {
    ...
    while (true) {
        ...
        pltrying = true;
        cookie = 1;
        while (pltrying &&
               cookie == 1);
        critical operations
        p2trying = false;
        ...
    }
    pltrying = p2trying = false;
    turn = 1;
}

process2 ( ... ) {
    ...
    while (true) {
        ...
        p2trying = true;
        while (turn != 1) {
            while (p2trying);
            turn = 1;
        }
        critical operations
        p1trying = false;
        ...
    }
    pltrying = p2trying = false;
    turn = 2;
}

```

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

39

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

47

**Attempt 7:** This does the same thing as the previous example but puts preferences to a specific process.

**Attempt 8:** A more elaborate problem from attempt 1 and 3. If the context switch happened in the while loop after the second while loop, while the \$PC of the other is before the true pltrying stage.

## A BETTER WAY

The best algorithms are relatively simple, but they only work for two processes. Generalizations are possible but:

- N must be known prior to the execution
- the algorithm becomes awfully complicated

Real-life solutions are based on stronger prerequisites than variable sharing:

- the OS can provide requisite tools for processes (available as system calls)
- some hardware tools can be used in desperate situations

What we want is a function call that will prevent preemption while in the critical section but also parameterize the critical section so that multiple critical sections can exist within a program. We also don't want busy waiting, so we want the waiting process to give way to other processes.

# TAXONOMY OF INTERRUPTS

- **Internal Interrupts:** Sometimes called traps or exceptions. They are caused by events occurring **within** the CPU. Things like:
  - syscall
  - division by zero
  - illegal memory reference
- **External Interrupts:** They are caused by event occurring **outside** the CPU. Things like:
  - I/O notifications
  - timers going off
  - reset button

## Critical sections revisited



Critical sections do not only apply to processes. The kernel has to deal with mutual exclusion problems involving interrupts.

Note that:

- I/O interrupts occur spontaneously
- naive solutions for critical sections won't work (imagine "busy waiting" in an interrupt)

Consequently, some hardware mechanism is needed.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

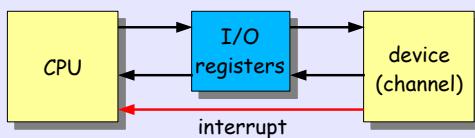
56

## External interrupts ...

... are **asynchronous** (w.r.t. the CPU activities). What does it mean?

Remember that I/O devices are in fact independent processors capable of sustaining non-trivial activities in parallel with the CPU.

Look at it this way: external interrupts provide handles whereby external devices can execute CPU code.



In principle, the I/O registers would suffice, but then the CPU would have to **poll** them for status change.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

55

## External interrupts are maskable

This means that the CPU can request to **hold** an external interrupt until it explicitly resumes its reception.

An interrupt occurring while it is masked remains pending and will be accepted when unmasked.

Note that some internal interrupts may also be maskable, e.g., **integer overflow, floating point underflow** (on some machines).

In contrast to external interrupts, a masked internal interrupt is irretrievably lost if it occurs. The philosophy of masking is different in the two cases:

- |                     |    |                                 |
|---------------------|----|---------------------------------|
| external interrupts | => | hold on for a moment            |
| internal interrupts | => | I don't care about these events |

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

57

Each interrupt have a priority level. Meaning the interrupting interrupt must have a priority higher than the current interrupt in order for the interrupt to take place. Otherwise, current interrupt is executed.

## INTEL X86 INTERRUPTS

All of Intel's interrupts are maskable except for one just in case. Also unlike the MC 68000, the default behavior of Intel's interrupt handling is that only one interrupt is handled at a time. There is no "priority" of interruption.

## Another view of mutual exclusion

```
process1 ( ... ) {
    ...
    enter_section (A) ;
    critical operations
    exit_section (A) ;
    ...
}
```

```
process2 ( ... ) {
    ...
    enter_section (A) ;
    critical operations
    exit_section (A) ;
    ...
}
```

**enter\_section (A)**

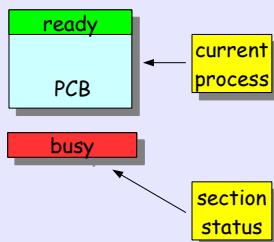
```
some_process ( ... ) {
    ...
    wait_for_event (A) ;
    ...
}
```

```
interrupt_hdlr ( ... ) {
    ...
    trigger_event (A) ;
    ...
}
```

This looks like a wait/wakeup mechanism, i.e., event passing.

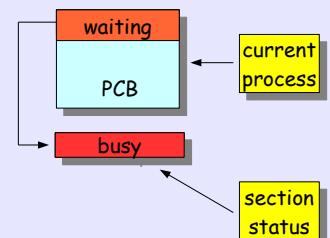
### How most systems would do it

- enter system call
- do not preempt
- check section status
- set section busy
- preempting OK
- return from syscall



### How most systems would do it

- enter system call
- do not preempt
- check section status
- mark as waiting
- preempting OK
- reschedule



Implementing "do not preempt" is easy: the kernel may simply make sure that no rescheduling will occur for a while. To make their lives easier, many systems do not allow rescheduling to occur at all while a process is in the Kernaland.

As you can see in this example the current PCB is now waiting as someone else is currently accessing shared data. The same idea works for devices.

## INTERRUPT MASKING

### Interrupt masking ...

... is the most powerful built-in synchronization mechanism available to the kernel.

The CPU is executing some code (activity), be it a piece of user program or something in the Kernaland, including an interrupt service routine.

The only way for this activity to lose the CPU, that is beyond its control, is when an (external) interrupt occurs.

Remember that internal interrupts do not occur unless the activity generates them (purposely, we hope).

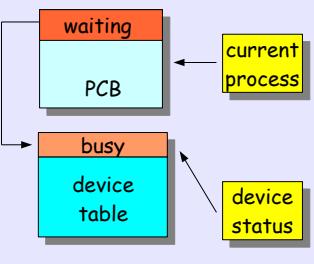
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

87

### Once again

- enter system call
- mask (some) interrupts
- check data available
- wait for data
- unmask interrupts
- reschedule



Nothing wrong can happen now. The interrupt, should it occur here, will remain pending until it is safe for it to be accepted.

Masking interrupts are important as some data may show up while the process is checking for data. The device status may change by an interrupt service routine so the process is waiting even though the device is ready to transmit but can't as the process is not checking for available data.

94

# 6: SYNCHRONIZATION ABSTRACTION

## SEMAPHORES

### COUNTING SEMAPHORES

---

```
wait(S):
    if(S > 0) S--;
    else wait until allowed to proceed;

signal(S):
    if(processes wait on S)
        allow one to proceed;
    else S++;

{ in main code }
repeat
...
wait(S);
critical section
signal(S);
...
until false;
```

#### WHAT DOES P AND V MEAN

P means "to test"

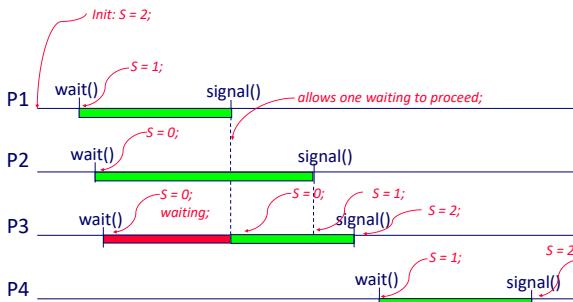
V means "to increment"

Where: S is the number of process(es) that can access the critical section. Note this looks very similar to the waiting algorithm in the previous chapter. If S is initialized to 1, then it is a binary semaphore.

This works well if you programmed it correctly, if programmed incorrectly like flipping the order of wait and signal then the semantics of the algorithm is lost.

## EXAMPLE OF A COUNTING SEMAPHORE

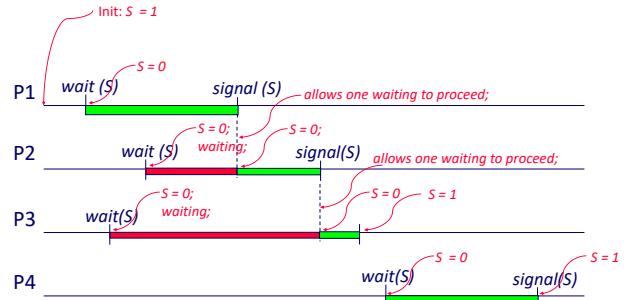
### A counting semaphore example



CMPUT 379, Winter 2019

Synchronization Abstractions 3

### A binary semaphore example



Note: we did not define (so far) any order in which multiple processes waiting on a semaphore are to resume execution.

CMPUT 379, Winter 2019

Synchronization Abstractions 4

## ATOMICITY & SEMAPHORE IMPLEMENTATION

Here is a basic implementation of a semaphore.

Note the two functions must be atomic, they cannot be blocked. Moreover, S can only be modified by one process at a time.

```

1  wait(S){
2      while (S <= 0); // busy wait
3          S--;
4      }
5
6  signal(S){
7      S++;
8  }
```

Because the standard implementation is inefficient as it can lead to busy waiting. The implementation below is better.

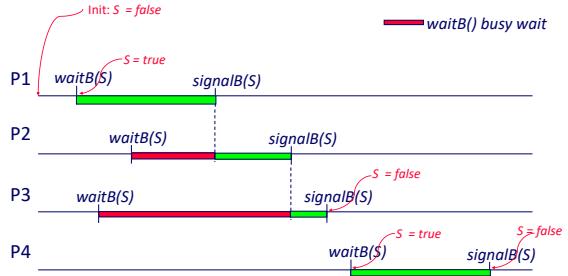
```

1  typedef struct {
2      int value;
3      struct process* list;
4  } semaphore;
5
6  wait(semaphore* S){
7      S->value--;
8      if(S->value < 0){
9          S->list.push_back(this->process)
10         block();
11     }
12 }
13
14 signal(semaphore* S){
15     S->value++;
16     if(S->value <= 0{
17         process = S->list.pop_front();
18         wakeup(process);
19     }
20 }
```

## BINARY SEMAPHORES

```
{ Test and Set version of
a binary semaphore }
waitB(S):
    while Test_and_Set(S) do no-op;
signalB(S):
    S = false;
```

### A binary semaphore example (T&S)



Note: with the busy wait implementation we have also given up any hope we will control the order of resumption.

CMPUT 379, Winter 2019

Synchronization Abstractions 7

This is a test and set semaphore that sacrifices efficiency for simplicity of implementation. This code is similar to waiting until time slice expiration.

#### NOTE

This is a valid solution for multi-core computers where multiple cores are accessing the same critical section at the same time.

## IMPLEMENTATION OF BINARY & COUNTING SEMAPHORES

### Binary semaphores (T&S implementation)

```
S.guard = false;
S.flag = true;
S.queue = EMPTY;

waitB(S):
    while Test-and-Set(S.guard) do no-op;
    if (S.flag) do {
        S.flag = false;
        S.guard = false;
    } elsedo {
        enqueue(S.queue);
        ... and S.guard = false; ...
    }

Taking matters under
control, we will build a
binary semaphore with T&S
but without busy waiting
(well almost).

signalB(S):
    while Test-and-Set(S.guard) do no-op;
    if (waiting(S.queue)) do
        makeready(S.queue);
    elsedo
        S.flag := true;
        S.guard = false;
```

### Implementing counting semaphores

(Using Binary Semaphores)

```
wait(S):
    waitB(S.S1);
    S.C := S.C - 1;
    if S.C < 0 then
        begin
            signalB(S.S1);
            waitB(S.S2);
        end
    signalB(S.S1);

signal(S):
    waitB(S.S1);
    S.C := S.C + 1;
    if S.C ≤ 0 then signalB(S.S2);
    else signalB(S.S1);
```

CMPUT 379, Winter 2019

Synchronization Abstractions 8

CMPUT 379, Winter 2019

Synchronization Abstractions 9

Adding some extra code can help with the busy waiting.

You can use two binary semaphores to implement a counting semaphore. One semaphore is used to protect the counter (value) and the other is a place to wait.

## PRODUCER-CONSUMER SEMAPHORE

---

During initialization let:

```
int n;  
semaphore mutex = 1;  
semaphore full = 0;  
semaphore empty = n;
```

Where n is the size of the pool of buffers.

## Producer-consumer (bounded-buffer) revisited

**repeat**

...

produce item in *nextp*

...

*wait(empty);*  
*wait(mutex);*

...

add *nextp* to buffer

...

*signal(mutex);*  
*signal(full);*

**until** false

**repeat**

*wait(full);*

*wait(mutex);*

...

remove item from *buffer* to *nextc*

...

*signal(mutex);*  
*signal(empty);*

...

consume item in *nextc*

...

**until** false

This is a circular buffer semaphore where "producers" add items into a finite buffer and "consumers" take items from the an array of buffer.

The mutex semaphore is used to prevent access to the pool of buffers when the pool is being modified.  
The empty and full semaphores are interconnected with each other.

The producer uses the empty semaphore to check if the pool is not full of filled buffers. If not then the producer tries to get access to the buffer and modifies it, decrementing the empty counter and signaling the full semaphore to increment their counter by one as the one of the buffers became full.

The consumer uses the full semaphore to check if the pool is not full of empty buffers. If not then the consumer tries to get access to the buffer and modifies it, decrementing the full counter and signaling the empty semaphore to increment their counter by one as one of the buffers became empty.

## READERS AND WRITER PROBLEM

---

```
1 // During initialization let:  
2 int read_count = n;  
3 semaphore rw_mutex = 1;  
4 semaphore mutex = 1;  
5  
6 // writer process  
7 do{  
8     wait(rw_mutex);  
9     /* writing is performed */  
10    signal(rw_mutex);  
11 }while (true);  
12  
13 // reader process  
14 do {  
15     wait(mutex);  
16     read_count++;  
17     if (read_count == 1)  
18         wait(rw_mutex);  
19     signal(mutex);  
20     /* reading is performed */  
21     wait(mutex);  
22     read_count--;  
23     if(read_count == 0)  
24         signal(rw_mutex);  
25     signal(mutex);  
26 } while(true);
```

This is the readers and writers problem where writers cannot update the data when there are readers reading the data. In this is a unique case as the preference of mutual exclusion is given to the reader as you can have multiple readers looking at the data without the worry of race conditions. Writers however, must wait until no process is reading the data in order for it to update.

## LIMITATION OF CLASSIC SEMAPHORES

---

Lack of enforcement schemes.

User processes can call wait and/or signal in arbitrary order:

- mutual exclusion is possibly violated, and/or,
- deadlock may occur

Processes must be willing to block when invoking wait().

The resumption strategy of signal is set and not controlled by the user.

- No additional control arguments.
- Process selection in signal should prohibit starvation,
- but fairness (proportional to "speed" of a process) is non-trivial.

# MONITORS

Monitors encapsulate the shared data into an object. Processes can access the object using procedure entry points or through specific synchronization calls (related to condition variables).

Conditional Variables are:

- Are not an ordinary variable. Similar to a structure.
- you have multiple variables
- signify "things" you can wait on
- They are not conditions (conditional expressions)
- Really they are just user-defined synchronization schemes

Many condition variables can be defined in a monitor.

The operation `x.wait()` means that the process invoking this operation is suspended until another process invokes `x.signal()`. `x.signal()` resume exactly one suspended process. If no process is suspended, then the `signal()` operation has no effect, i.e., the state of `x` didn't change.

## MONITOR SEMANTICS

---

If a process executes `x.signal()`, should it be allowed to execute beyond the `x.signal()` if a process was already waiting in `x.wait()` and was resumed as a result of `x.signal()`?

1. The process that issued the `x.signal()` waits until the process that was waiting in `x.wait()` either exits the monitor or waits on another condition. (Hoare)
2. The process that was waiting on the `x.wait()` waits until the process that issued the `x.signal()` leaves the monitor or waits on a condition.
3. The process that was waiting on the `x.wait()` resumes immediately but the one that called `x.signal()` must leave the monitor (Concurrent C)

## MONITOR IMPLEMENTATION OF A BINARY SEMAPHORE

---

You can use monitors to implement a binary semaphore.

```
var busy: Boolean, x: condition;

procedure wait():
    if busy then x.wait;
    busy = true;

procedure signal():
    busy = false;
    x.signal;

init:
    busy = false;
```

One key difference between this implementation and the "classic" implementation is that the signal becomes a no-op if there isn't anyone waiting.

## SEMAPHORE IMPLEMENTATION OF MONITORS (HOARE)

The complication is that a process that signals on a condition variable must, itself get suspended from executing (hence has to wait) because the signaled process will now resume execution within the monitor. This is a different "kind" of waiting than the waiting to (originally) enter the monitor when a monitor method was called. We introduce one new semaphore (next) to deal with this.

1. You have processes waiting to enter the monitor (wait on mutex)
2. You have processes waiting on the "next" semaphore
3. You have processes waiting on the condition variable (x.sem)

### Semaphore

mutex      serialize access to operations in monitor ( = 1)  
nextsem    suspend processes after executing SIGNAL ( = 0)  
x.sem      suspend a process executing WAIT on condition  $x( = 0)$

### int

next      number of processes waiting due to signal ( = 0)  
x.count    number of processes waiting due to waiting on  $x_i$

```
1 function_j(){  
2     wait(mutex);  
3     <operation body>  
4     if (next > 0)  
5         signal(nextsem);  
6     else  
7         signal(mutex);  
8 }  
  
1 x.wait(){  
2     x.count++;  
3     if (next > 0)  
4         signal(nextsem);  
5     else  
6         signal(mutex);  
7     // process waits here  
8     wait(x.sem);  
9     x.count--;  
10 }  
  
1 x.signal(){  
2     if (x.count > 0){  
3         next++;  
4         signal(x.sem);  
5         wait(nextsem);  
6     }  
7 }  
8 }
```

## JAVA EXAMPLE

Objects provide mechanism that could be used to implement monitors. A monitor class needs to use only private data fields, and use synchronized methods.

Each object possesses a lock and a wait queue. The lock is accessed by the synchronized (obj) ... block. In this form the argument to synchronized defines the object whose lock we are interested to access. A method can be defined with `asynchronized` attribute. Really a shorthand for synchronized (this) ... of the outermost block (containing the entire code of the function).

Each object contains a single anonymous condition-like structure. The condition is accessed via calls to `wait()` / `notify()`. `wait()` causes the current thread to wait until another thread invokes the `notify()` method or the `notifyAll()` method for this object.

`wait()`, `notify()` or `notifyAll()` can be called on an object only by a thread that currently holds the lock to the object.

`wait()` the thread calling it is placed in the wait queue of the object. The lock to the object is released.

`notify()` if a thread is waiting at the object's wait queue, it is made ready. The thread waits to re-obtain the lock on the object. (Could be blocked from acquiring it because the thread that executed the `notify()` has not yet released it, or because another thread got it in the meantime.) After the lock is obtained, the thread resumes at the point following the `wait()`. The logical condition that was true when the `notify()` was invoked is not necessarily true when the thread that was waiting resumes running.

`notifyAll()` is like `notify()` but all threads in the wait queue of the object are made ready. Still, each of the ready threads needs to re-acquire the lock before continuing after the `wait()`.

## SEMAPHORE

```
private class Semaphore {  
    private int S;  
    public Semaphore(int initial)  
    {  
        S = initial  
    }  
    public synchronized void waitSema()  
    {  
        S--;  
        while(S < 0) wait();  
    }  
    public synchronized void signalSema()  
    {  
        S++;  
        if(S <= 0) notify();  
    }  
}
```

## BOUNDED BUFFER AND PRODUCER-CONSUMER

```
Class BBuffer  
{  
    final int size = 10;  
    int[] thebuffer = new int[size];  
    int indexPut = 0, indexGet = 0;  
    int count = 0;  
  
    public synchronized void put(int item)  
    {  
        while (count == size-1) wait();  
        buffer[indexPut] = item;  
        indexPut = (indexPut+1)%size;  
        count++;  
        if (count == 1) notifyAll();  
    }  
  
    public synchronized int get()  
    {  
        int item;  
        while (count == 0) wait();  
        item = buffer[indexGet];  
        indexGet = (indexGet+1)%size;  
        count--;  
        if (count == size-1) notifyAll();  
        return item;  
    }  
}
```

## DEADLOCK

```

Class Swapper {
    private int value;
    public synchronized int getValue()
    {
        return value;
    }
    public synchronized void setValue (int i)
    {
        value = i;
    }
    public synchronized void swapValue
    (Swapper s)
    {
        int temp = value;
        value = s.getValue();
        s.setValue(temp);
    }
}

```

## LIMITATION OF MONITORS

- Unlike semaphores, signal()s in monitors do not accumulate.
- Condition variables are not linked directly to logical (Boolean) expressions.
- Evaluation of the logical condition (if any) is left up to the programmer.
- A monitor's state (that led to a particular "condition") may change (depending on the monitor semantics) either
  - between signal() and resumption of a wait()-ing thread or
  - between signal() and resumption of the thread that signal()-ed.
- Users are free to call into a monitor in any (not always "correct") order.
- Nesting of calls into monitors needs to be addressed.
- Visibility of the monitor internal variables needs to be addressed

## ATOMIC TRANSACTION

### Atomic transactions

*Gain access to resources...*  
...  
 $account[i] := account[i] - sum;$   
...  
**CRASH!!!!**  
...  
 $account[j] := account[j] + sum;$   
...  
*Release allocated resources...*

CMPUT 379, Winter 2019Synchronization Abstractions 27

Change should not be committed before the end of the transactions. An "all or nothing" property.  
In the event of a crash, recovery is possible:

- log-based recovery
- checkpoints.

Concurrency introduces new challenges.

## SERIALIZABILITY

### SERIALIZABILITY

The outcome of the concurrent execution of atomic transactions should be equivalent to the sequential execution of the same transactions in an arbitrary order.

### Serializability example

<i>(serializable execution)</i>		<i>(an equivalent sequential execution)</i>	
T0	T1	T0	T1
-----	-----	-----	-----
read(A)			read(A)
write(A)			write(A)
	-----	-----	-----
	read(A)	read(A)	read(B)
	write(A)	write(B)	write(B)
-----	-----	-----	-----
read(B)			read(A)
write(B)			write(A)
	-----	-----	-----
	read(B)	read(B)	read(B)
	write(B)	write(B)	write(B)

CMPUT 379, Winter 2019Synchronization Abstractions 30

## CONFLICTING OPERATIONS

- Write/write (overwriting uncommitted data)

Example: transaction T0 overwrites value A with A'. Transaction T1 overwrites value B with B'. Then T0 overwrites value B' with B". Finally, T1 overwrites A' with A". The result is non-serializable. (Note all updates are without previous reads, i.e., "blind")

- Write/read (reading modified, but uncommitted, data)

Example: Transaction T0 reads A and writes A'. Transaction T1 reads A' and writes A". Then transaction T0 does something else, but does not commit (aborts). The state is rolled back to A. Now it is as if T1 has read a value that never existed!

- read/write (unrepeatable reads)

Transaction T0 reads value A. Transaction T1 updates A to A' and commits it. T0 then re-reads the value and gets A'.

## GENERAL PROBLEM OF SYNCHRONIZATION

Our first priority is correctness, i.e., concurrent operations should leave the system in a consistent, i.e. "correct", state, but there are more things we need to address:

Livelock, or starvation: the situation in which some processes are making progress toward completion but some others are locked out of the resource(s). (could manifest as unfairness)

Deadlock: the situation in which two or more processes are locked out of the resource(s) that are held by each other. Progress is then impossible among the deadlocked processes.

## ATTAINING CONCURRENCY AND SERIALIZABILITY

Associate each data item with a lock.

- must possess lock to access the data item
- Lock Modes:
  - Shared: allows reading only, others can also lock as shared (to prevent read/write conflicts)
  - Exclusive: allows reading and writing, others cannot hold the lock (to prevent write/read and write/write conflicts)
- Locks are not enough, you must know when to release them. The simplistic approach: obtain all at the beginning release all at end.
- A better approach: Two phase locking: First phase: growing phase (may obtain but not release), Second phase: shrinking phase (may release but not obtain).

# 7: DEADLOCKS

## DEADLOCK

A process is said to be in a the state of deadlock when it is waiting for something that will probably never happen.

## RESOURCE ALLOCATION

Traditionally, the deadlock problem is discussed in terms of resource allocation: a process holds some resources and needs more to complete. A serious problem occurs if there is no sensible (non-destructive) way to remove the resources from the process. By stretching our idea of a resource, we can define all (non-trivial) kinds of deadlock as resource allocation problems. Remember: critical sections are resources.

## HOW TO AVOID DEADLOCK?

There are many ways to avoid a deadlock

1. provide enough resources for everybody, such that they will never get into problems
2. teach the "processes" to be courteous and yield
3. teach them how to fight if other parties are not courteous
4. impose a protocol (restrictions) on the resource allocation scheme. Like a traffic signal.

## RESOURCE TYPES

### TYPES OF RESOURCES

**Preemptible:** can be taken from the process and later returned without affecting the process's integrity. For example, CPUs are generally preemptible.

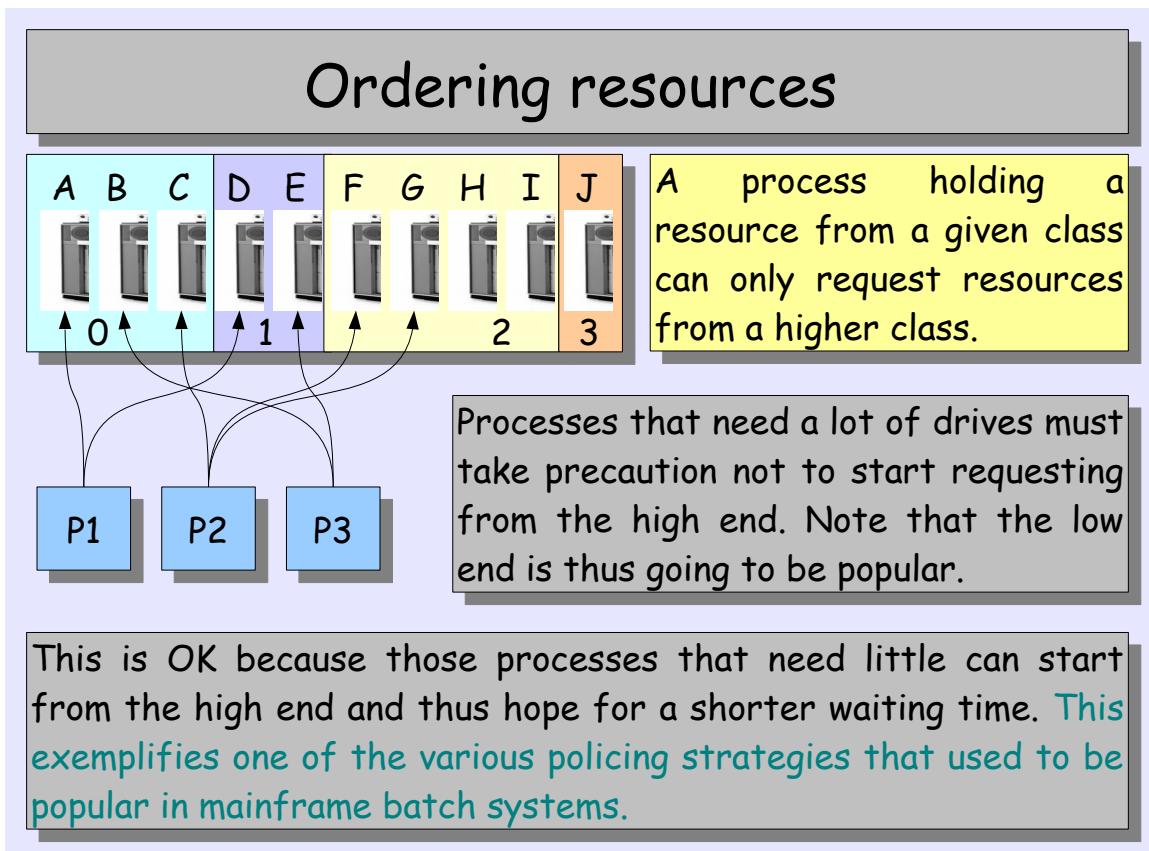
**Non-preemptible:** once a process acquires the resource, it must remain allocated to the process until explicitly released. For example, Tape drives, USB storage keys, are non-preemptible resources.

# NECESSARY CONDITIONS FOR A DEADLOCK

## DEADLOCK CONDITIONS AND THEIR DENIALS

- **Mutual Exclusion:** Resources are not shared. Processes holding them require mutual exclusion (exclusive access).
- **Hold + wait:** Processes are allowed to hold resources while waiting for additional resources. **Prevention** Processes cannot request resources incrementally. You must receive everything you are ever going to need at the beginning. If your request cannot be fulfilled immediately, you have to release everything, and re-request everything back together with the new item.
- **No Preemption:** Resources cannot be removed from processes and later returned unless the process releases them explicitly. **Prevention** If your request cannot be fulfilled immediately, you have to release everything, and re-request everything back together with the new item.
- **Circle:** One can identify a circular chain of processes and resources, with each process holding a resource needed by the previous process and requesting a resource held by the next process. **Prevention** You must request resources in a certain order. This order is followed by everybody.

# ORDERING RESOURCES



You can have different ordering depending on the resource requirement of the process. Higher requirements of resources can get their requirement from the low end while lower requirement of resources can get their requirement from the high end.

# BANKER'S ALGORITHM

Suppose that there is only one resource type with  $M$  units. These are the rules:

## RULES

1. Each process (job) declares in advance the maximum number of units that it will be needing at any given time. This declaration cannot be more than  $M$ .
2. The system aborts a process that violates its declaration.
3. A process may request an arbitrary number of units at once, as long as the request does not violate its declaration. The process may have to wait, but the system guarantees that as long as processes eventually complete (and thus release what they've been holding), all processes will happily execute to successful completion.

## SAFE VS UNSAFE STATE

### SAFE VS UNSAFE STATE

A system is in a safe state if and only if the system can fulfill the maximum possible amount of resources of the "smallest" process ready to be executed.

If it cannot then the state is unsafe and a deadlock can occur.

## GENERAL CASE AND ALGORITHM

$n$  the # of processes

$m$  the # of resource Types

$av[j]$  the # of resource units of type  $j$  still available

$Max[i,j]$  the max demand of process  $i$  for resource  $j$

$Al[i,j]$  current allocation of resource  $j$  to process  $i$

$N[i,j]$  remaining need for process  $i$  for resource  $j$

$N[i,j] = Max[i,j] - Al[i,j]$

$Av[j] = T[j] - \sum_i Al[i,j]$

Where  $T[j]$  is the total number of units of resource  $j$  available in the idle system.

```
1 void allocate (int j, int k) {
2     if (k > N[current ,j])
3         abort (current);
4     while (1) { // lock
5         if (k <= Av[j]) {
6             Av[j] -= k;
7             Al[current ,j] += k;
8             N[current ,j] -= k;
9             if (safe ())
10                 grant (j, k); // unlock
11                 return;
12             }
13             Av[j] += k;
14             Al[current ,j] -= k;
15         }
```

```
15     N[current ,j] += k;
16     }
17     wait (RESOURCE_RELEASE);
18     schedule (); // unlock
19   }
20 };
```

### SAFETY CHECK

```
1 bool safe () {
2     int W[m], i, j; bool F[n], alldone;
3     for (j=0; j<m; j++) W[j]=Av[j];
4     for (i=0; i<n; i++) F[i]=false;
5     Restart:
6     alldone = true;
7     for (i=0; i<n; i++) {
8         if (!F[i]){
9             for (j=0; j<m; j++)
10                 if (N[i,j] > W[j]) break;
11             if (j == m) {
12                 for (j=0; j < m; j++) W[j] += Al[i,j]
13                 F[i] = true;
14                 goto Restart;
15             }
16             alldone = false;
17         }
18     }
19     return alldone;
20 }
```

## EXAMPLE OF A SAFE STATE

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that **there is a way to allocate whatever is left to the processes as to execute them all to completion**.

Example:	Process	Max	Current	M=12
	P0	10	5	
$4-2 \leq 3$ $3+2 = 5$	P1	4	2	
	P2	9	2	

Left

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

18

Process 1 can be processed as it only requires 2 more "resources" to reach its max. It is consumed and its resources released back into the pool. This repeats for process 0 then 2.

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that **there is a way to allocate whatever is left to the processes as to execute them all to completion**.

Example:	Process	Max	Current	M=12
	P0	10	5	
$10-5 \leq 5$ $5+5 = 10$	P1	4	2	

Left

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

20

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that **there is a way to allocate whatever is left to the processes as to execute them all to completion**.

Example:	Process	Max	Current	M=12
	P0	10	2	
$9-2 \leq 10$ $10+2 = 12$	P1	4	2	

Left

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

22

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that **there is a way to allocate whatever is left to the processes as to execute them all to completion**.

Example:	Process	Max	Current	M=12
	P0	10	5	
	P1	4	2	
	P2	9	2	

Left

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

19

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that **there is a way to allocate whatever is left to the processes as to execute them all to completion**.

Example:	Process	Max	Current	M=12
	P0	10	2	
	P1	4	2	

Left

10

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

21

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that **there is a way to allocate whatever is left to the processes as to execute them all to completion**.

Example:	Process	Max	Current	M=12
	P0	10	2	
	P1	4	2	

Left

12

This configuration is safe.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

23

## EXAMPLE OF AN UNSAFE STATE

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that there is a way to allocate whatever is left to the processes as to execute them all to completion.

Example:

Process	Max	Current	M=12	Left
P0	10	5		
P1	4	2		
P2	9	3		2

How about this one?

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski 24

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that there is a way to allocate whatever is left to the processes as to execute them all to completion.

Example:

Process	Max	Current	M=12	Left
P0	10	5		
P1	4	2		
P2	9	3		4

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski 26

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that there is a way to allocate whatever is left to the processes as to execute them all to completion.

Example:

Process	Max	Current	M=12	Left
P0	10	5		
P1	4	2		
P2	9	3		4

Note that **unsafe state does not imply deadlock!** The process's **Max** declaration may have applied to some moment in the past, and the process may not request that many items any more.

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that there is a way to allocate whatever is left to the processes as to execute them all to completion.

Example:

Process	Max	Current	M=12	Left
P0	10	5		
P1	4	2		
P2	9	3		2

$4 - 2 \leq 2$   
 $2 + 2 = 4$

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski 25

### Safe vs. unsafe state

The system makes sure to keep the resource allocation in a **safe state**. It grants a request only if the resulting state is safe, meaning that there is a way to allocate whatever is left to the processes as to execute them all to completion.

Example:

Process	Max	Current	M=12	Left
P0	10	5		
P1	4	2		
P2	9	3		4

?

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski 27

# DEADLOCK DETECTION

You see a stale configuration of processing, each of them waiting for some resource, and each of them holding some resource. You ask yourself "are they deadlocked or just temporarily stuck on something"?

n	# of processes that appear stuck
m	# of resource types
Av[j]	# of resource units of type j still available
Al[i,j]	current allocation of resource j to process i
Rq[i,j]	# of units of resource j requested by i

## DIFFERENCE BETWEEN THIS AND BANKER'S

While this looks similar to Banker's algorithm there are subtle differences that allows this to simply check if the stuck processes can make any progress at all. This also doesn't require the processes to declare anything in advance. We just want to check if the stuck process can make progress at all.

## DEADLOCK DETECTION ALGORITHM

```
1 bool deadlock () {
2     int W[m] , i , j;  bool F[n] , alldone;
3     for (j=0; j<m; j++) W[j]=Av[j];
4     for (i=0; i<n; i++) F[i]=false;
5     Restart:
6         alldone = true;
7         for (i=0; i<n; i++) {
8             if (!F[i]) {
9                 for (j=0; j<m; j++)
10                     // replace N[i,j] with Rq[i,j]
11                     if (Rq[i,j] > W[j]) break;
12                     if (j == m) {
13                         for (j=0; j < m; j++) W[j] += Al[i,j];
14                         F[i] = true;
15                         goto Restart;
16                     }
17                     alldone = false;
18                 }
19             }
20         return !alldone; // negate this from safe
21 }
```

## IS DEADLOCK DETECTION IMPORTANT?

The popular realistic cases of deadlock result practically exclusively from synchronization, i.e., acquiring multiple locks in the wrong order.

Remember: always acquire multiple locks in the same order in all processes (or threads).

Sometimes you don't know whether the other processes follow your protocol. What then? You can force enforcement I don't know.

## ANOTHER USEFUL OPERATION

```
1 // S->P(); + S->V(); S->T(); Try the semaphore
2   , return immediately (true or false)
3 void getTwoLocks (semaphore *sem1 , semaphore *
4   sem2) {
5     while (1) {
6       sem1->P();
7       if (sem2->T()) return;
8       sem1->V();
9       sem2->P();
10      if (sem1->T()) return;
11      sem2->V();
12    }
13 }
```

# 8: SCHEDULING

In this chapter we will talk about different scheduling algorithms, how they work, their pros and cons.

An analogy can be made where the customers are the processes and the server is the entity that is tasked to handle or process the customer. This can be abstracted to CPU bursts which is handled by the CPU scheduler to be processed by the CPU, etc.

## BASIC PERFORMANCE CRITERIA

**Turnaround time:** How much time it takes on average to handle a customer.

**Throughput:** How many customers can I handle per unit of time?

**Customer satisfaction:** Maximize the number of customers that are happy with my service.

There is a careful balance between turnaround time and throughput.

## SHORT TERM SCHEDULING

**CPU bursts are the customers.** We care about the turnaround time of CPU bursts.

SJF says that processes with short CPU bursts (i.e., I/O bound processes) should be given priority in acquiring the CPU.

Easy to say. How can we know what the duration of a process's next CPU burst is going to be?

## GUESSING THE CPU BURST

### EXPONENTIAL MOVING AVERAGE

A standard way to calculate a snapshot dynamic average of (potentially) infinitely many samples, such that old samples are aged out is through this formula:

$$\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n \quad (1)$$

Where:

$t_k$  is the k-th sample

$\tau_k$  is the k-th average

$\alpha$  is the recency coefficient (0-1] where a higher value means history matters very little and a lower value means history matters a lot.

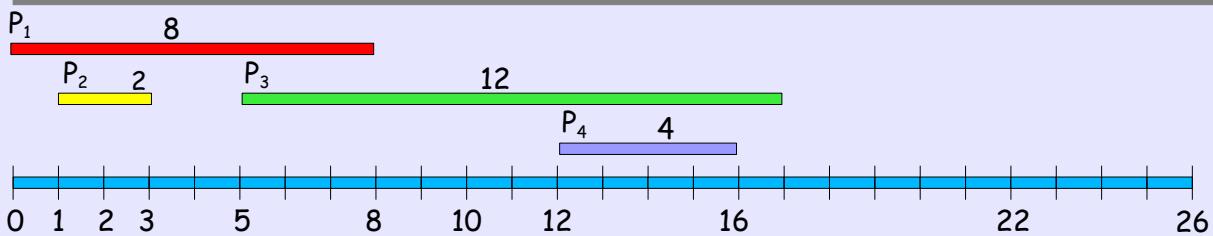
# SJF: SHORTEST JOB FIRST

## SHORTEST JOB FIRST

Give priority to the customers with the shortest processing time. This policy guarantees the shortest possible average turnaround time.

## NON-PREEMPTIVE VS PREEMPTIVE

### Example



Non preemptive SJF: average turnaround =  $(8+9+17+14)/4 = 12$

Preemptive SJF: average turnaround =  $(10+2+21+4)/4 = 9.25$

Thus, preemptive SJF gives a better turnaround time. This isn't surprising, as shorter jobs, not being known in advance, cannot be noticed (and properly accounted for) by the non-preemptive variant when earlier (and longer) jobs are started.

The time taken for each process in both the non preemptive and preemptive SJF algorithm  
 $T(\text{process}) = \text{time it actually finished} - \text{time it was suppose to start}$

## PRIORITY

Generally, scheduling schemes may be based on the notion of priority associated with a process.

### PRIORITY FOR SJF

$$P = \frac{1}{\tau} \quad (2)$$

Where  $\tau$  is the predicted duration of the next CPU burst.

### OTHER FACTORS

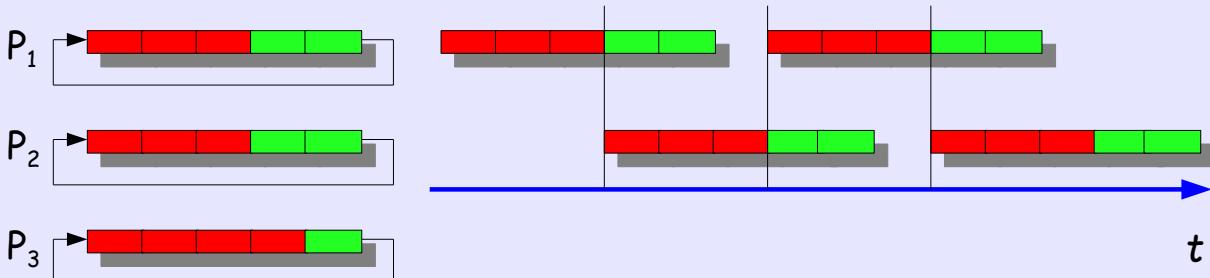
There could be other factors that affect the priority of the process such as internal factors like CPU burst length, memory size, objective importance of the process, etc or external factors who's attributes are assigned by humans.

## STARVATION

The problem with both SJF algorithms is that they are starvation prone. The current iteration of the algorithm does not allow other "low" priority processes to be executed.

### Starvation

A scheduling policy may be starvation prone, i.e., a constant availability of high-priority jobs may result in a starvation of low-priority jobs. For example, SJF (preemptive or not) is starvation prone:



The standard way to eliminate starvation is to add a time component to the priority, e.g., the age factor.

# ROUND ROBIN

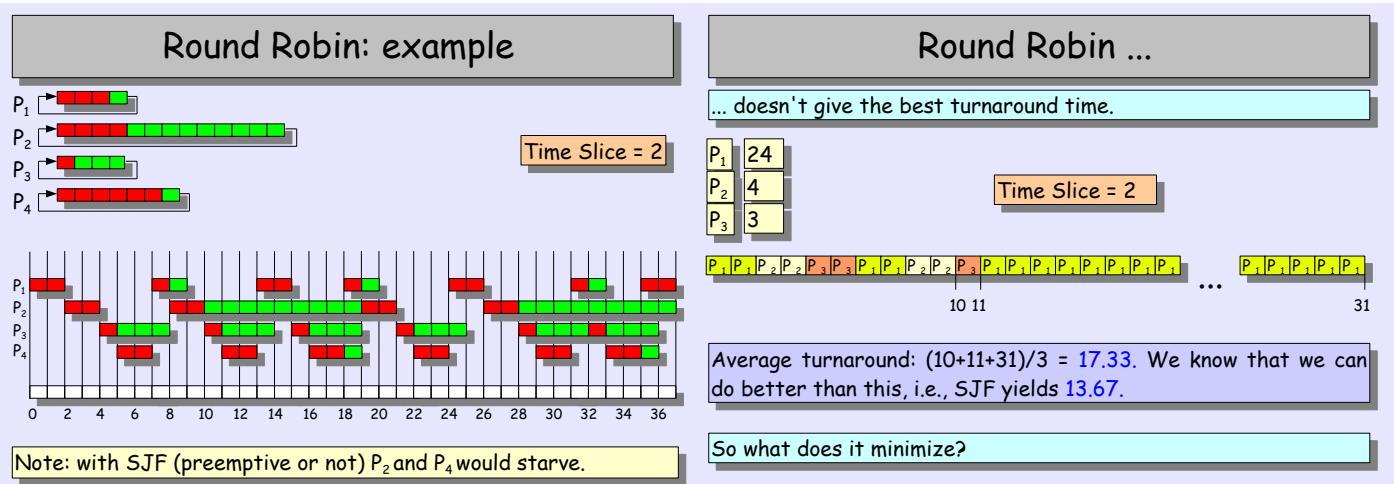
Average turnaround time is not the only performance criterion of interest. For an interactive job (with a human user sitting in front of a terminal), the primary criterion is the (expected) waiting time for the program's reaction.

## TIME SLICE

A process can lose the CPU in two ways:

- It hits an I/O burst, i.e., its CPU burst runs out naturally.
- Its time slice expires.

The time slice can be the same for all; also, some processes may be getting a larger share of the CPU â€¢ reflecting their higher importance



Round robins guarantees that if a process "has something to say," it won't have to wait for more than  $(N - 1) \times S$  where  $N$  is the number of processes and  $S$  is the time slice. Theoretically, the best response time is achieved when  $S$  is the shortest possible, i.e., single instruction time. Practically, the context switch overhead imposes a limit on how small  $S$  can be. Typically,  $S$  is between a millisecond and a tenth a second.

# FPPS: FIXED PRIORITY PREEMPTIVE SCHEDULING

What is the best scheduling policy for high priority, extremely I/O bound processes? Such processes are called "reactive"

## FPPS

Processes are rigidly ordered in some fashion. The first process that is ready gets the CPU

We know for a fact that the CPU bursts are always short. There is no need to guess or estimate. Therefore, fixed (high) priority is OK. Also, the exact ordering of those processes (their actual priority) doesn't really matter, as long as they are ahead of everything else.

This type of scheduling is particularly important for kernel processes.

To perform actions that cannot be properly performed from an interrupt service routine. For example:

- something that does I/O on its own and looks like a thread (paging-related daemons, file-system components, especially for networked file systems)
- something that contains complicated actions of lower priority than interrupts; while those actions are triggered by interrupts, it makes better sense to carry them out elsewhere.

Note that interrupts may block other interrupts, and thus impair then system's responsiveness, if they are "too heavy".

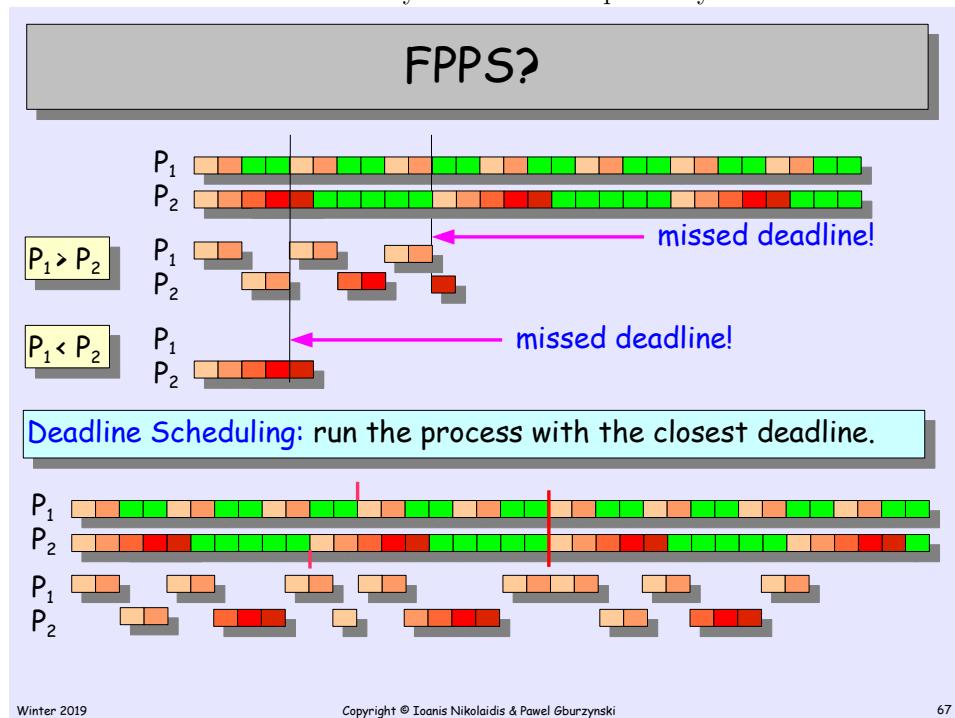
In the latter case, if the action looks like a single burst, it need not be implemented as a single process.

Therefore, during an interrupt service where it would return to User-land it will run the bottom halves of the lower priority actions.

## HARD DEADLINES

---

In some (RT) applications, average performance is irrelevant. What matters is the obedience of hard deadlines. A process may have a certain amount of time to complete its work and then slack off, give the CPU to someone else. However, they must finish their work before its next ready time which is set at a certain time interval. Note different processes could have deadlines where they need to be completed by.



Regardless of priority both processes will miss their deadline for completion.

If P<sub>1</sub> must be completed as soon as it becomes available then P<sub>2</sub> will miss their deadline as the scheduler is giving P<sub>1</sub> too much time. Likewise if P<sub>2</sub> must be completed as soon as it becomes available then P<sub>1</sub> will miss their deadline for the same reason.

To fix this issue the scheduler must complete the process that is closest to its deadline.

In this case instead of P<sub>1</sub> starting its third round the scheduler allows P<sub>2</sub> to finish as its deadline was closer than P<sub>1</sub>.

# MULTI-MODE SYSTEMS/MULTIPLE POLICIES

## HAVING MULTIPLE POLICIES

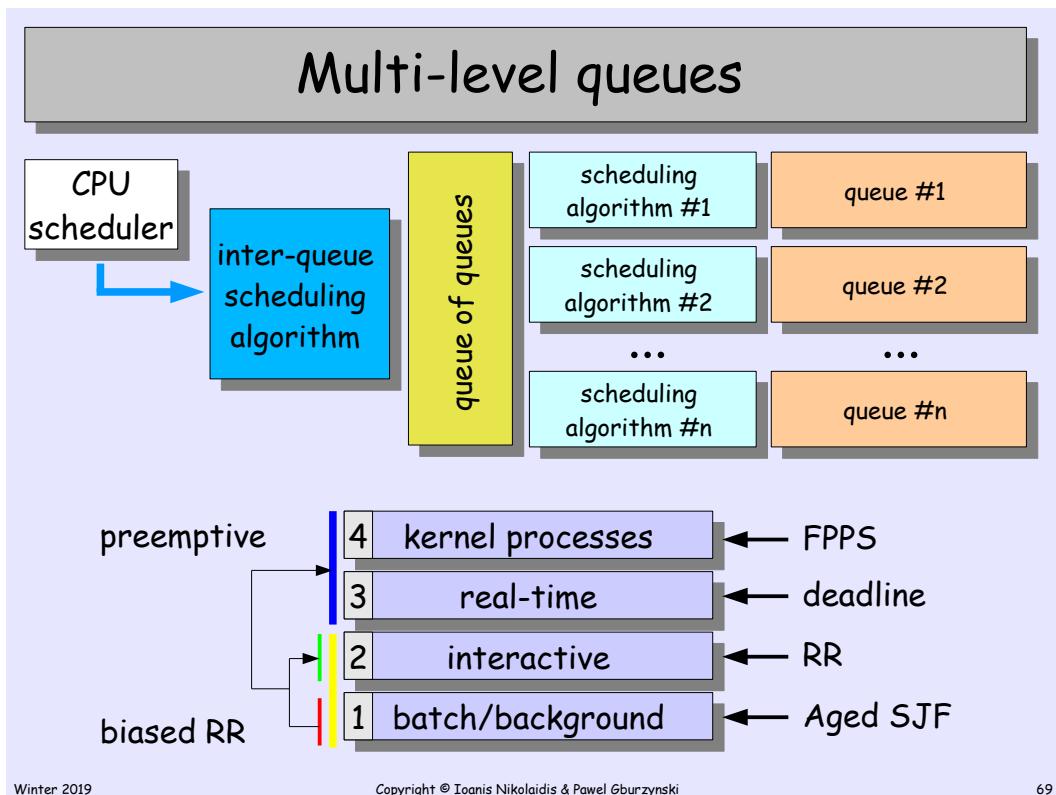
**SJF:** Batch (or background) jobs. Turnaround time is the measure that matters. Aging may be added to avoid starvation.

**Round Robin:** Interactive jobs. Response time is the objective.

**FPPS:** kernel jobs. They should have the highest priority, and their relative importance doesn't really matter.

**Deadline Scheduling:** Real-time jobs with hard constraints. Obedience of deadlines is the objective.

## MULTI-LEVEL QUEUES

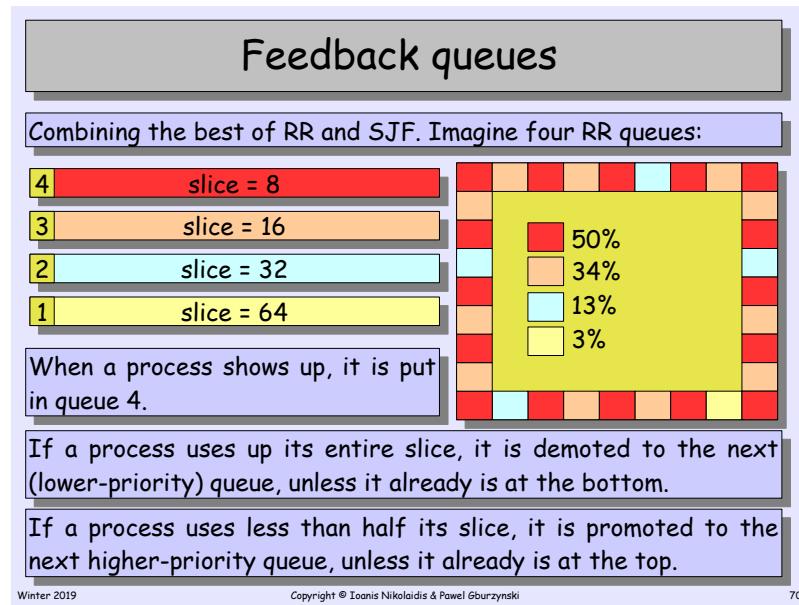


Different tasks based on priority have different scheduling policies.

Kernel processes have the highest priority, they get all the cpu time. Then real-time, followed by interactive, and finally batch/background processes have the lowest priority compared to the rest.

Some batch/background processes can be interactive depending on the workload.

Kernel processes cannot be controlled by the user. They are fixed by the developer. The rest can be controlled at some level.



A single policy that is a combination of two policies we discussed before.

Processes that uses less or equal to the smallest time slice will be called more often achieving the benefits of a SJF policy. But to prevent starvation there are round robin queues that have larger time slices but are called less often for processes that take up more cpu time. This trickle down to larger queues with smaller frequencies.

## REAL-TIME VS STANDARD SYSTEM

### Standard system

Capacity	high throughput, statistically good utilization of all equipment and low overhead
Responsiveness	low average turnaround time, low average response time
Overload	fair performance degradation

### Real-time system

Capacity	"scheduleability," number of context switches per second, number of interrupts per second
Responsiveness	worst-case delay, low variance, deterministic behavior
Overload	stability, unaffected performance for critical tasks

### PRIMARY SOURCES OF PROBLEMS

- **Scheduling policies**, e.g., aging introduces non-determinism, rigid priorities are starvation-prone, deadline scheduling requires extra specification and is not for everybody.
- **Context switch**, e.g., extra actions upon context switch may reduce capacity and increase non-determinism
- **Synchronization tools and paradigms**, e.g., critical sections in the kernel may impair scheduleability (context switches and interrupts).
- **Memory management**, e.g., sophisticated contemporary techniques of memory mapping introduce a lot of non-determinism.

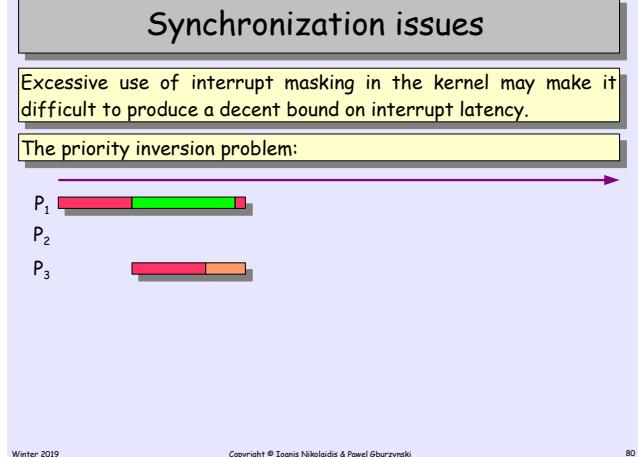
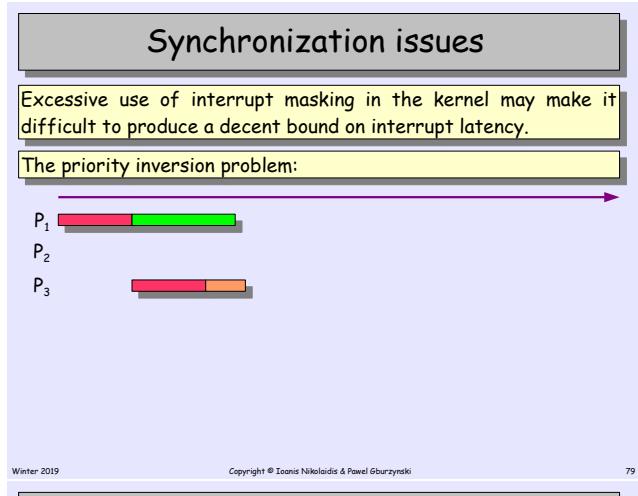
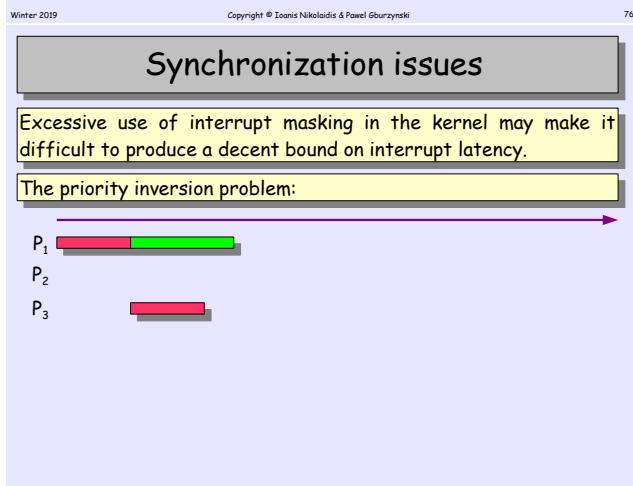
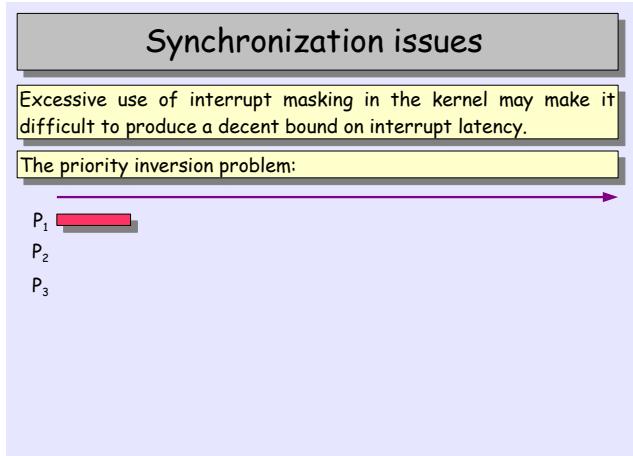
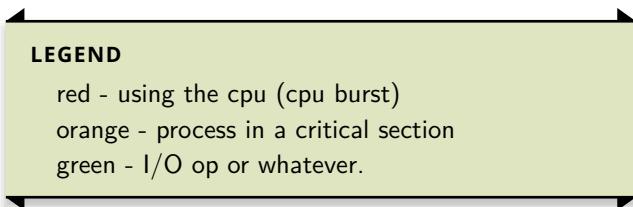
Synchronization tools and paradigms are a problem in real-time systems as it loses determinism. The process that gets the lock are not known until it happens. Same with memory management, it abstracts memory management at the cost of having non-determinism.

## EXAMPLES OF NON REAL-TIME PARADIGMS

- **Non-preemptibility of processes in the Kernel-land** (many standard UNIX systems work this way): Processes cannot be rescheduled until they leave the Kernel-land. Usually, they spend very little time in there, but who can tell the bound?
- **Virtual memory, including lazy loading, also DLLs**: A program running happily at full speed may suddenly have to wait for I/O for no good reason
- **Disk inodes for representing devices (as on UNIX systems)**: Access to an otherwise extremely fast device may stall because its inode has disappeared from the in-memory cache

## SYNC ISSUES

### PRIORITY INVERSION



P1 has priority over P2 and P3 and tries to gain access to the critical section but can't as P3 has the lock.

## Synchronization issues

Excessive use of interrupt masking in the kernel may make it difficult to produce a decent bound on interrupt latency.

The priority inversion problem:



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

81

P<sub>3</sub> continues with P<sub>1</sub> waiting to gain access to the critical section.

## Synchronization issues

Excessive use of interrupt masking in the kernel may make it difficult to produce a decent bound on interrupt latency.

The priority inversion problem:



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

83

P<sub>2</sub> is introduced. Since P<sub>2</sub> has a higher priority than P<sub>3</sub> it blocks P<sub>3</sub> and uses its CPU burst.

## Synchronization issues

Excessive use of interrupt masking in the kernel may make it difficult to produce a decent bound on interrupt latency.

The priority inversion problem:



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

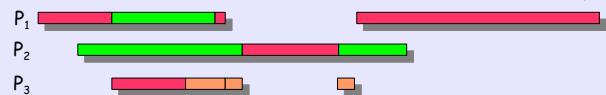
85

P<sub>2</sub> finishes its CPU burst and allows P<sub>3</sub> to execute, but P<sub>3</sub> only needed a little extra time to finish its ownership of the critical section.

## Synchronization issues

Excessive use of interrupt masking in the kernel may make it difficult to produce a decent bound on interrupt latency.

The priority inversion problem:



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

86

P<sub>1</sub> can now get access to the critical section.

## Synchronization issues

Excessive use of interrupt masking in the kernel may make it difficult to produce a decent bound on interrupt latency.

The priority inversion problem:



Within this area, P<sub>1</sub> has been effectively demoted to the priority level of P<sub>3</sub>.

A process holding a critical section should be temporarily promoted to the highest priority of a process waiting for it.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

88

# 9: MEMORY MANAGEMENT

There are different levels of abstractions regarding memory.

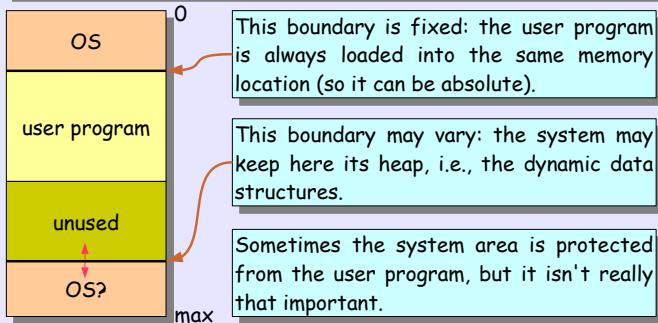
There is the physical, logical, and the virtual. The last two are the process's view of memory.

## MEMORY ALLOCATION (MANAGEMENT)

### Memory allocation (management)

(a historic perspective)

In a single-programming system, the issue is rather simple, e.g., we may have something like this:



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

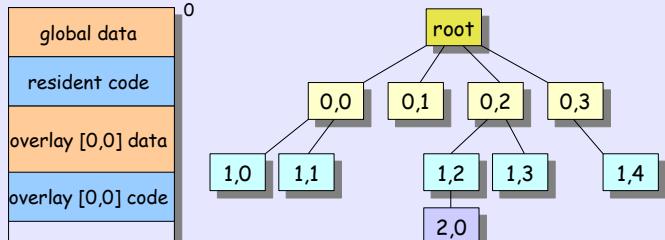
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

2

### Overlays

An early idea aimed at running large programs in small memory:



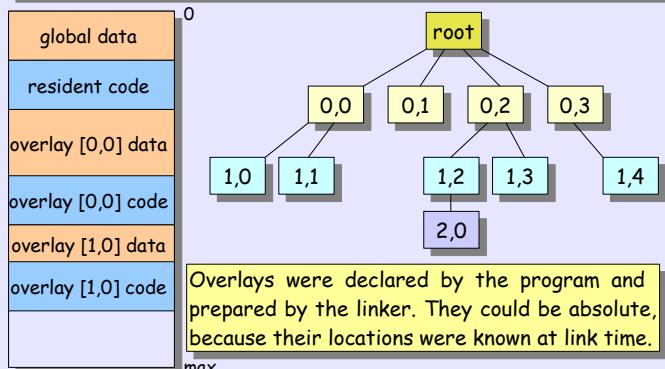
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

3

### Overlays

An early idea aimed at running large programs in small memory:



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

4

## OVERLAYS

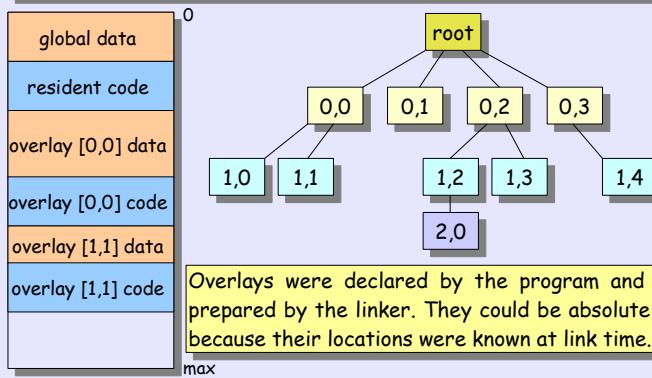
One of the first ways to manage memory is thought overlays. The program is partition by the programmer where he or she can override certain parts of memory if that portion is not being used. Similar to scoping.

The actual overlaying is called by the programmer in the parent overlaying node.

The overlay is called in a DFS like fashion, and this overlaying process was not automated.

## Overlays

An early idea aimed at running large programs in small memory:



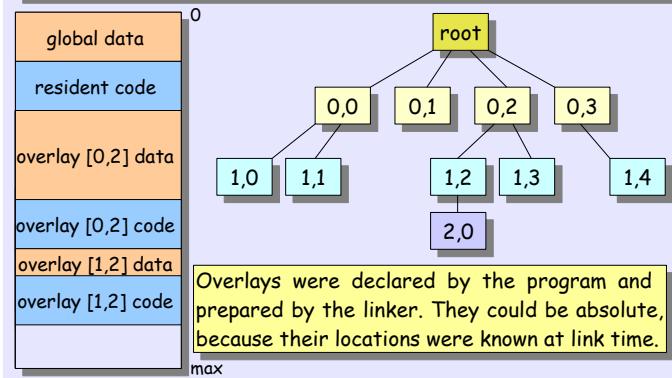
Overlays were declared by the program and prepared by the linker. They could be absolute, because their locations were known at link time.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Overlays

An early idea aimed at running large programs in small memory:



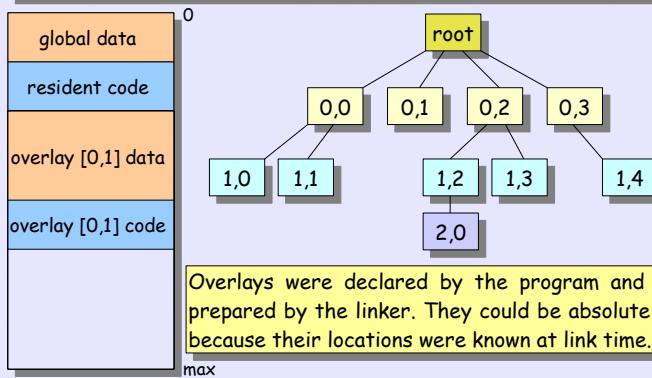
Overlays were declared by the program and prepared by the linker. They could be absolute, because their locations were known at link time.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Overlays

An early idea aimed at running large programs in small memory:



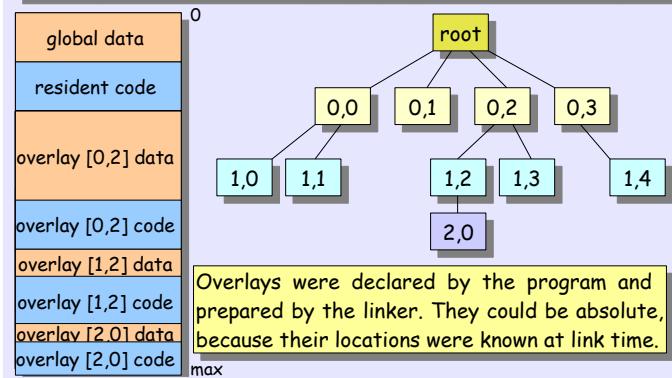
Overlays were declared by the program and prepared by the linker. They could be absolute, because their locations were known at link time.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Overlays

An early idea aimed at running large programs in small memory:



Overlays were declared by the program and prepared by the linker. They could be absolute, because their locations were known at link time.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

etc, you get the idea.

## MEMORY LAYOUTS OF EARLY PROGRAMS

Memory layout during the early days of computing were rather simple. It wasn't obvious that stacks are generally useful and that code and data should be separated. The prevailing programming language were FORTRAN for efficiency and COBOL for convenience. They were both static, neither allowed for recursion officially. Consequently, nobody cared about a built-in stack. Few of the popular CPU architectures supported it. Memory was viewed as a single, continuous, mostly static, chunk.

## WHY SEPARATE DATA AND CODE?

### A digression

**Q:** What is the single most important advantage of separating code from data?

**A:** Reentrancy, i.e., the ability to share code among multiple processes.

Prerequisites: 1) being able to enforce the read-only status of the code area, 2) having separate pointers to code and data.

Of course, the mechanics of program control, e.g., subroutine calls, should not rely on storing anything in the code area.

Recall fork!

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## TRAMPOLINES

### A digression: Trampolines

Even these days the purists have to yield sometimes. Recall Lazy Linkage. In some cases, a dynamically created function (typically put onto and run from the stack) is demonstrably the most efficient way to accomplish something tricky.

```
int A (int u, int v) {
    int c;
    int B (int w) {
        return w + c;
    };
    ...
    return B (v);
}
```

B is a nested function (such functions are legal in modern C (like gcc, for example). Note that, when called, it must know where A's local variables are.)

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

### Example

Hardware mechanism for subroutine call on CDC 6000 machines:

caller

...  
call F  
...

F callee

store at F the image of a branch instruction to C+1 and branch to F+1

...  
branch

The code image of the called function is modified.

There is a single static place (per function) where the return address is stored. Thus, there is no way to call the function recursively, unless you implement your own stack by hand.

All memory references issued by a program went to a single preallocated chunk of physical memory, so it wasn't easy to share.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

This basically tells you that stacks are very useful.

### Sometimes one misses the old days...

One reason for modifying code on the fly was efficiency. Its legitimacy can be disputed, but ...

```
void crunch (int N, double *M) {
    for (int i = 0; i < N; i++)
        if (some_global_option)
            M [i] = 1.0/M [i];
        else
            M [i] = M [i]/M_PI;
}
```

Such modifications were popular among library functions that would auto-tune themselves to global options, hardware parameters, or self-measured perceived performance.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

### A digression: Trampolines

Even these days the purists have to yield sometimes. Recall Lazy Linkage. In some cases, a dynamically created function (typically put onto and run from the stack) is demonstrably the most efficient way to accomplish something tricky.

```
int A (int u, int v) {
    int c;
    int B (int w) {
        return w + c;
    };
    ...
    return B (v);
}
```

B is a nested function (such functions are legal in modern C (like gcc, for example). Note that, when called, it must know where A's local variables are.)

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

### A digression: Trampolines

Even these days the purists have to yield sometimes. Recall Lazy Linkage. In some cases, a dynamically created function (typically put onto and run from the stack) is demonstrably the most efficient way to accomplish something tricky.

```
int A (int u, int v) {
    int c;
    int B (int w) {
        return w + c;
    };
    ...
    return B (v);
}
```

B is a nested function (such functions are legal in modern C (like gcc, for example). Note that, when called, it must know where A's local variables are.)

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## A digression: Trampolines

Even these days the purists have to yield sometimes. Recall Lazy Linkage. In some cases, a dynamically created function (typically put onto and run from the stack) is demonstrably the most efficient way to accomplish something tricky.

```
int A (int u, int v) {
    int c;
    int B (int w) {
        return w + c;
    };
    ...
    return B (v);
}
```

B is a nested function (such functions are legal in modern C (like gcc, for example). Note that, when called, it must know where A's local variables are.

Everything is fine until we want to have pointers to such functions.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

24

## Look at this:

```
int A (int u, int v) {
    int c;
    int B (int w) {
        return w + c;
    };
    int (*fp) (int);
    fp = B;
    c = u;
    return C (fp, v);
}
```

```
int C (int (*fp) (int), int q) {
    return (*fp) (q);
}
```

B is blind without an extra pointer that tells it where c can be found. Thus, it needs such a pointer whenever run. So how to represent pointers to functions like B? Note that the same issue concerns pointers to member functions in C++.

But note this: in whatever context B is legitimately called, its A must be present on the stack. The pointer to A's stack frame is a "momentary constant" of that context.

Winter 2019

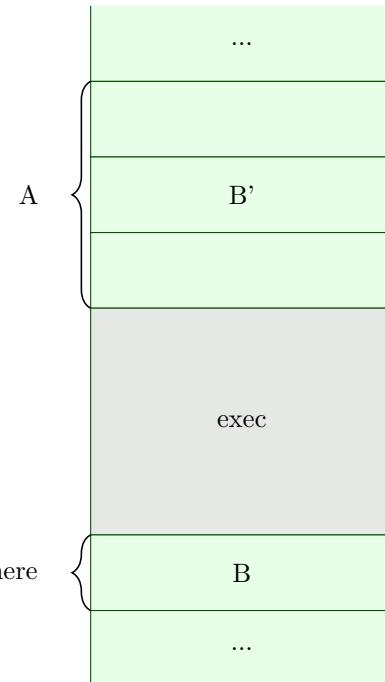
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

25

The Most efficient way out is to have on A's stack a trivial chunk of code. This is what it does:

1. Load the pointer to A's frame into the place (a register) when the called function expects it
2. Jump to the function according to its bare pointer (provided in some register)

Such chunks of code (called trampolines or thunks) are used in some other places, notably, for invoking signal handlers.

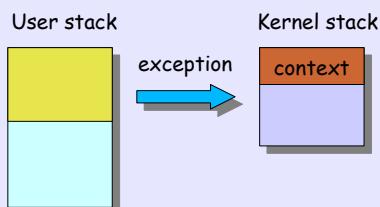


B' is responsible for executing B on the stack. Moreover, it is responsible for executing the jump and the passing of the function pointer to B.

## SIGNAL HANDLERS

Because signal handles are executed in the context of the receiving process started by the schedule.

## What about those signal handlers?

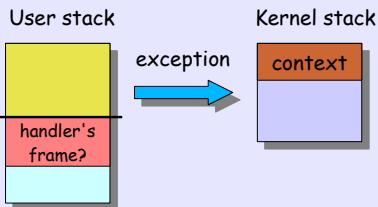


Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

28

Once the TS expiration is reached to the process. The important contents of the process is stored into the kernel's stack. Things like PC, things that are needed to resume.

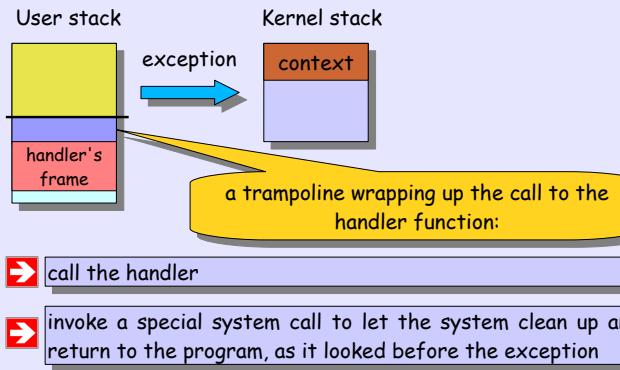
## What about those signal handlers?



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## What about those signal handlers?



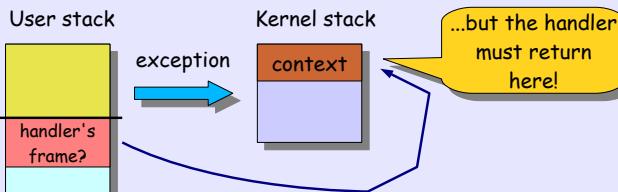
30 Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

34

The kernel needs to create an activation frame for the signal handler in the user stack, as the user does not have control.

## What about those signal handlers?



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## THE GOAL OF MEMORY ALLOCATION SCHEMES

The main goal of memory allocation is:

- **Efficiency and flexibility combined with the full dynamics of multi-programming.** Ideally, if there is free space to run a new program, it should be able to run it regardless of the condition of memory.
- **Separation of concerns** The system deals with the physical memory, while the program deals with the logical memory, the process only sees and deals with the logical memory.
- **Protection coexisting with collaboration opportunities.** collaboration like sharing should only be done in controlled and well defined environments.

The trampoline contains executable codes that execute system calls that tells the kernel that the receiving process did receive the signal. It also chain multiple signal handlers together and call the return signal system call to the kernel to tell that the current process is finished handling system calls.

# FIXED PARTITION

At the hardware level, memory is not uniform.

But at the software level, memory should be uniform as it's easier to work with.

Uniform model with addressing, but we specialize regions into segments. The general idea is that these segments which are continuous, can be placed anywhere in any order within memory.

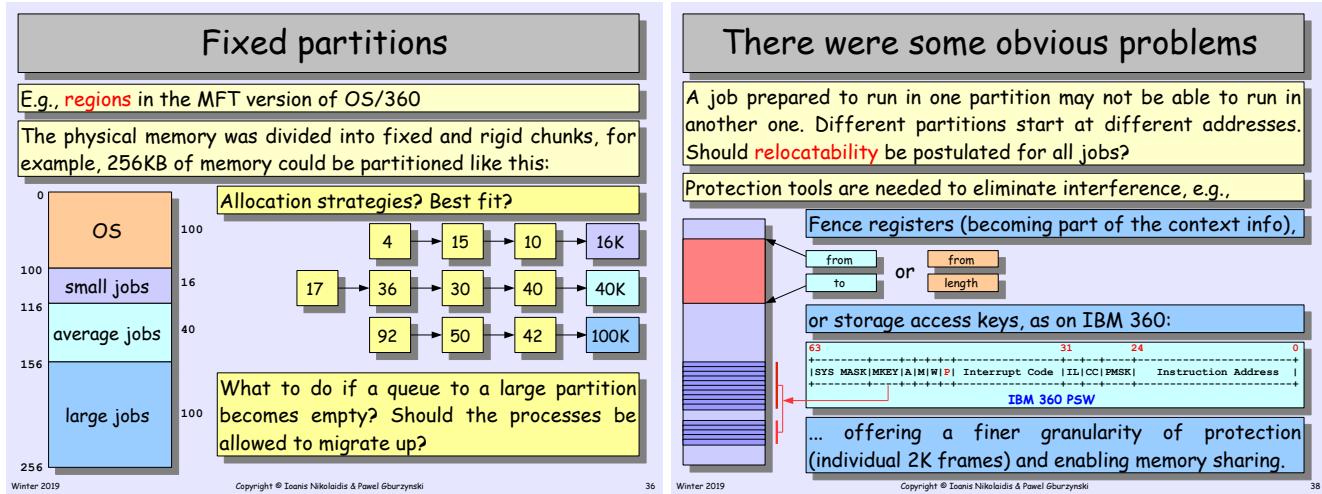


Figure 3: Fixed partitions and its issues

These types of memory allocation were used in the 60s. Programming languages during these time didn't support recursion or other run-time dependent factors that affects execution size. So the size of the execution of the program can be determined during compile time.

The problem with this type of memory allocation is that only a fixed number of processes can be running at any given time.

Fence registers are used as a way to prevent illegal memory access of the current process. These registers are part of the the context info and are known to the CPU.

For the storage access key method each segment of memory have a key associated with the region. Each process will have a key to memory location where it have access. Mismatched keys are considered illegal access. Catch all keys or all access keys are also valid

# FRAGMENTATION

## TYPES OF FRAGMENTATION

The partitioning of memory into fixed sizes poses two problems.

- **Internal Fragmentation:** the unused leftovers of partitions
- **External Fragmentation:** the fact that the partitions exist at all. i.e. the program is bigger than the partition even though you have enough physical memory.

What we are trying to solve is the elimination of fixed size partitions using dynamic partitioning of memory.  
What we want:

- Flexibility with the size of code.
- relocatability (seen this from position independent code)
- protection
- sharing (controlled)

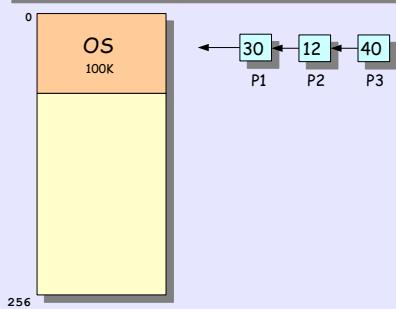
All of this has to be done in a space (no waste) and time (with CPU support "in real time", and "instantly") efficient manner.

## VARIABLE PARTITION

### Variable partitions ...

... e.g., MVT on IBM/360

The OS creates partitions automatically as jobs come (and go):



Winter 2019

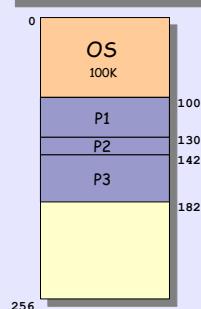
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

40

### Variable partitions ...

... e.g., MVT on IBM/360

The OS creates partitions automatically as jobs come (and go):



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

41

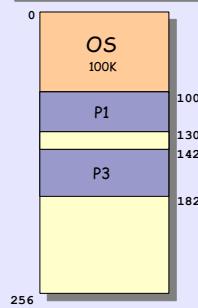
Processes are coming in

They are loaded into memory

## Variable partitions ...

... e.g., MVT on IBM/360

The OS creates partitions automatically as jobs come (and go):



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## GAP FILLING STRATEGIES

- **Best Fit:** The program is put into the tightest hole that can accommodate it. Trying to make the smallest hole possible.
- **Worst Fit:** The program is put into the largest hole.
- **First Fit:** Just put it into the first available spot.

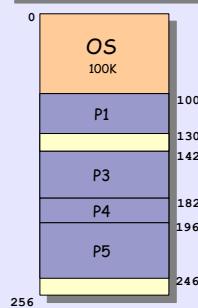
## BEST FIT VS WORST FIT

The logic behind worst fit is that after fitting in the program into memory you are likely to have enough space for another program. While best fit is more likely to have wasted space. First fit is a don't care policy.

## Variable partitions ...

... e.g., MVT on IBM/360

The OS creates partitions automatically as jobs come (and go):



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## COMPACTON

Can we move programs in memory to reclaim a single continuous free chunk?

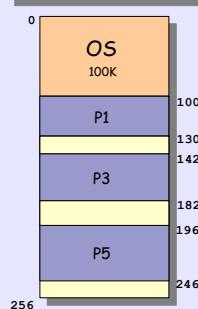
There is no trivial solution to this problem as very few programs are actually movable, **Recall: out discussion between relocability, position independence, and movability.**

This leads to the idea of abstracting physical address to logical addresses.

## Variable partitions ...

... e.g., MVT on IBM/360

The OS creates partitions automatically as jobs come (and go):



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## CDC 6000 EXAMPLE

The CDC 6000 have two special (protected) registers **RA** and **FL**. The two registers are reloaded as part of the context switch whenever the CPU is assigned to a process. This process does add a bit of complexity but accomplishes one important goal: we can begin to differentiate between logical and physical addresses in memory.

Here is the hardware implementation in C.  
The logical addresses of CDC 6000 starts from 0 - MAX.

```

1 // checks if the logical address is
2 // within the program
3 // p is the physical address
4 // RA (relative address) is the
5 // physical address where the process being
6 // stored in memory begins (starts from 0)
7 // FL is the length of the address
8 if (i < FL){
    p = RA + i;
}

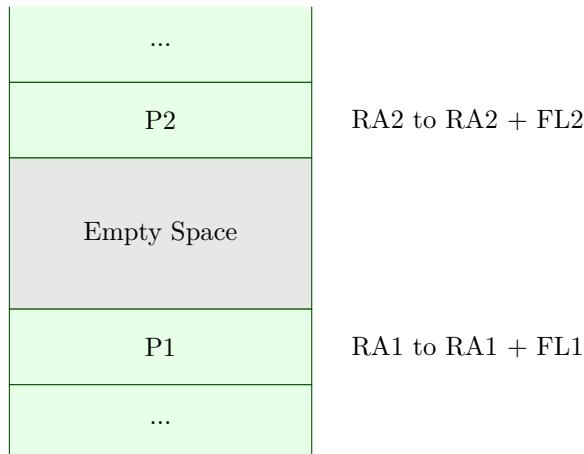
```

## MEMORY SHARING

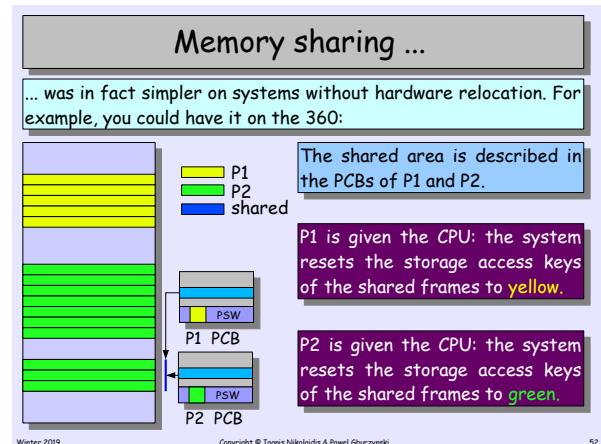
Due to the memory protections provided by the CPU there is no way to share memory between different processes. You can emulate it using system calls and pass messages around, but that involves passing things that shouldn't be done by the kernel in terms of responsibilities. One can try to overlap but that leads to a mess.

Another problem is that everything in the process is writeable. There is no protection within the process itself.

Below is the memory representation of processes



## MEMORY SHARING



On the IBM 360 allows sharing based on an ownership key. Moreover, there was a global ownership key where all processes can gain access.

## SWAPPING

Turning memory into a preemptible resource. Programs can be suspended and swapped out, then swapped in later. This is useful particularly with movable partitions. This allows you to have a pool of programs in execution be larger than the physical size of memory. This does lead to another level of scheduling where the program in swap needs to compete to get loaded into memory.

This allows the system to accommodate important programs without having to kill any of those currently in execution.

## WHEN TO SWAP OUT?

- The program has exceeded its inter-swap interval, and there are other programs waiting to be swapped in (fair processing)
- The program has reached a stage where it has to wait for a longish time. Not necessarily an I/O request, but "truly external". This is because a program's memory can become a non-preemptible resource.
- A high priority program requires immediate execution and there is no other way to accommodate it.
- With some partition systems, like one with compactions, a program can grow or shrink. When a program grows too much it may need to be swapped out until a more favorable situation arises.

# PAGING

Pretty much every OSes that is worth taking about have paging as a feature in memory management.

## FRAMES AND PAGES

**Frames** are physical memory that is divided into fixed-size blocks. The logical address space of a program is divided into chunks of the same size and are called **pages**.

There is a mechanism for mapping any page into any frame. Some parts of the logical space may be unmapped (holes are OK).

### Typical page sizes

Intel x86	4096 bytes
MC68451	256 bytes minimum
IBM-370	2048 or 4096 bytes

The page size is always a power of two. This makes the interpretation of a memory address easier. Below is the format for memory access in a 32-bit system.

## LOGICAL MEMORY ADDRESS



The **page number** is the offset from the head of the page table to get the correct entry for the frame number. The **offset** is the offset within the page/frame.

## PAGE TABLE ENTRY



The **frame number** is the frame index in physical memory.

The **attributes** tells the status of the frame. i.e., is it **valid** (is the page mapped at all), **accessible** (R/W, R/O, X/O)

The mapping of pages into frames is described in a special data structure called the **page table**. When memory is being referenced the page number is transformed into a frame number (the physical address) by the page table. This operation is done in hardware, essentially at every memory reference (at least in principle). The page table is part of the process context where it describes the complete memory layout of the process.

The **Page Table Pointer** pointing to the head of the page table. They are stored in a CPU register for the current process that is being executed.

The benefits of this are: **No external fragmentation, very little internal fragmentation. Logical address space is separated from physical memory, natural protection, and easy memory sharing.**

## CACHING

---

To cache access to the memory we use the Translation Lookaside Buffer. Once a page is accessed the entry is copied into the TLB so lookup is faster. During a context switch, we need to flush the buffer as the values are garbage, only valid for the process that was using it.

## TLB

---

The entry in the TLB is the frame pointer to the physical memory as we want the physical translation. This is because we want locality, in general the next instruction is going to be near the frame stored in the TLB.

## ADDRESS TRANSLATION

---

There are situations where the system must operate in physical memory. The actual translation of the address can be switched on or off by the system by changing a certain bit in a protected register.

Reason why?

- Accessing special locations, e.g., corresponding to registers of peripheral devices.
- Supplying addresses (buffers) to peripheral devices as parameters of physical I/O operations.
- Setting up a page table

### PERIPHERAL DEVICES

Note: generally, peripheral devices cannot take advantage of address translation as these translations are context-relative (different processes see different logical address spaces). When a peripheral device is running, the CPU context usually describes an unrelated process.

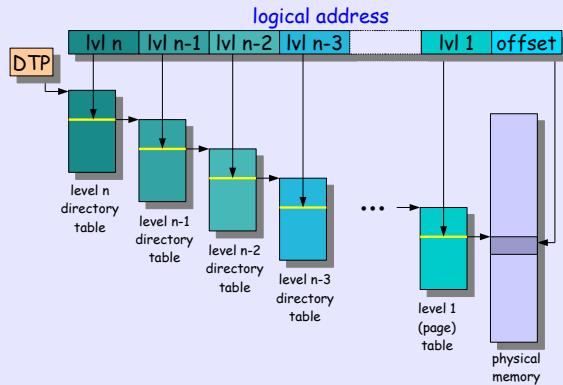
## MULTI-LEVEL TABLES

---

The actual (used) portion of a program's logical address space does not need to be continuous. Also, different chunks can be assigned different accessibility attributes. Thus, it makes sense to assume that the logical address space has the same generic layout for all programs (large and small).

Problem is that some programs are so small that having a fixed size page table for all programs will be overkill. To solve it we use hierarchical design. This design can be slow for reading as you have to read n levels of redirections but the buffer (TLB) fix this issue.

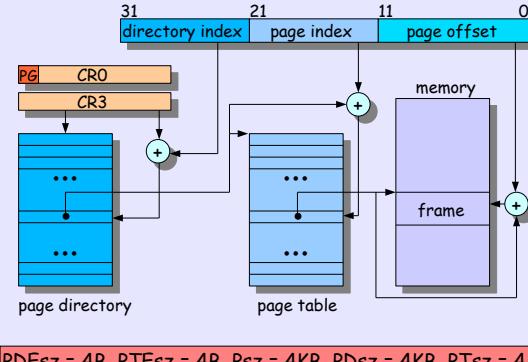
## Generally, we may have ...



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Example: 32-bit Intel



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

61

The n-th level contains the offset for the n-1-th level in the directory table to gain access to the page table. From there using the offset allows us to get access to the frame without storing a whole bunch of invalid page tables.

The 32-bit Intel is an example of how a hierarchical design can work. This is a two level entry.

Page directory entries that are invalid will not have a page table as they would be completely empty. Therefore, you will only have page tables that have a valid entry.

### INTEL EXCEPTION

Intel also does segments as well as this form of memory management.

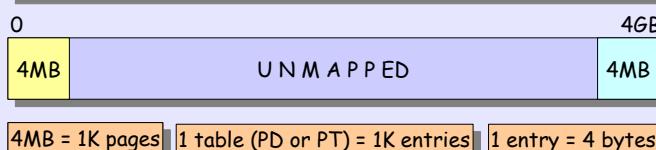
## WHAT DOES IT BUY US?

### What does it buy us?

Quite a bit of complication for one thing. But suppose that a program wants a little something at the beginning of the address space + a little something at the very end.

With one level, we have:  $PTS = PTEsz \times (4G/Psz) = 4MB$  of tables.

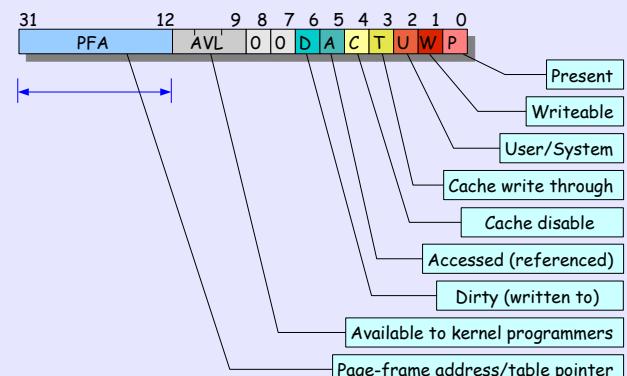
With two levels, we have:  $PDSz = 4KB$ ,  $PTsz = 4KB$ . Total =  $3 \times 4KB = 12 KB$  of tables. With this much, we can map a 4MB piece at location 0 and a 4MB piece at the end (e.g., for the stack).



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

### PDE/PTE format



Winter 2019

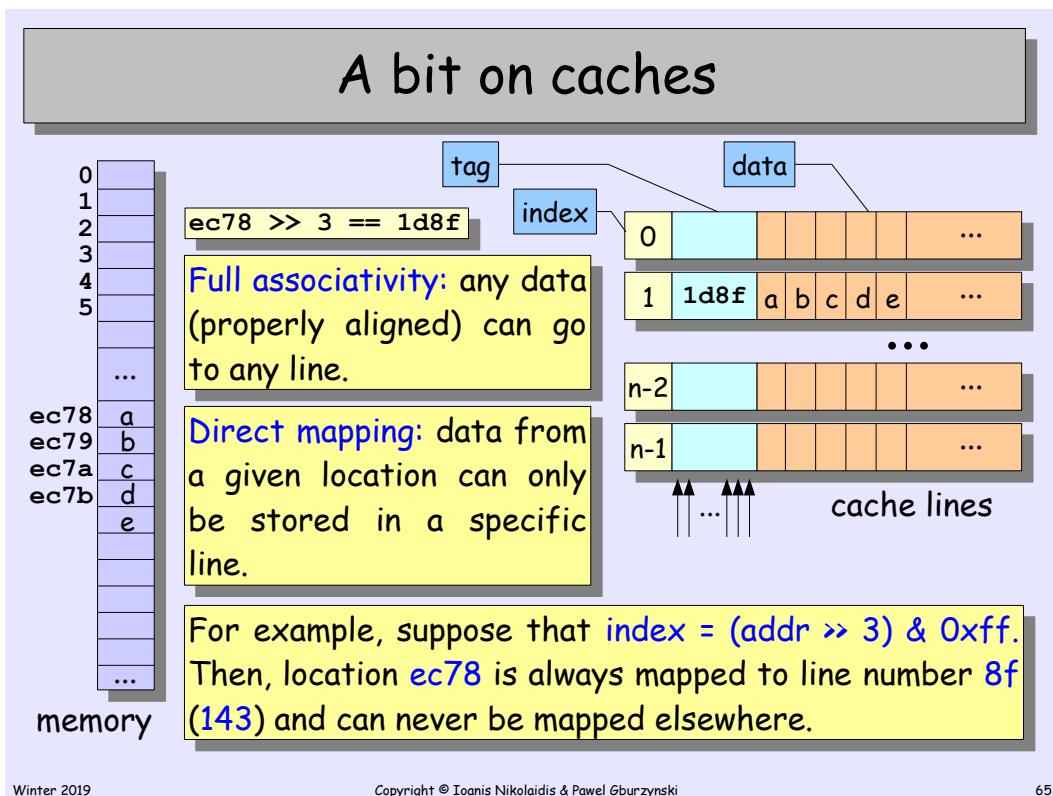
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

63

## A FEW NOTES

The size of the PFA determines the maximum amount of physical memory that the system can use. Attributes can be specified on a per-page or per-directory basis. Generally, it makes sense to associate accessibility attributes with chunks larger than pages. Sometimes, in two-level hierarchies, the directories are called **segments**. This reflects the fact that they provide a coarser grain of allocation - one that better corresponds with logical segments of programs. With all these levels of indirection, the cost of a memory reference is bound to skyrocket, unless something is done. This is taken care of by some extra hardware called the **TLB**.

## CACHING



The tag is the prefix in memory where the content in memory is stored. Generally the cache line is fixed and is a power of two, in this case its 8 bytes. The write back policy is to update the entry in cache, and not to the memory unless the cache needs to be flushed. The write through policy writes it both to cache and memory. Both are important in certain cases.

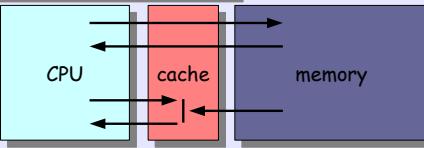
### CHARACTERISTICS OF A CACHE

- **Cache Size:** The number of lines and line size.
- **Associativity rules:** Fully associative, direct mapped, some hybrid.
- **Replacement Policy:** How the lines are recycled. Note associativity rules may restrict this policy.
- **Write policy:** write-through vs write-back – write to memory every time the cache is been updated vs write to memory when the cache is flushed and has been modified.

## LOCKING

**Spin locks and memory cache**

A spin lock typically looks like this: `while (T&S (&flag));`



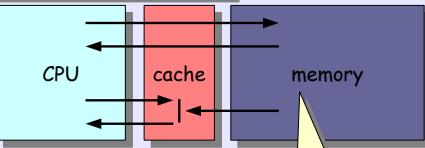
Test & Set usually bypasses the cache and occupies the bus. Sometimes this may be a better way:

```
void spinlock (int *flag) {
    while (1) {
        if (!T&S (flag)) break;
        while (*flag != 0);
    }
}
```

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzyński

**Spin locks and memory cache**

A spin lock typically looks like this: `while (T&S (&flag));`



Test & Set usually bypasses the cache and occupies the bus. Sometimes this may be a better way:

```
void spinlock (int *flag) {
    while (1) {
        if (!T&S (flag)) break;
        while (*flag != 0);
    }
}
```

If the value has changed, the cached copies in other CPUs' caches must (should) be invalidated.

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzyński

The during a spin lock the cpu will consume the entire memory's bandwidth during this situation as it wants fresh information, not information from the cache. This is undesirable.

The test and set is done only once reading directly from memory, if the core didn't win the competition then read from the cache. Once the flag in memory is updated the cache in the cpus are updated/invalidated so that they can start to read from the memory.

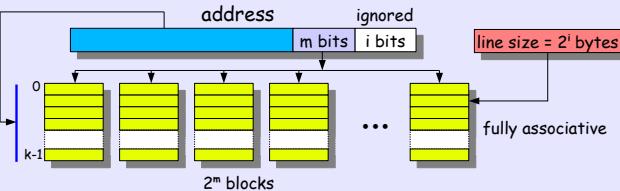
## ASSOCIATIVITY TYPES

**Associativity types**

Full associativity is nice but expensive to implement. Direct mapping is simple, but gives lousy performance.

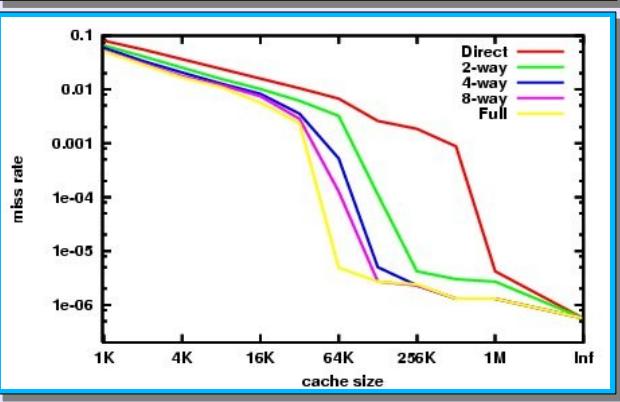
In practice, we often see some compromises between the two associativity extremes, e.g., k-way caches.

A cache is **k-way**, if any given memory location can be cached in any of **k** specific lines. Usually, **k** is a power of two.



Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzyński

**Associativity tradeoffs (example)**



miss rate

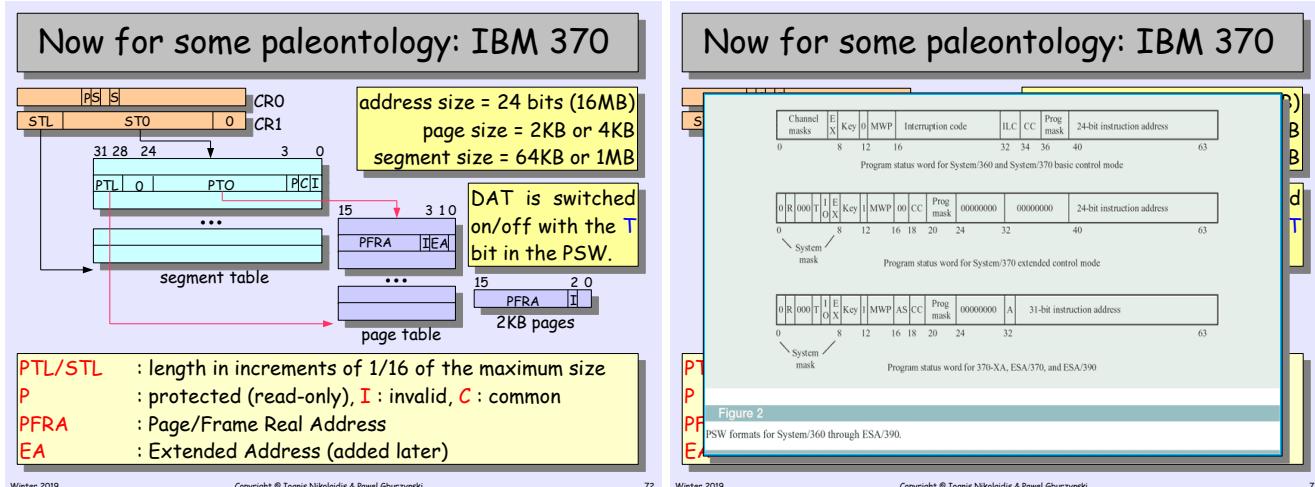
cache size

Direct  
2-way  
4-way  
8-way  
Full

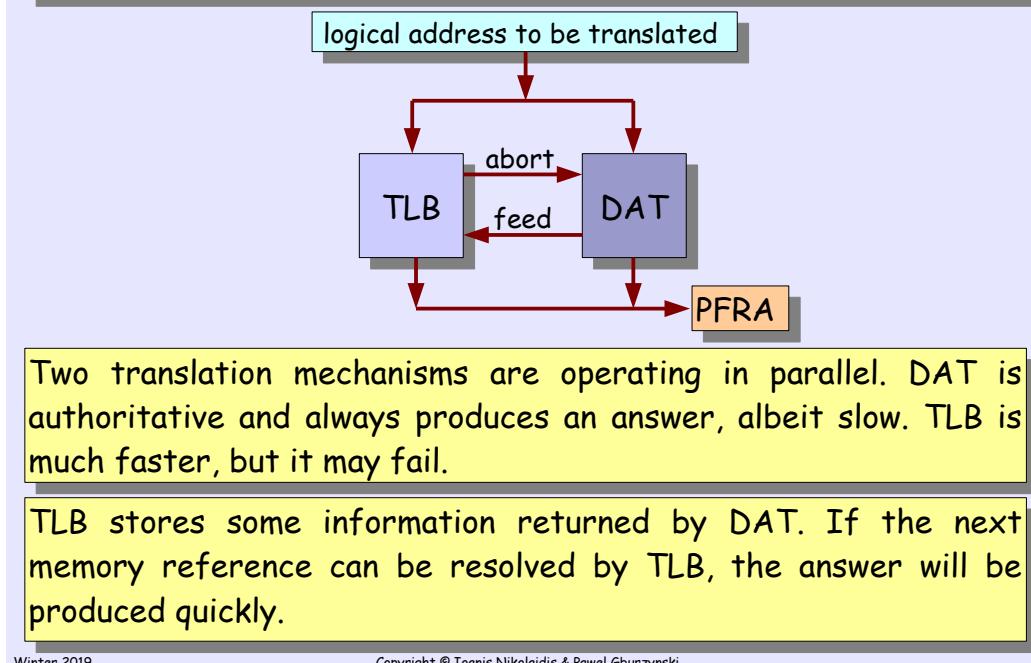
Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzyński

You get the best performance with full associativity. But is more expensive.

# IBM 370 TLB



## TLB on the 370: the mother of all TLBs

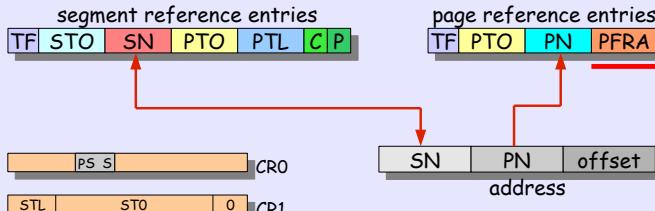


The IBM 370 have two forms of translation mechanisms, the faster TLB which may fail, and the slower but always correct DAT. These two mechanisms operate in parallel when the cpu is trying to access memory. The TLB will always finish faster if the memory address is stored within its cache. However, if it does not finish then the DAT will resolve the translation and store the result into the TLB for caching purposes.

## HOW THE TLB WORKS

### TLB on the 370 ...

... is a specialized cache capable of storing a limited number of entries of two types:



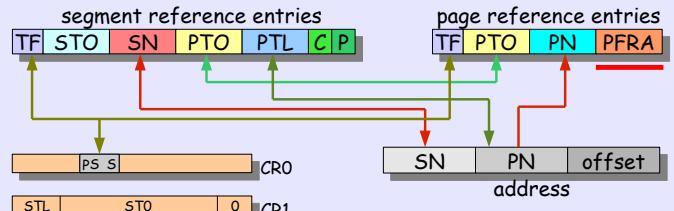
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

76

### TLB on the 370 ...

... is a specialized cache capable of storing a limited number of entries of two types:



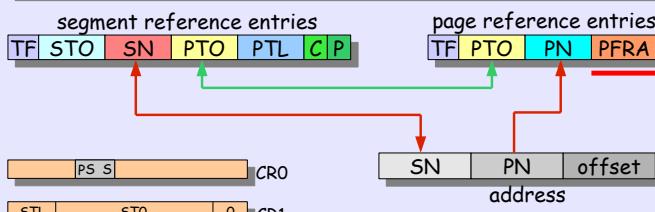
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

79

### TLB on the 370 ...

... is a specialized cache capable of storing a limited number of entries of two types:



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

77

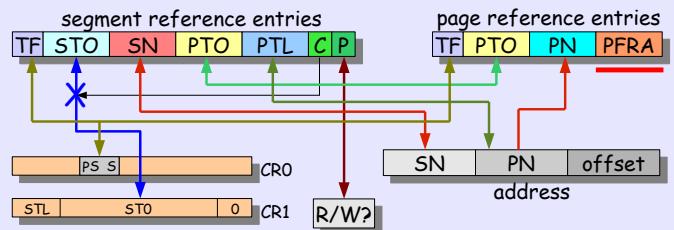
We maintain the PFRA and the PTO in the TLB. There must be sanitity checks done to be sure that data is part of the process.

Is the STO in the segment reference entry match the CR1, if so then its a cache hit.

This is done very quickly (sub nano seconds) but very expensive.

### TLB on the 370 ...

... is a specialized cache capable of storing a limited number of entries of two types:



One more little twist: If C is set, the STO matching is not verified. This accounts for the so-called common segments.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

78

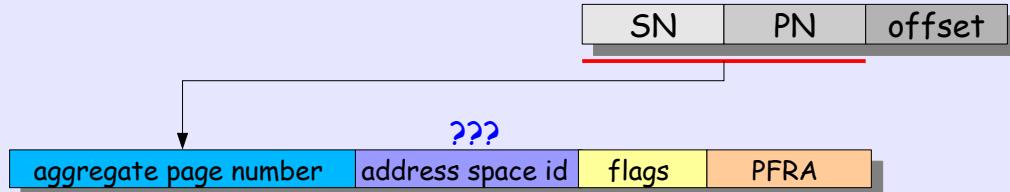
is the PN less than the PTL? This is part of the sanitity check.

P is the premission of the segment reference

82

### In principle, it can be simpler

An entry can only store information about pages, ignoring the segments:



This is in fact how most TLBs are implemented today. Note that the associativity string (APN) is longer than in the previous case. One can argue that using two entry types (segment/page) improves the accuracy of the TLB - for the same expense understood as the total amount of storage.

In the defense of today's implementation, space is cheap.

## COHERENCE OF TLB CACHES

Entries in the TLB may become **obsolete** (and incorrect) when the actual mapping changes, e.g., a page is unmapped or remapped. This tends to happen quite often. Upon a context switch, completely new tables become active. There are two solutions to this problem.

1. **Completely invalidate the TLB.** This used to happen on the old Intel CPUs. Writing to CR3 automatically invalidates the TLB.
2. **Retain all entries, but make it possible to distinguish different address spaces** (the 370 approach)

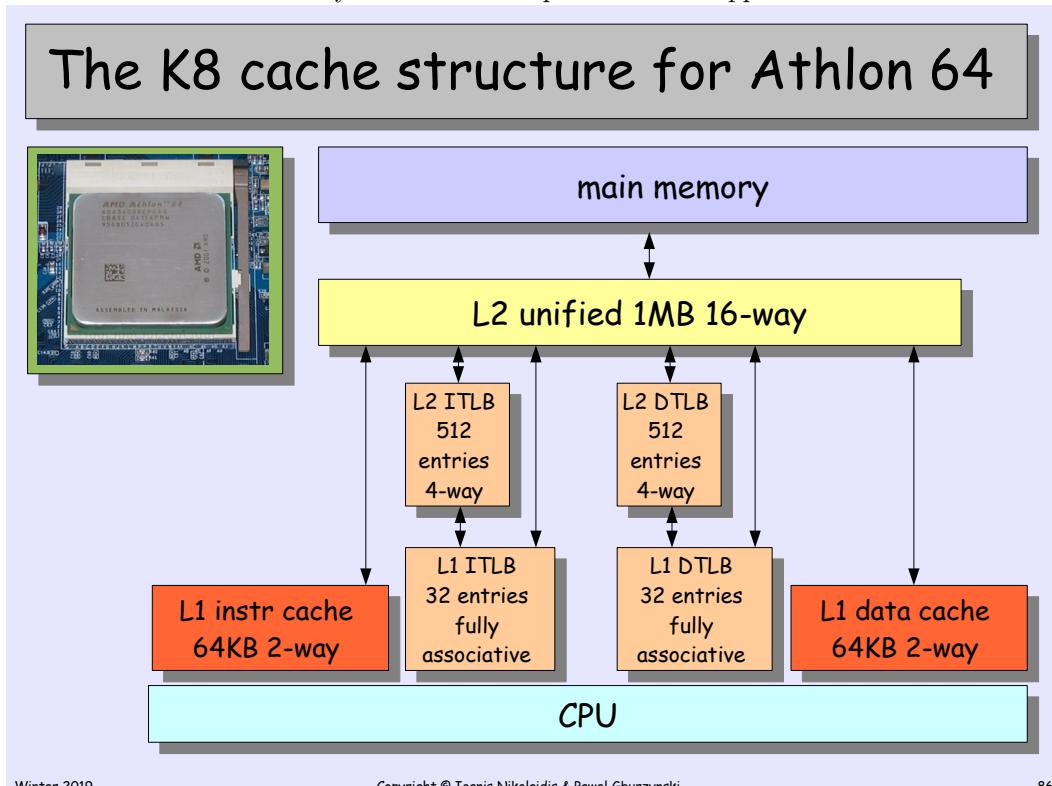
The latter is more efficient in the face of frequent context switches. Additionally, the entries pertaining to the common segment are never invalidated (as long as its mapping doesn't change)

## COMMON SEGMENT

A popular approach to organizing the logical address space of a process is to include the kernel memory in it.  
This is the case in older Linux kernels.



We not only execute the kernel on the process's stack, but also within the process's address space. The kernel part is in fact a common segment: it occurs in all address space (always in the same location). Currently, modern 32-bit systems tend to depart from this approach.



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

86

Starting from a clean state will cause the cpu to access the main memory quite often, during which will fill the cache. But once the loading is done then subsequent calls will probably be faster due to the locality of data and instructions.

## DYNAMICS OF ADDRESS SPACES

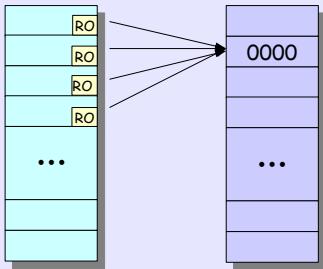
By now, we should be comfortable with the idea that the logical address space of a process contains holes. Suppose that a program gets loaded into memory. How does it get loaded. DLL are one method, but what about the stack? How can we know how much of it the program is going to need? Can we somehow figure it out automatically (without forcing the program to tell us anything) when the program is going to use some memory, so that we can bring it in as needed?

## A SIMPLE TRICK

A situation of an uninitialized array.

### A sample trick

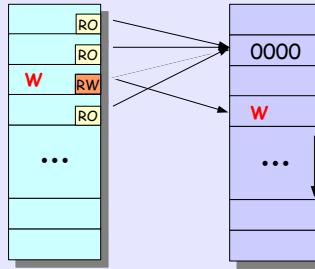
Suppose that a program declares a huge array filled with zeros.



If the array happens to be constant, we can trivially save by using only one frame for the entire array. What if it isn't?  
Mark all the pages as read only.

### A sample trick

Suppose that a program declares a huge array filled with zeros.



If the array happens to be constant, we can trivially save by using only one frame for the entire array. What if it isn't?  
Mark all the pages as read only.  
Whenever the program writes something to the array, the system will receive an exception (formally an error).

The system will automatically correct it, i.e., copy the zero frame to a new frame, mark it as read/write, and resume the program.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

88

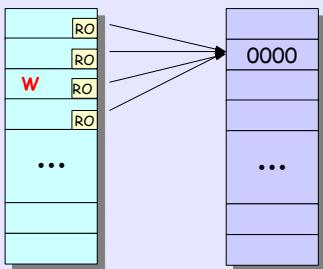
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

90

Imagine an array of uninitialized data. To prevent redundant frames, the elements of the arrays all point to the same frame. The elements are read only as they don't want to change the default value.

### A sample trick

Suppose that a program declares a huge array filled with zeros.

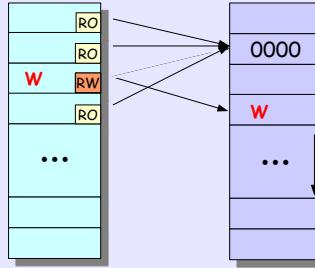


If the array happens to be constant, we can trivially save by using only one frame for the entire array. What if it isn't?  
Mark all the pages as read only.  
Whenever the program writes something to the array, the system will receive an exception (formally an error).

For that page the system remaps it to a new frame where it can now be written.

### A sample trick

Suppose that a program declares a huge array filled with zeros.



If the array happens to be constant, we can trivially save by using only one frame for the entire array. What if it isn't?  
Mark all the pages as read only.  
Whenever the program writes something to the array, the system will receive an exception (formally an error).

The system will automatically correct it, i.e., copy the zero frame to a new frame, mark it as read/write, and resume the program.

This is in fact how uninitialized static data area (.data) is set up.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

91

Lets try to writing something into the 3rd element.

But all the pages are marked read-only.

# COPY ON WRITE

Recall fork. This is how the operation is carried out. The child inherits exactly the address space of the parent. No memory is copied. All memory that is not to be actually shared is marked as Read only + Copy On Write. For as long as the two processes **do not write to some area**, they are going to share the frame. This is why fork is so simple and inexpensive physically, while being so powerful as a concept. This hints at a more general idea. Let us call it CRAPE for "Creatively Responding to Apparent Program Errors." The program errors we have in mind are called **page faults**.

## MORE TRICKS

When a program pushes something onto a nonexistent page of its stack, the system allocates a blank new frame and resumes the program. **No need to load the entire program at once.** As the program runs through its code, the system receives an exception whenever the program references a new (so far untouched) page. Then it brings the new code fragment into memory and resumes the program. **No need to allocate the entire data at once.** by the same token as before. Suppose that the system wants to **emulate** something (e.g., a device) as a bunch of memory locations. It can easily intercept all references to a given page range and fool the program into thinking that is actually reads from and writes to memory.

### Memory mapped stuff, typically files

Here is a standard system call available on UNIX systems:

```
void *mmap (void *start, size_t length, int prot,
            int flags, int fd, off_t offset);
```

It allows you to map a file fragment into memory. In a nutshell, this is how a program is loaded for execution. Recall this:

```
PHDR off 0x00000034 vaddr 0x08048034 paddr 0x08048034 align 2**2
  filesz 0x000000c0 memsz 0x000000c0 flags r-x
INTERP off 0x00000004 vaddr 0x080480f4 paddr 0x080480f4 align 2**0
  filesz 0x00000013 memsz 0x00000013 flags r--
LOAD off 0x00000000 vaddr 0x08048000 paddr 0x08048000 align 2**12
  filesz 0x0000059d memsz 0x0000059d flags r-x
LOAD off 0x0000005a0 vaddr 0x080495a0 paddr 0x080495a0 align 2**12
  filesz 0x00000114 memsz 0x00000134 flags rw-
DYNAMIC off 0x0000005ec vaddr 0x080495ec paddr 0x080495ec align 2**2
  filesz 0x000000c8 memsz 0x000000c8 flags rw-
NOTE off 0x000000108 vaddr 0x08048108 paddr 0x08048108 align 2**2
  filesz 0x00000020 memsz 0x00000020 flags r--
```

Let us call this idea **lazy loading**. The program's components get loaded gradually as they are (implicitly) demanded.

### Processes versus threads

Threads are strongly related processes that typically share all memory (and most other resources) except for the stack.

Some systems do not differentiate much between processes and threads. Some others economize on internal data structures.

processes

threads

# PREREQUISITES FOR CRAPE

## Prerequisites for CRAPE

In order for CRAPE to work, the respective program faults must be 100% correctable. This means that:

The fault (an internal interrupt) should return enough information for the system to identify the exact source of the problem.

The partial effects of the last instruction (the one that has triggered the fault) must be accountable for.

The diagram shows two instruction sequences. On the left, for an ADD instruction (A + B → C), the steps are: fetch instruction, fetch (A), fetch (B), compute, and store at C. On the right, for a LOAD instruction (R1, R2, OFF(R1)), the steps are: fetch instruction, load R1, and load R2. A red arrow points from the ADD sequence to the LOAD sequence. To the right of the LOAD sequence is a memory model showing three memory locations: a green block above, a yellow block in the middle, and an orange block below. A text box states: "If R1 is overwritten before the fault, the instruction cannot be restarted." At the bottom left is "Winter 2019", at the bottom center is "Copyright © Ioannis Nikolaidis & Paweł Gburzynski", and at the bottom right is "96".

There are two approaches to solve this issue

1. Before executing a potentially unrestartable instruction, check if all the data it needs are in memory. If not, trigger an exception before starting to execute the instruction.
2. While executing the instruction, keep track of its partial effects. There are two alternative solutions regarding what to do after a fault.
  - (a) Undo the partial effects before triggering the fault.
  - (b) Save the partially done state, and when the instruction is restarted (or resumed in this case) - just keep going

The first one is tough, the second one easier

## Simple instructions cause no problems

... assuming that the CPU has been designed with CRAPE in mind.

LOAD R1-R7, OFFSET(R1)

If the range of the data block is well contained, say, no more than one page, it makes sense to check beforehand.

Alternatively, the CPU can load the data into intermediate (internal) registers first and only move them to the target registers at the successful completion.

MOVE (R2)+,-(R3) | load (R2); R2++; R3--; store; R2--; R3++; fault

Here the CPU can keep track of the increments/decrements and undo them before triggering the fault (PDP 11/45).

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

## Instructions that operate on large data

MVCL rd, rs | move characters long (IBM 370)

rd	destination address
rd+1	destination length
rs	source address
rs+1	pad source length

When the instruction is interrupted, the registers reflect its partial state of execution, so it will be continued when re-executed.

MOVC3 len, src, dst | move characters 3-arg (Vax 11/780)

Here, the arguments can be general, e.g., address constants. When the instruction is interrupted:

R0 = remaining count, R1 → next src byte, R3 → next dst byte

The PSW has a special flag == "instruction partially done". When set, the instruction takes arguments from the registers instead of using the "original" ones.

98

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

99

R2 is used as a location of the address. load the content of R2, and then in one step auto increment R2. R3 is the store location, decrement R3 then load the content. This basically copy an array.

The first implementation use additional registers to keep track of partial completion, the second implementation sets a flag to tell that the instructions are incomplete

## On the Intel ...

... instructions like those are implemented via prefixing, e.g.,

mov ecx, length  
mov esi, src  
mov edi,dst  
**repe movsb**

The prefix (**repe**) means: repeat until **ecx** reaches zero (decrementing **ecx** after each execution).

Although **repe movsb** is assembled as a single instruction, it can be safely interrupted, as the respective registers keep track of its partial execution.

Contemporary CPUs avoid too fancy instructions in terms of the complications of their arguments (not functionality), such that the extent of changes before a fault can be contained and rendered harmless from the viewpoint of CRAPE.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

100

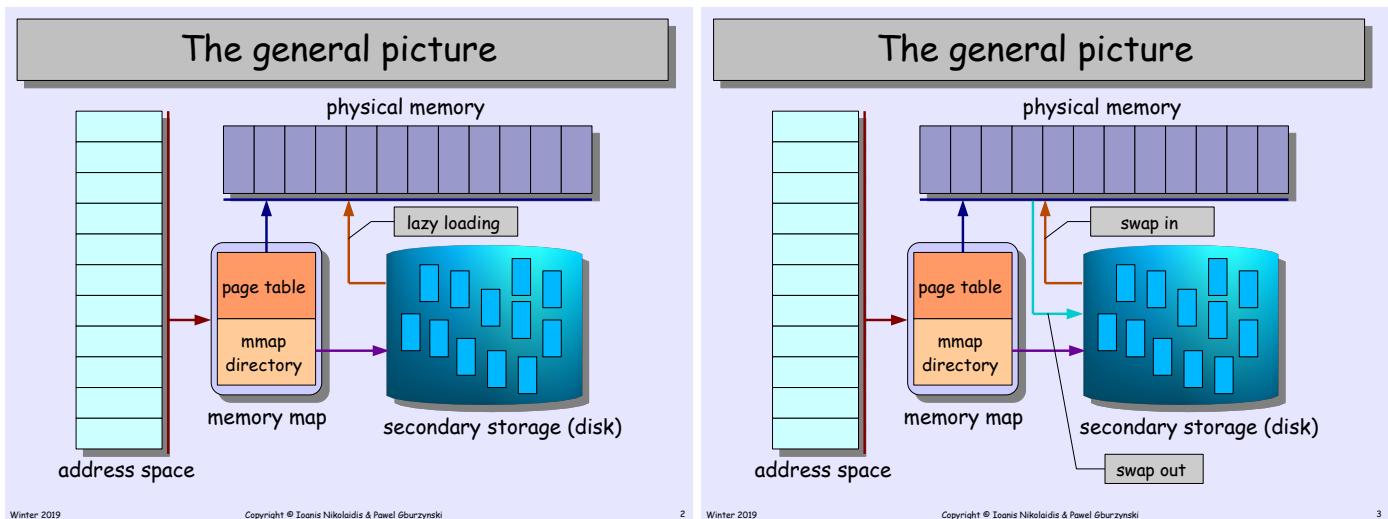
On Intel, the ecx contains how many steps are left when the process suffered a fault.

# 10: VIRTUAL MEMORY

We determined that swapping is indispensable:

Physical memory is never tied up to the current configuration of processes in execution. Priority jobs can always be admitted. Physical memory is made preemptible. Thus, memory-related deadlocks can be avoided.

In a partition-based system, a process can be either completely swapped out or completely swapped in. A partially swapped-in state is useless (from the viewpoint of the process). With CRAPE, we can be more flexible. For example, lazy loading upon (implicit) demand looks like a partially swapped-out state. The process can run even though some of its stuff is not in memory.



## SWAPPING OUT

Note that the static components of the program's address space, i.e., code + constants, are already "swapped out" before the program starts. They are never going to change. Thus, only the dynamic components, i.e., data + stack, will need a special swap area in secondary storage. Also, only these components will ever have to be physically swapped out.

The idea behind this whole scheme is this:

A partially swapped out program may still be able to run. It may not be needing some parts of its allocated memory area at this very moment. Thus, we may be able to run more programs in whatever physical memory we have than it would seem at first sight.

## SWAPPING AT THE PAGE-LEVEL

---

When a program demands a new page, the OS may have to find room for it. Basically, it needs to decide what page should be evicted.

The question now becomes how do we do this optimally, i.e., minimize damage, minimize the cost of eviction, etc.

## HOW MUCH MEMORY DOES A PROCESS NEED

---

The system may come up with some idea and assume that the process will be able to live within the confines. This imposes another restriction on the degree of the fanciness of machine instruction.

The question now becomes: What is the maximum number of pages that may be needed simultaneously by the most fancy machine instruction?

## PAGE REPLACEMENT POLICIES

### GOAL

To find the page with the highest likelihood of not being immediately needed. We want to maximize the time during which the victim page will not be demanded back. This interval determines the page fault rate in the system.

This is easier said than done due to the random nature of operating systems. The best we can do is make an educated guess.

There are two cases for a page to be purged. It can be **clean** or **dirty**. In the former case, it doesn't have to be sent back to secondary storage. The dirtiness of the page can be determined by a flag that is sent once the page has been modified. Some logical segments are always clean like code, constants, etc. But the data page can be clean as well: typically at least 80% of data references are for reading rather than writing.

## SOME IDEAS

---

**Random replacement:** The victim is selected at random. This policy does not get good numbers, but not horrible to implement.

**Random replacement+:** The victim is selected at random from the pool of clean pages first before considering dirty pages. This idea is generally bad as it will purge out all code/constant pages after a while. We want to find a page that is likely not going to be needed for a long time. Then look at its clean/dirty status. **Not the other way around!**

# FIFO

Select the oldest page in memory. This one is easy to implement.

The system maintains a list of pages. Upon a page fault it manipulate the link list by evicting the first page in front of the list. The new page is appended onto the back of the list.

The operation is constant  $O(1)$

The rationale makes sense: a page that has been in memory for a long time is likely to have been explored and thus not needed.

## PROBLEMS WITH FIFO

However, it does not account for permanently needed pages. An example is shown below.

Being old does not mean useless. Sometimes following the doctrine of the Dawi is a good thing.

```
1 double sum;
2 int i;
3 double HugeArray [N];
4 // more code
5 sum = 0;
6 for(i = 0; i < N; ++i){
7     sum += HugeArray[i];
8 }
9 // more code
```

`i` and `sum` are permanently needed, even though the traversed and explored fragments of `HugeArray` may be not.

This is rather typical. In any computing problem there are data that are needed all the time like counters.

# LRU (LEAST RECENTLY USED)

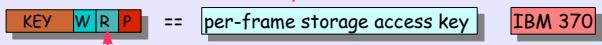
Purge the page that has not been used for the longest time. This is not difficult conceptually and can use the same data structure from FIFO. However, its behavior is different. It manipulate the link list at **every** access. For every access you have to bring that data to the back of the list. Therefore, it is more complicated than FIFO. When we need to purge a page, the head points to the page that hasn't been used for the longest time. Unfortunately, this calls for an expensive action to be performed after every memory reference. This is why LRU can only be approximated.

However, if you create the data structure simply then each operation can be quick.

Moreover, "LRU" support distinction of referenced vs not-referenced pages.

## How to approximate LRU?

The system needs a way to tell whether a page has been referenced within some time interval.



If there's no 'referenced' bit, CRAPE comes to the rescue:

Instead of clearing the referenced bit, mark the page as invalid without evicting it.

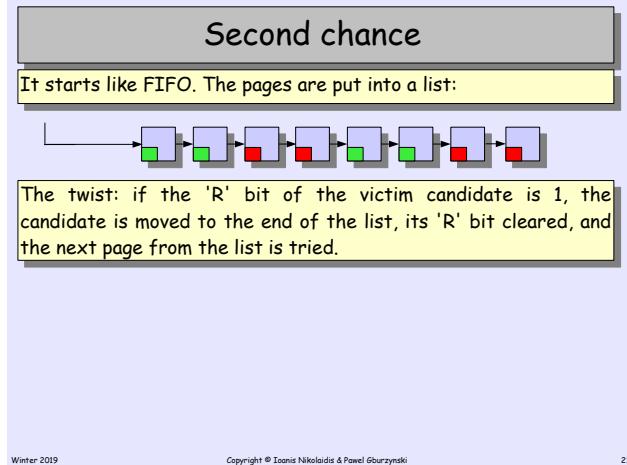
On a page fault, mark the page as valid and note the fact that it has been referenced.

## TIME INTERVAL

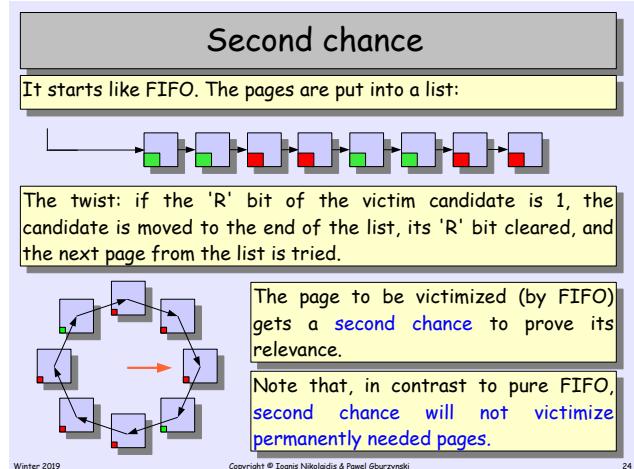
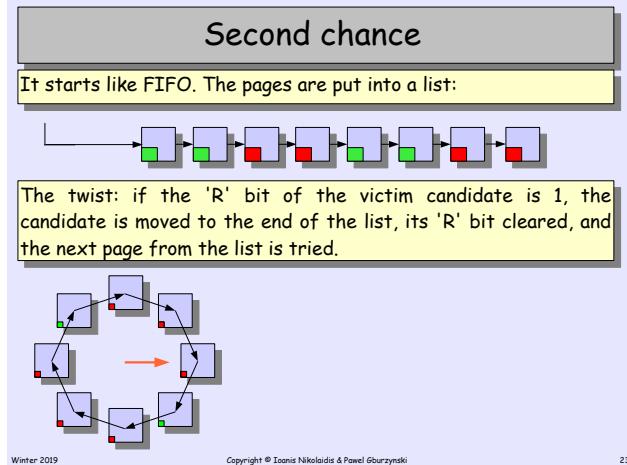
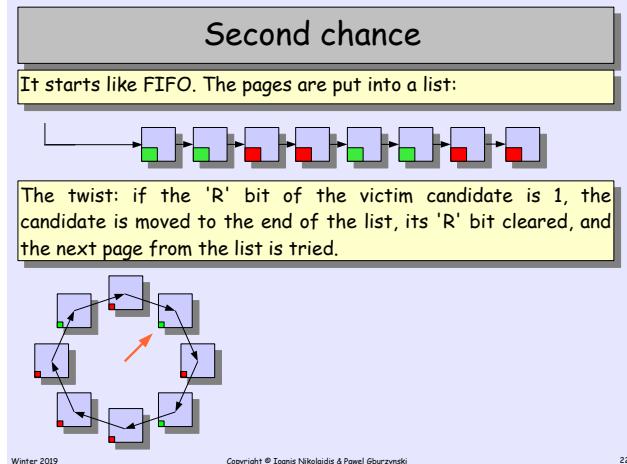
An idea to solve this LRU approximation is to go through the list at a fixed time interval, or in this case # of memory references. If it finds a page whose 'R' bit is set, it moves the page to the end and clears the 'R' bit. This isn't very smart. What if the last page gets referenced after having been checked and just before the fault? We can do better.

## SECOND CHANCE (FIFO OPTIMIZATION)

We stated that FIFO has problems storing permanent pages. Second chance will help fix this problem.



The twist: if the 'R' bit of the victim candidate is 1, the candidate is moved to the end of the list, its 'R' bit cleared, and the next page from the list is tried. Worst case: having to traverse the entire list again, only happens on page faults. The behavior is LRU-like.



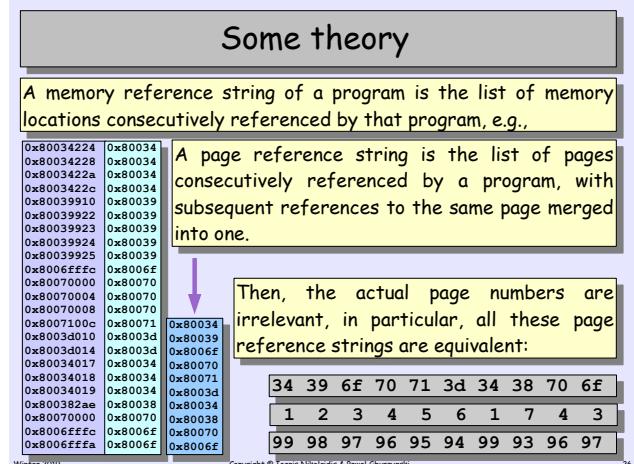
The page to be victimized (by FIFO) gets a second chance to prove its relevance. This is different from pure FIFO as the second chance will not victimize permanently needed pages.

## PERFORMANCE

Second chance is a very good approximation of LRU, especially for its cheap price. In some variation it is the most popular policy being used in actual OSes.

However, there are other components involved. With multi-programming, it is also important how to divide the available physical frames among multiple processes.

## SOME THEORY



But we can make it reasonable.  
Can we Eliminate page faults? NO.

# MONOTONIC STRATEGIES

Suppose that the system has admitted a program for execution. It says: you run with a local replacement strategy. Here is your real memory,  $k$  frames total. This is how much I can afford you at the moment. Needless to say, I don't care how grandiose ideas you may have regarding the size of your **virtual memory**.

The program runs and generates too many faults for the OS to be comfortable with.

So it says: OK, I have managed to scramble a few more frames. So now I will give you a bit more. Hopefully this makes your life easier.

It is natural to assume that more physical memory (for the same program) will translate into fewer page faults.

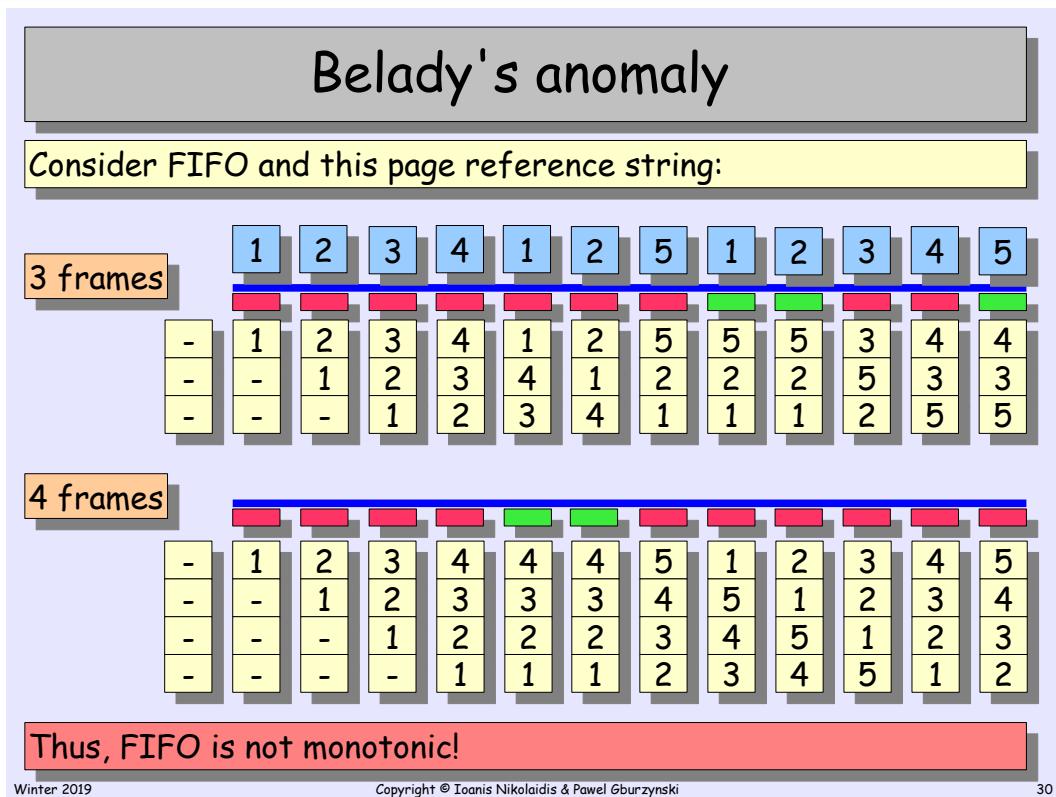
## DEFINITION

A strategy is **monotonic**, if no program i.e. page reference string (under a local variant of the strategy) will **ever** generate more faults when run in more physical memory.

An example is LRU.

This is because the set of  $k+m$  most recently used pages includes  $k$  most recently used pages.

A program running in  $k+m$  frames cannot generate more faults than when run in  $k$  frames: **at any time, it has the same frame as before + something extra**.



As you can see FIFO is not monotonic as there are more faults with 4 frames (10 faults) than with 3 frames (9 faults).

## WHAT CAN A PAGE FAULT MEAN?

1. An error (an actual illegal memory reference)
2. A CRAPE trick that only requires some change of flags, etc.
3. Request for a brand new page, e.g., the stack has grown.
4. Request for a static page (code) that must be fetched from the executable file (or a mmap'ped file).
5. Request for a data page that was once swapped out and is now demanded back.

What caused it.

1. This is a programming error, fix your damn program.
2. No need to do any page transfer at all -> no blocking
3. A frame may have to be evicted. This may require sending it back to the swap area -> possible blocking.
4. Same as 3 + A page has to be brought in -> mandatory blocking. Possibly one more I/O transfer to send the evicted page pack to the swap area.
5. Same as 4

It's easy to determine error 0, the OS looks at the page table and see if the process owns that page. error 3 and 4 are what we are going to look at.

## HOW MUCH MEMORY?

### How much memory?

Most programs need little physical memory to run efficiently, compared to their total memory accessed over the entire lifetime.

The phenomenon captured in curves like this is called **locality**.

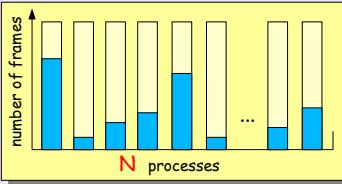
The sharp upturn at the vertical axis is called **thrashing**.

Thrashing means that a program's performance does not degrade gracefully as the amount of physical memory available to it is reduced.

In other words, the partitioning of memory among multiple programs in execution is not a trivial issue. The problem is that different programs exhibit different degrees of locality.

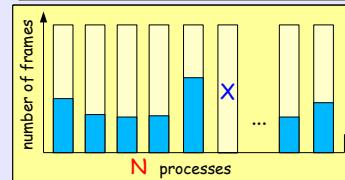
# OPTIMIZATION PROBLEM

## The optimization problem



The system would like to fill its entire physical memory with processes (as much as possible, anyway).

## The optimization problem



The system would like to fill its entire physical memory with processes (as much as possible, anyway).

**X** == complete swap-out.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

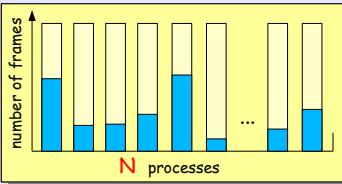
34

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

37

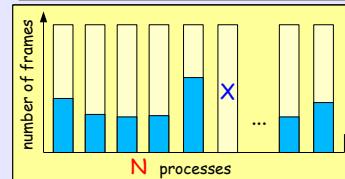
processes may need to be swapped out in order to reduce the page fault rate.

## The optimization problem



The system would like to fill its entire physical memory with processes (as much as possible, anyway).

## The optimization problem



The system would like to fill its entire physical memory with processes (as much as possible, anyway).

What are the objectives?

Winter 2019

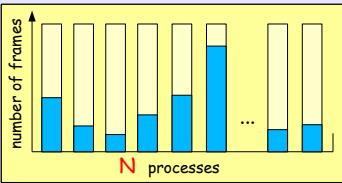
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

35

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

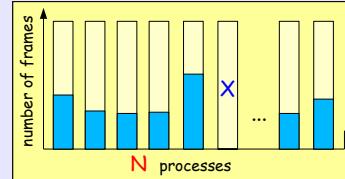
38

## The optimization problem



The system would like to fill its entire physical memory with processes (as much as possible, anyway).

## The optimization problem



The system would like to fill its entire physical memory with processes (as much as possible, anyway).

What are the objectives?

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

36

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

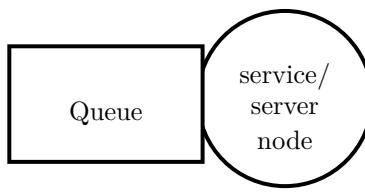
40

The goal is to maximize all three resources, CPU, Memory, Paging Device. You want to maximize the paging device as you want to guarantee to maximize the memory. Basically if you under utilize the paging device means you don't have enough processes.



In other words "fairness and neatness": all programs generate page faults close to the same "optimum" rate  $f_{opt}$ .

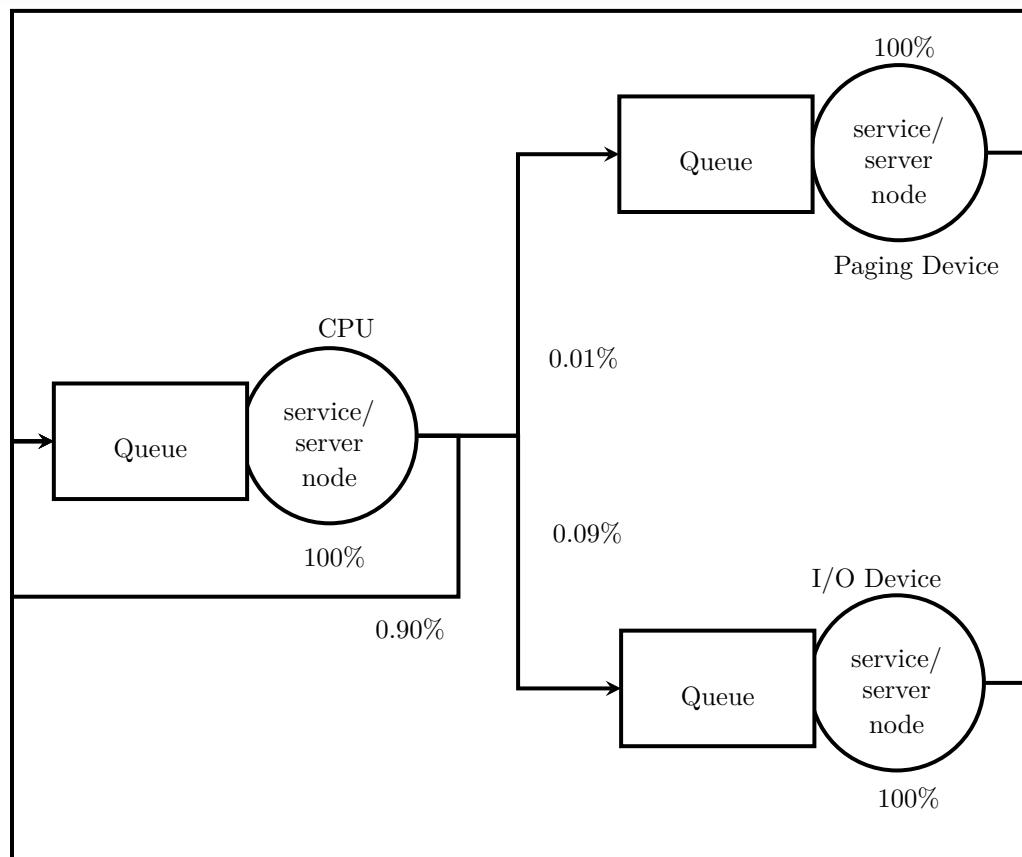
# QUEUING CONVENTION



You have a queue and a server that is processing the queue. Jobs or in this case processes are pushed to the back of the queue and the server processes the process at the head of the queue. Once the job is done it is either put back to the same queue, goto another queue, or is removed from the system.

## NOTATION

Server: Takes time to complete -> server. Example, the CPU is a server as it is serving or executing something.



Here is a more comprehensive example.

## GLOBAL VS LOCAL THRASHING

### GLOBAL THRASHING

Is where you have too many programs running and the system cannot handle it. The way to fix this is by suspending and swapping out programs.

## OPTIMUM PAGE FAULT FREQUENCY

The paging device has some "average" processing time of a page fault, say  $t_f$ . Note that different faults may have different actual processing times. The predominant component is the time to carry out a disk I/O transfer. Let  $f$  denote the observed page fault frequency. Note that it is physically impossible for  $f$  to be steadily greater than  $\frac{1}{t_f}$ . This is because the number of processes and the queue to the paging device are finite.

The goal of the paging device is to always be busy, i.e., being utilized, not being idle. We want  $f_{opt}$  to be close to  $\frac{1}{t_f}$ . Now,  $t_f$  is essentially equal to the cost of a disk transfer, thus:

$$f_{opt} \approx \frac{1}{t_{io} \times (1 + \alpha)}$$

$t_{io}$  is the average time to perform a disk transfer

$\alpha$  is the dirty factor

## SO HOW MUCH MEMORY?

$$\Delta = \frac{1}{f_{opt}}$$

$\Delta$  is the number of memory references that a program should be allowed to carry out successfully before it hits a page fault. This gives us a guideline for memory allocation. i.e.,  $\Delta$  is the number of frames that a program should have. Of course this is all wishful thinking as the future is unpredictable, but it's a target nonetheless. Of course we can also look at the program's history as well.

## EXAMPLE

**The working set model**

The working set of a given program at time  $t$ , denoted  $W(t, \Delta)$ , is defined as the collection of pages referenced by the program within the last  $\Delta$  memory references preceding  $t$ .

Assume  $\Delta = 10$

1 2 5 6 7

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

Time is the number of memory references. The blue box holds what memory I've referenced in the past 10 references.

**The working set model**

The working set of a given program at time  $t$ , denoted  $W(t, \Delta)$ , is defined as the collection of pages referenced by the program within the last  $\Delta$  memory references preceding  $t$ .

Assume  $\Delta = 10$

1 5 6 7

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

shift time forward and update the list

**The working set model**

The working set of a given program at time  $t$ , denoted  $W(t, \Delta)$ , is defined as the collection of pages referenced by the program within the last  $\Delta$  memory references preceding  $t$ .

Assume  $\Delta = 10$

1 2 5 6 7

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

**The working set model**

The working set of a given program at time  $t$ , denoted  $W(t, \Delta)$ , is defined as the collection of pages referenced by the program within the last  $\Delta$  memory references preceding  $t$ .

Assume  $\Delta = 10$

1 2 3 5 6 7

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

47

**The working set model**

The working set of a given program at time  $t$ , denoted  $W(t, \Delta)$ , is defined as the collection of pages referenced by the program within the last  $\Delta$  memory references preceding  $t$ .

Assume  $\Delta = 10$

1 2 3 4 5 6 7

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

48

**The working set model**

The working set of a given program at time  $t$ , denoted  $W(t, \Delta)$ , is defined as the collection of pages referenced by the program within the last  $\Delta$  memory references preceding  $t$ .

Assume  $\Delta = 10$

1 2 3 4

Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski

49

**10: VIRTUAL MEMORY**

87

## The working set model

The working set of a given program at time  $t$ , denoted  $W(t, \Delta)$ , is defined as the collection of pages referenced by the program within the last  $\Delta$  memory references preceding  $t$ .

Assume  $\Delta = 10$

2 6 1 5 7 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 4

2 3 4

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

50

## The working set model

The working set of a given program at time  $t$ , denoted  $W(t, \Delta)$ , is defined as the collection of pages referenced by the program within the last  $\Delta$  memory references preceding  $t$ .

Assume  $\Delta = 10$

2 6 1 5 7 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 4

2 3 4

Of course,  $\Delta$  is always much higher than 10. For example, a typical average disk access/transfer time is, say, 20 msec, so let us suppose that  $\Delta$  is 35 msec. This may translate into well over 1M memory references on a contemporary CPU.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

51

## A SIMPLE SOLUTION

Assume that time is measured in memory references (and nothing more).

Let

- $N$  be the number of processes in execution
- $M$  be the total amount of physical memory available to processes
- $|W_i(t, \delta)|$  denote the size of the working set of process  $P_i$  at time  $t$

### TOTAL DEMAND FOR FRAMES AT TIME $t$

$$D_t = \sum_{i=0}^{N-1} |W_i(t, \delta)| \quad (3)$$

$D_t < M$  the system may be able to admit a new program for execution.

$D_t > M$  the system cannot fulfill the demands of the

present programs and may have to consider a complete swap-out of some.

Unfortunately there is no way for the system to know exactly the working set size of a process. Even if the system could magically know that size, it wouldn't make sense to reconsider the program admission or swapping-out after every memory-reference.

Thus, this is all just theory. But it gives us ideas and hints for a practical implementation. For example, the system can monitor the page fault frequency and respond to its fluctuations.

$f > f_{opt} + \epsilon$  add more memory to the program

$f < f_{opt} - \epsilon$  take some memory away from the program.

In any case, it makes no sense to be jumpy with drastic decisions.

## THE TRIPLE POINT GOAL

The goal of is to have 100 percent utilization of the CPU, paging device, and memory. That is you want all three components of the system to be doing something and not idling. If you system cannot reach this goal then:

- you may not have enough jobs to run
- your system may be misconfigured

In general it doesn't make to buy:

- more memory than you need
- a faster CPU than you need
- A faster (and larger) disk than you need

# HOW IT ALL LOOKS LIKE IN REAL LIFE

The theory is OK, but its practical embodiment tend to be somewhat confusing to the uninitiated. Here are the reasons:

- Pages from processes' address spaces are not the only type of stuff stored in memory (file cache is another important thing)
- It makes sense to separate reclaiming free memory from servicing page faults.
- Creative strategies of optimizing disk i/o operations (e.g., read ahead) add finesse to our options.

## RECOMMENDATIONS

- Follow a strategy approximating LRU. Second chance gives you a generic hint: before you victimize, make sure to have looked at the 'R' bit.
- Detect thrashing the easy way (monitoring page fault frequency, looking at the utilization of the paging device). Swap when thrashing!
- Don't worry about memory partitioning. It will take care of itself, if you treat all processes as a single (global) system.

We shall start from an overview of some simple generic (but realistic) approach, which corresponds roughly to old BSD UNIX, and then comment on a few more modern refinements.

## EXAMPLE

### An example

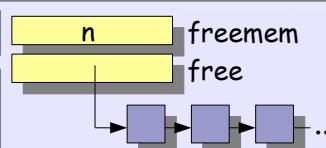
The system maintains a list of free pages:

The page fault service function doesn't execute any replacement policy:

```
page *pagein () {
    while (1) {
        mask ();
        if (free != NULL)
            break;
        waitforevent (&freemem);
        unmask ();
        schedule ();
    }
    p = free;
    free = p->next;
    freemem--;
    unmask ();
    return p;
}
```

The system tries hard to avoid a situation when `free == NULL`. This is not supposed to happen (although it is not impossible).

The free list is maintained by a special system process called `page daemon`, which makes sure that the list never shrinks below a certain minimum.



## The page daemon

```
void pagedaemon () {
    while (1) {
        if (freemem >= lotsfree) {
            delay (some_standard_longish_interval);
            continue;
        }

        victimize_some_page ();

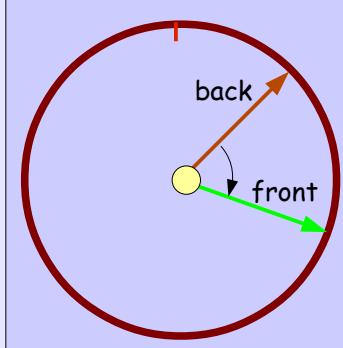
        if (freemem < desperate)
            try_to_swap_out ();

        delay (victimize_interval (freemem - desperate));
    }
}
```

Conceptually, the shortest `victimize_interval` (which gets shorter as `freemem` gets closer to `desperate`) corresponds to  $f_{opt}$ . If free memory continues to shrink at this interval, it means thrashing.

## VICTIMIZER

## The victimizer

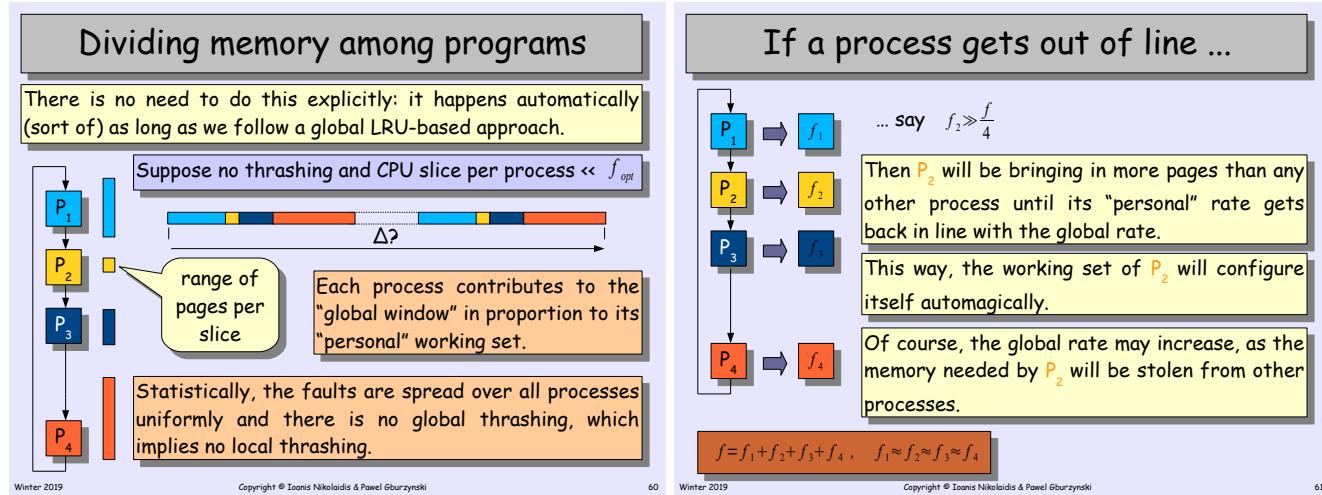


The `front hand` clears the "referenced" bits, and the `back hand` looks at them again. If the back hand finds a page with the referenced bit cleared, that page is victimized.

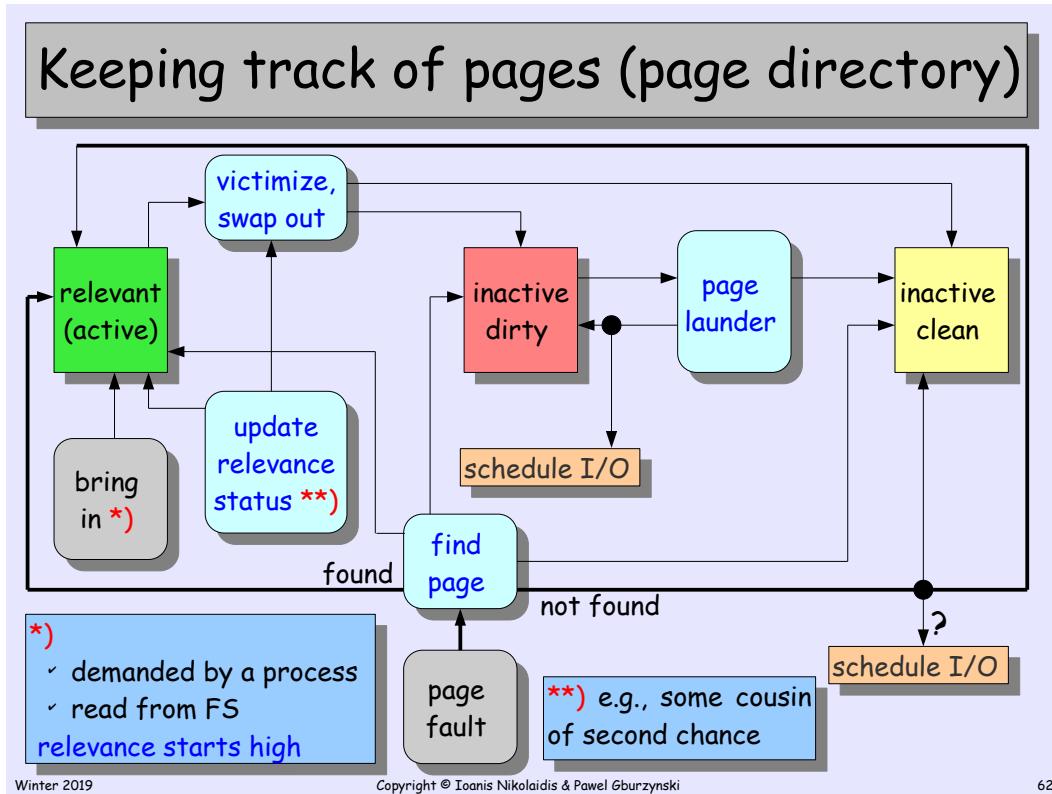
Originally, there was a single hand, and the strategy looked even closer to `second chance`.

Note that the victimizer only starts when `freemem < lotsfree`, which means that it may start in an old stale configuration. Thus, the system should refrain from victimizing any pages until the `back hand` looks at something that has been actually given a (second) chance.

# DIVIDING MEMORY AMONG PROGRAMS



## KEEPING TRACK OF PAGES (PAGE DIRECTORY)



Note that:

When a page is moved from Active to Inactive (dirty or clean), it doesn't have to be unmapped from its process's address space! It will be unmapped when actually "grabbed" (i.e., re-cycled). Migration to states Relevant and Dirty happens automatically upon reference.

For as long as the page is not actually re-used (i.e., filled with something different) the tag reflecting its current content is stored in the page directory. Thus, the page can be found when demanded and easily "un-victimized". The act of victimizing a page is never such a simple (one-step) operation as it would seem from the simplified models. In particular, the action of laundering a dirty page takes time during which things can happen.

# MODERN TRENDS

How do we deal with 64-bit address spaces? How to implement page tables for them?  
Let's do some calculation:

Let

$A$	be the number of address bits (e.g., 64)
$l$	be the number of bits in a page offset (e.g., 12)
$k$	be the number of levels
$u = (A - l)/k$	# of address bits per level
$2^u$	# of entries in a table
$E$	is the size of one entry (8 for a 64-bit address)
$E \times 2^u$	the size of a table is
$1 + 2 \times (k - 1)$	total # of all tables

## THE OVERHEAD

$$O = E \times (1 + 2 \times (k - 1)) \times 2^{A-l/k} \quad (4)$$

## Some examples

$A=64, \quad l=12, \quad E=8$	
$k$	overhead
1	$3.6 \times 10^{16}$
2	$1.6 \times 10^9$
4	470K
6	35K
8	10K
10	5K
12	3.5K

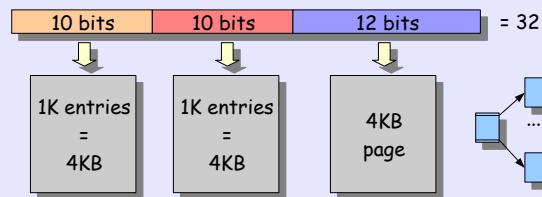
$A=32, \quad l=12, \quad E=4$	
$k$	overhead
1	$4.2 \times 10^6$
2	12K
4	896

$A=64, \quad l=14, \quad E=8$	
$k$	overhead
5	72K
4	320K
6	28K

Of course, not all possibilities are realistic, but the trade-off is clear. Now, would you really want to have 5-level tables?

## Consider x86



Overhead for the extremely hollow address space is 3 tables, i.e., 12 KB.

This is good for up to 4GB of logical address space and up to 4GB of physical memory. That used to be a lot, but isn't any more.

## GOING BEYOND 32-BITS

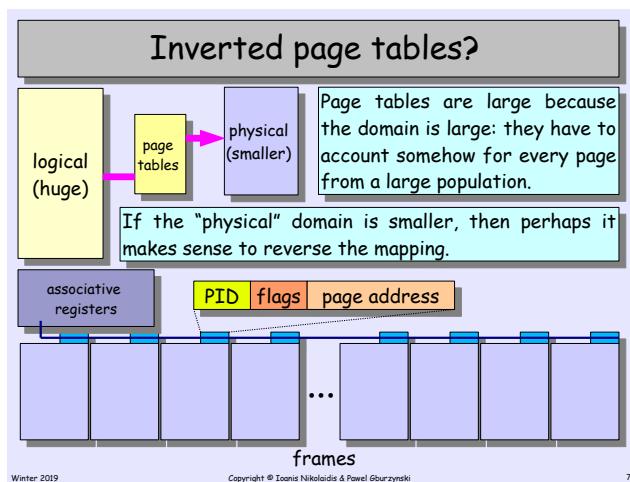
The physical memory limit is determined by the PDE/PTE (page table entry) format, and this is also the first problem to overcome. Here is the page table entry for a 32-bit system.

31	Page Frame Address	12      9 8 7 6 5 4 3 2 1 0	AVL 0 0   D A C T U W P
----	--------------------	-----------------------------	-------------------------

The lower 12 bits were used for attributes for the page table entry.

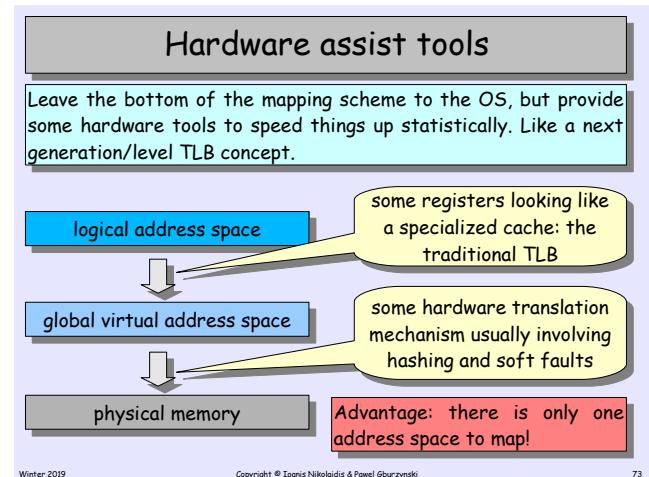
To expand the memory address range we can use a Page Directory pointer that points to the page directory we want to reference. From there it is just another multi-level page table. This however can get out of hand as you are increasing the level of re-directions that the computer needs to handle for each reference.

## INVERTED PAGE TABLES



The cost of forcing the hardware to go through five levels of tables may be comparable to the cost of doing something smarter entirely in software.

## HARDWARE ASSIST TOOLS



## WHAT ABOUT HAVING NO TABLES AT ALL

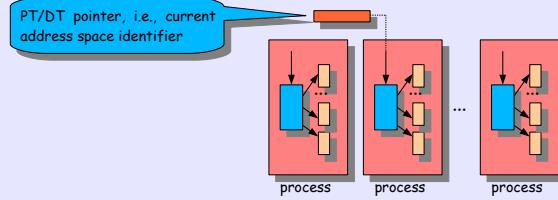
The worst problem with inverted page tables is that they make memory sharing difficult. Remember the TLB. It resolves the vast majority of all references without resorting to the tables at all. So this is the idea:

Make TLB misses perceptible by the kernel (as correctable faults), and let the system take care of filling the entries. The system can implement the mapping the way it pleases. Instead of hardware mapping we use software.

## PARADIGM SHIFT

### The paradigm shift

With the traditional approach, different sets of tables are used for each logical address space.



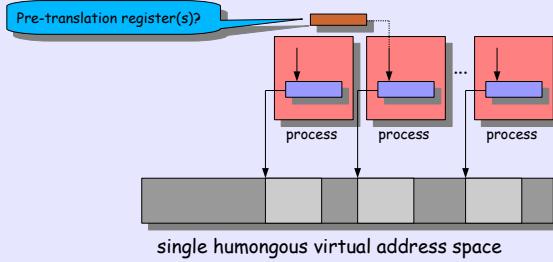
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

74

### But once we agree that ...

... the virtual address space is humongous (and the traditional approach doesn't work any more), we may as well do it (conceptually) like this:



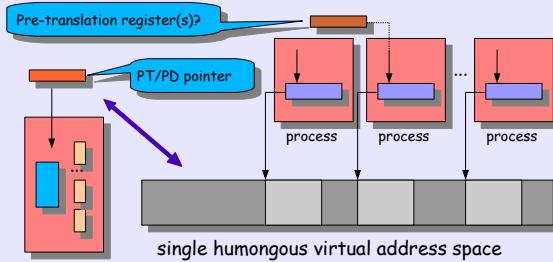
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

75

### But once we agree that ...

... the virtual address space is humongous (and the traditional approach doesn't work any more), we may as well do it (conceptually) like this:

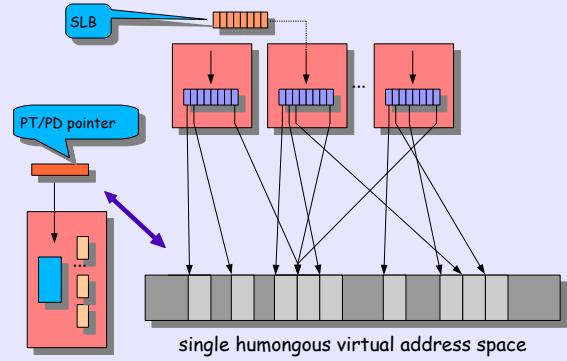


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

76

### The conceptual view



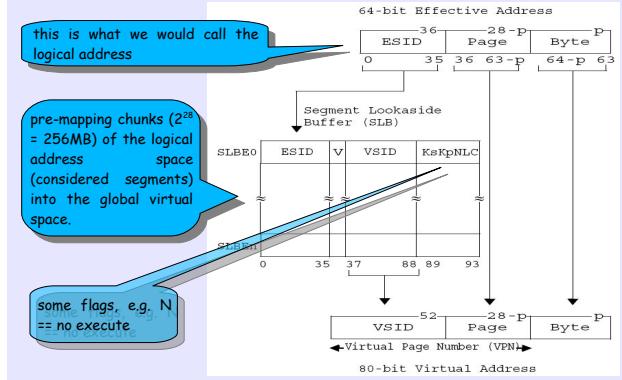
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

79

## POWER PC

### PowerPC address translation

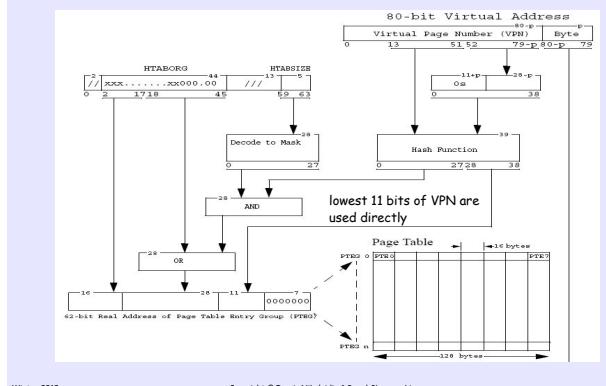


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

78

### Now for the virtual to physical part



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

80

## THE RULES

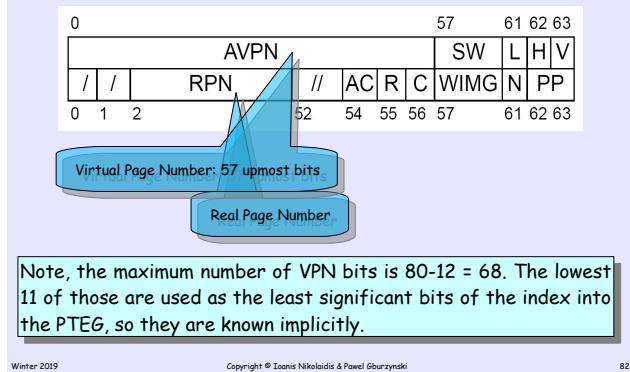
The hash function produces two PTEG (Page Table Entry Group) pointers. In fact, two slightly different hash functions are applied. They return pointers to the primary and secondary PTEGs.

If the page is found in the primary group, the frame address is taken from there, if not, the secondary group is scanned. If the page is not found there, we have a fault.

A single group contains descriptions of up to 8 pages that may hash to the same value.

Of course, there is a (traditional) TLB (applied after the SLB), which statistically circumvents all this mess. Note that the TLB is applied globally.

### PTE layout



Winter 2019 Copyright © Ioannis Nikolaidis & Paweł Gburzynski 82

**How much does it cost.** The recommended size of the hashed page table is half the total number of real pages to be accessed (i.e., that many PTEGs).

The penalty, if the table is too short for a 4 GB real memory and 4 KB page sizes, is about 1.6 percent of all memory is reserved for the PTEG which is pretty good, considering that the table is global

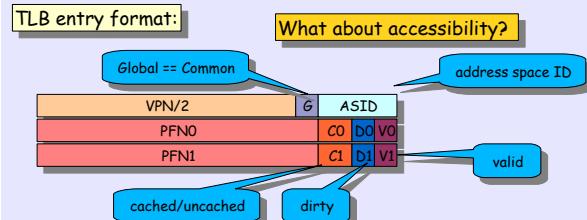
## MIPS

MIPS does not have a hardware DAT! TLB misses triggers interrupts (exceptions) that the OS must respond to.

### A drastic approach: the MIPS

There is no hardware DAT! TLB misses trigger interrupts (exceptions) that the OS must respond to.

TLB entry format:



Such entries are explicitly settable (and readable) with (privileged) machine instructions.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

85

## IA-64

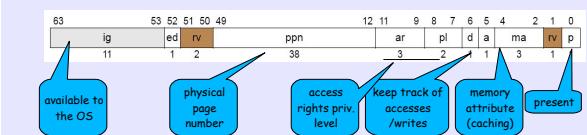
The DAT is optional on IA-64. The architecture allows the use of the VHPT walker.

### The VHPT walker

Two options: a per-process (single-region) straightforward linear page table, or a single global hashed table:

Region size is  $2^{51-12}$  (assuming 4KB pages) =  $2^{49}$  = nearly GADZILLION pages. Clearly, you cannot store a page table like this in memory ... unless ... it is stored in virtual memory.

PTE in short (linear) format:



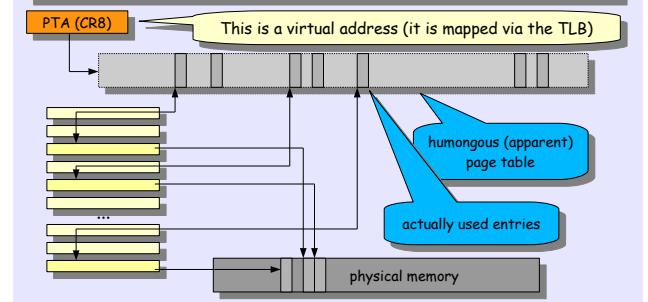
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

87

### So how does it work?

You put into the TLB, in addition to entries describing recently referenced actual pages, entries pointing to the recently relevant fragments of the page table:



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

88

## The hashed variant: PTE

offset	63	52 51 50 49	32 31	12 11	9	8	7	6	5	4	2	1	0
+0		ig	ed	rv	p	n	ar	pl	d	a	ma	rv	p
+8			rv		key		ps					rv	
+16	ti			tag									

This is a bit pattern which, in combination with the hash index, uniquely identifies the page (including its region).

In contrast to PPC, the walker examines only one entry (the one directly pointed to by the index). If the page doesn't match, it triggers a fault, and the OS will have to take over.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

89

## SUPERPAGES

### Superpages

Some contemporary architectures (including PPC and IA) allow you to have pages of variable size. This means that TLB entries can describe physical memory chunks of different sizes, e.g.,

PFRA SID LEN FLAGS

A large continuous range of physical memory can be described with much fewer entries. This means:

- reduced size of tables
- more TLB hits

But also more burden for the OS: promoting/demoting pages, fragmentation.



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

91

# 11: FILE SYSTEMS

A file is an identifiable collection of information, typically stored on a non-volatile medium for safe archival and retrieval.

## TAPE-BASED FILE SYSTEM

In the beginning each file is stored on a separate tape. Physical labels were used to label files.

## PROBLEMS

- **Poor utilization of tapes** Either every tape must be large enough to accommodate the largest conceivable file, or we have to deal with tapes of different sizes and force the program to specify the needed size when it creates a file.
- **Poor protection.** what if the operator mounts the wrong tape?
- **Files are sequential.** No direct access. Difficult to update file "in-place".
- **inconvenience.** You had to physically load each file by hand.

There were some protection methods like physical ring locks or **logical vsn label** on tape.

## BLOCKING

Practically all storage media require information to be blocked, but not always for the same reason. Tapes would have gaps between blocks so that the read head can stop before the start of the block so that the inertia doesn't go too far ahead. In this case a larger block leads to better tape utilization.

## DISKS

The evolution of tapes. Disks are also blocked. A single physical block (sector) is the minimum unit of information that can be transferred (and addressed).

## SURFACE MAPPING

P	the number of platters
2P	is the numbers of heads
C	the number of cylinders
T	the number of sectors per track
$S = 2PCT$	the number of sectors
$i$	the sector index of the disk from 0 to S-1
$c - i/(2PT)$	is the cylinder number
$u = i \bmod T$	is the sector number within the track
$h = (i/T) \bmod 2P$	is the head number

## HOPSCOTCH

0	1	2	3	4	5	6	7	8
0	5	1	6	2	7	3	8	4

Sometimes the read would be staggered as each block will undergo some error checking and correcting so the hopscotch will stagger the data such that by the time the error checking and correction is done for the current block, the next block the head is passing by would be the next block in sequence.

**Design issues in a file system.** The user-visible part: file id and location (hierarchy). operations on files, efficiency, reliability, resilience, integrity. Internals: organization of disk storage, caching, integration with virtual memory, reliability, resilience, integrity.



# REPRESENTING FILES

A disk can be envisioned as a simple array of sectors. Every sector is directly accessible, independently of other sectors. Files are sequences (or sets) of sectors described in some data structures. Ideally, those data structures should facilitate.

- fast location of subsequent sectors for sequential access
- fast location of any sectors for random access
- efficient file growth and shrinkage
- efficient usage of disk storage - no fragmentation
- reliability: low risk of serious damage from minor malfunctions

## CONTINUOUS ALLOCATION

You have a single contiguous segment for a file. This makes it easy to describe a file, have fast access, and allow fast access to both sequential and random ways.

However, the file size should be known in advance, it may be difficult to grow a file, and fragmentation. As naive as it appears, this approach is quite popular in some applications. A program may receive a disk partitions to implement its own "file" covering a consecutive range of sectors.

## DISK PARTITIONS

### Disk partitions

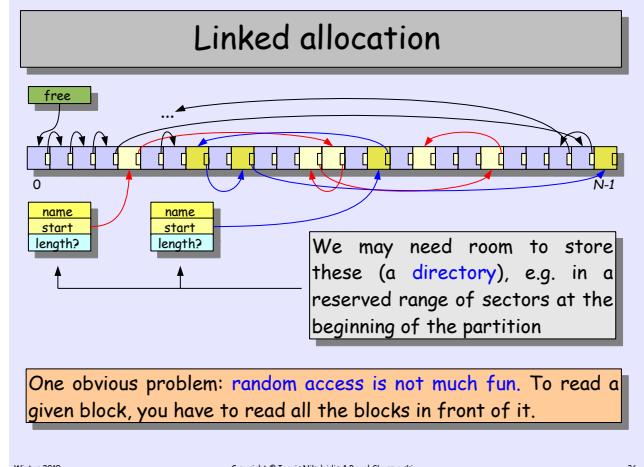
Typically, a disk is partitioned (sometimes trivially, i.e., into a single partition) before a file system is put on it.

You may want to have several file systems on the disk (possibly of different types).

Some programs (e.g., high-performance databases) may prefer to operate on raw partitions and implement their own filesystems.

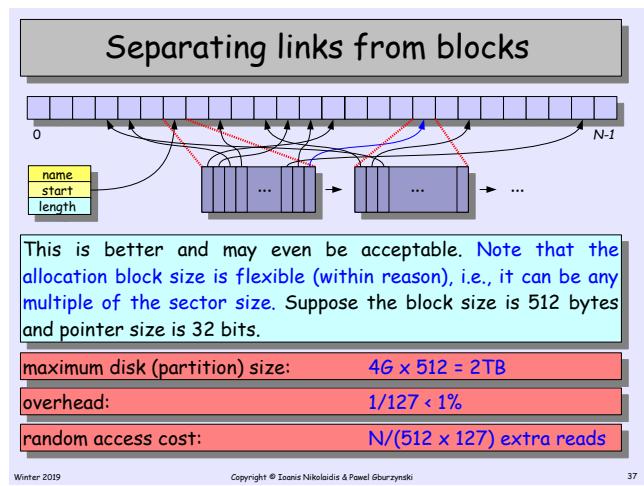
/sbin/fdisk /dev/hda						
Disk /dev/hda: 4311 MB, 4311982080 bytes						
15 heads, 63 sectors/track, 8912 cylinders						
Units = cylinders of 945 * 512 = 483840 bytes						
Device	Boot	Start	End	Blocks	Id	System
/dev/hda1	*	1	217	102501	83	Linux
/dev/hda2		218	7525	3453030	83	Linux
/dev/hda3		7526	8912	655357+	82	Linux swap

## LINKED ALLOCATION



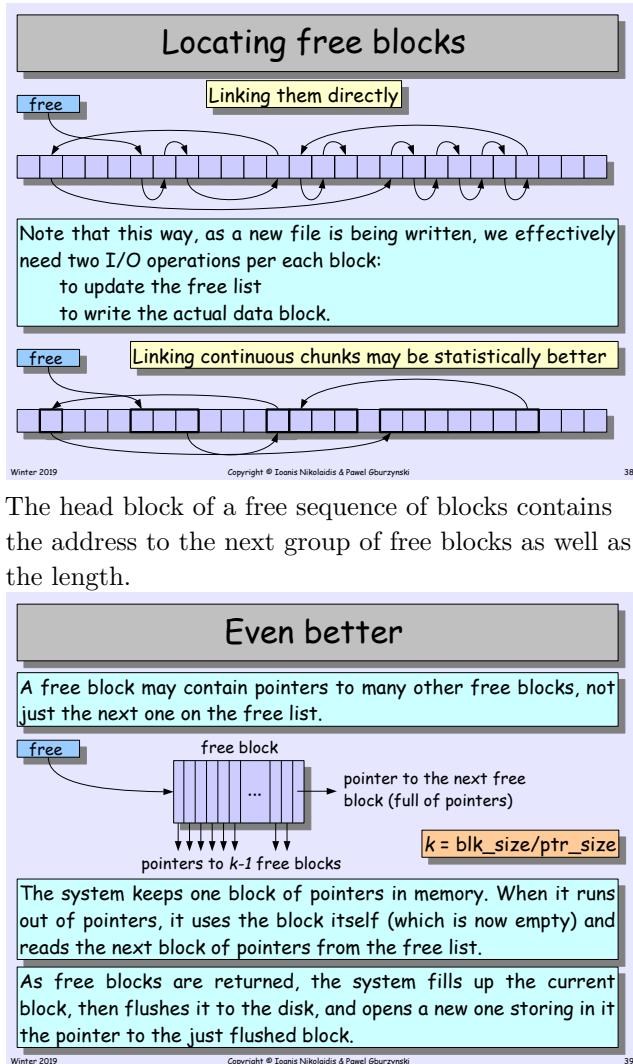
great way to solve fragmentation but has no direct or random access, has poor sequential access.

## SEPARATE LINKS FROM BLOCKS



Adding a level of redirection allows us to reduce fragmentation by having elements within the head of the file to contain offsets to blocks of data. Note: these blocks may not be in order but are called in order by the sequential calling of the page offset entries in the head page of the file. 512 is the sector size as well as the block size.

## LOCATING FREE BLOCKS

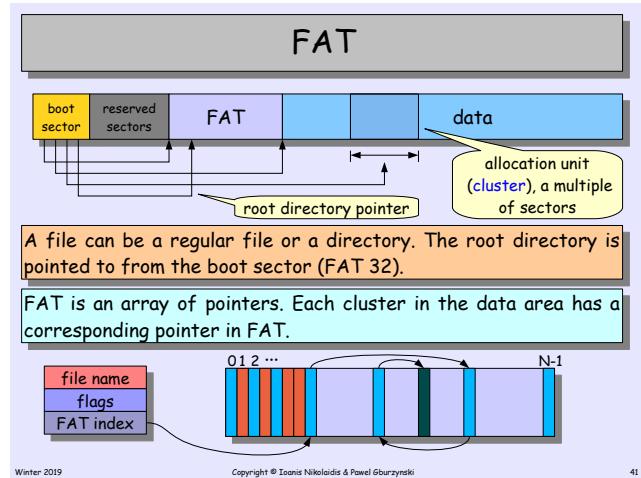


## BITMAP

Each allocatable block has one status bit in the bit map. If the status bit is 1, it means that the block is empty (or perhaps full). The bit maps must be eventually written to the disk, but the system can cache them in memory. Otherwise, a block acquisition/release would always require an extra I/O operation to update a bit in the map.

However, there are trade-offs: Block size is reduced which leads to less fragmentation, but more complexity and higher access time. And Sync intervals which leads to less overhead, but a higher likelihood of a disastrous damage.

## FAT



## WHAT IS THE BEST BLOCK SIZE

This statistic is outdated but in 1993 the average file size is 22k.

Most random files were smaller than 2k, pick a random byte and it probably in a file larger than 512k.

In general 89% of files take up 11% of the disk space, but 11% of files take up 89% of the disk space.

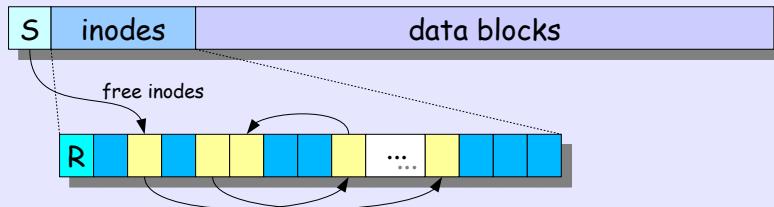
Today, the average file size is much bigger.

# UNIX FILE SYSTEM

## The "UNIX" file system

Let us start from the somewhat aged original. Later, we will comment on contemporary modifications and departures.

The first block of the partition is called the **superblock**. It contains the parameters of the file system, e.g., the actual block size and the boundary between **inodes** and **data blocks**.



The **inodes** part is an array of entries of fixed size. Each entry describes one file (or is empty).

Whenever a file is opened, its position is set to byte number zero. Whenever k bytes are read or written, the file position is advanced by k. Random access is implemented by making it possible to change the file position explicitly.

## FILE IDENTIFICATION

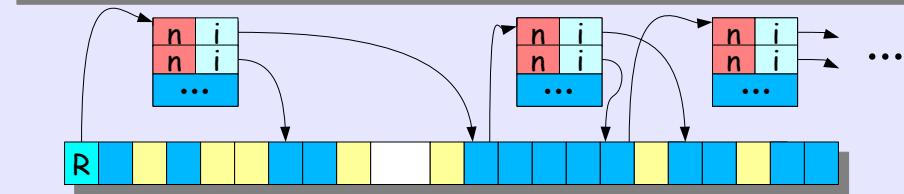
## File identification

Inodes directly represent files. If we didn't mind the inconvenience, we could identify files by inode numbers. The system in fact does it internally. And you also can, if you are root.

The familiar tree-like hierarchy is imposed on the inodes by the following two properties:

A file pointed to by an inode can be a **directory**. This is a special file that associates names with some inode numbers.

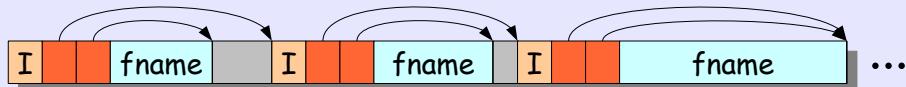
The **root directory** of the file system is stored in the file described by inode 0.



## DIRECTORY STRUCTURE

### Directory structure

A directory entry contains precisely two items of information: **file name + inode number**. The entries are linked:



Note that some **fragmentation** problems may arise as files are created, destroyed, and the directory entries are recycled.

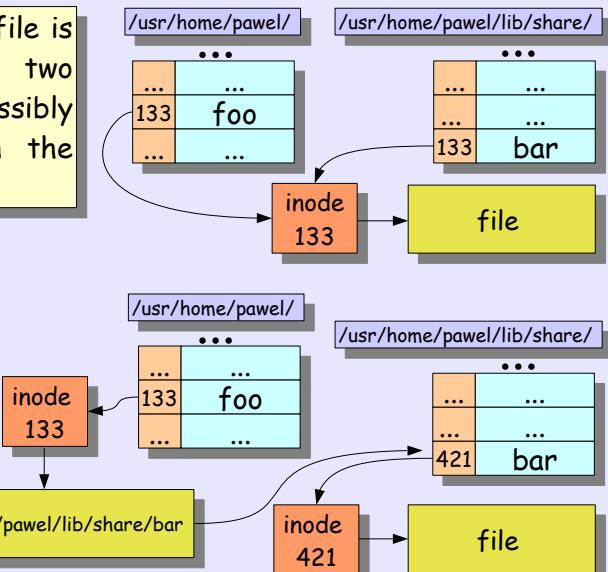
The important point is that a **directory entry contains no information about the file**. All that is kept in the inode. This brings about a few interesting features.

## LINKS

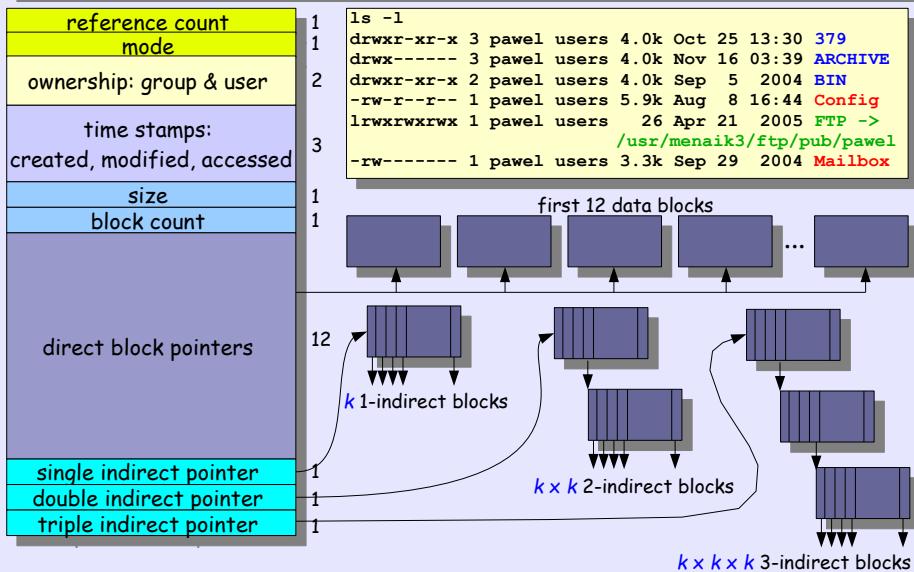
### Links

**Hard links:** the same file is accessible under two different names, possibly at distant places in the hierarchy.

**Soft links:** the link is stored in a (special) file. This is more flexible, e.g., a soft link can cross file systems.



## Inode layout + block list



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

52

We can have multiple indirect pointers blocks

## Maximum file size

Assume 512-byte blocks, with  $k = 128$ . We get:

$$F_{\max} = 12 \times 512 + 128 \times 512 + 128^2 \times 512 + 128^3 \times 512 > 1\text{GB}$$

This doesn't look like infinity any more. Note that by switching to 1K blocks (still decent), we multiply this number by  $2^4 = 16$ .

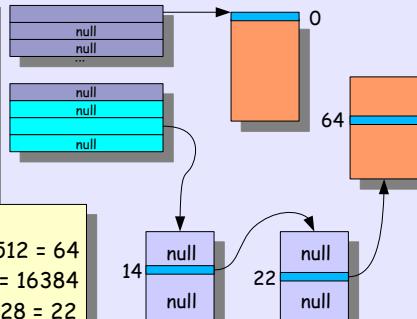
Hollow files are nicely accommodated. For example, consider this:

```

...
fd = open ("junk", O_CREAT |
           O_WRONLY | O_TRUNC, 0660);
write (fd, &byte, 1);
lseek (fd, 1000000, SEEK_SET);
write (fd, &byte, 1);
close (fd);
...

```

1000000/512 = 1953.... = 1954 blocks  
 $1954 - 12 = 1942$   
 $1942 - 128 = 1814$   
 $1814 / 128 = 14.17...$



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

53

## FREE SPACE MANAGEMENT

The original system used a free list where you have a list of fixed sized arrays of size  $k-1$  elements which are pointers to free blocks. The head of the list is stored in the superblock. When the file-system was mounted, the first free block was read into memory, providing immediate access to (up to)  $127 + 1$  free blocks.

## MODERN ENHANCEMENTS

---

Dual block size	Low overhead for large files, low fragmentation for small ones
Grouping	Splitting inodes into chunks and replicating the superblock
Bitmaps	instead of free lists. They are also grouped
Journaling	Neat tricks to practically eliminate the possibility of a serious damage when somebody puts the proverbial plug.

We really want small blocks to reduce fragmentation. However, larger blocks reduce the overhead on disk transfers, and thus improve the performance of heavily I/O bound programs like databases.

### DUAL BLOCK SIZE

1 block = k fragments, e.g., 4KB blocks, each consisting of four 1K fragments. Bit maps are used for representing block status. Individual bits are applied to fragments rather than blocks, for example.

blocks	0	1	2	3
fragments	0-3	4-7	8-11	12-15
blocks	0000	0011	1100	1111

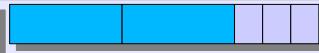
Suppose 1 means free. Only block 3 is free. Block 0 is used entirely. Blocks 1 and 2 are fragmented. They cannot be used as blocks, but their unused fragments are available for allocation to small files

## Trick.

### The trick

There is no difference between a large file and a small one (in its representation on the disk). Any file consists of a number of blocks + possibly a number of fragments (never more than  $k - 1$ ) at the end.

For example, assume that the file length is 11,000 bytes. The file takes two blocks and three fragments.



$$11,000 = 2 \times 4096 + 2 \times 1024 + 760$$

Winter 2019

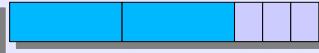
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

59

### The trick

There is no difference between a large file and a small one (in its representation on the disk). Any file consists of a number of blocks + possibly a number of fragments (never more than  $k - 1$ ) at the end.

For example, assume that the file length is 11,000 bytes. The file takes two blocks and three fragments.



$$11,000 = 2 \times 4096 + 2 \times 1024 + 760$$

Winter 2019

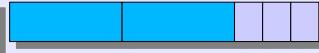
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

60

### The trick

There is no difference between a large file and a small one (in its representation on the disk). Any file consists of a number of blocks + possibly a number of fragments (never more than  $k - 1$ ) at the end.

For example, assume that the file length is 11,000 bytes. The file takes two blocks and three fragments.



$$11,000 = 2 \times 4096 + 2 \times 1024 + 760$$

Winter 2019

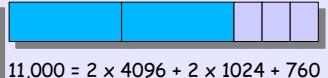
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

61

### The trick

There is no difference between a large file and a small one (in its representation on the disk). Any file consists of a number of blocks + possibly a number of fragments (never more than  $k - 1$ ) at the end.

For example, assume that the file length is 11,000 bytes. The file takes two blocks and three fragments.



$$11,000 = 2 \times 4096 + 2 \times 1024 + 760$$

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

62

### The trick

There is no difference between a large file and a small one (in its representation on the disk). Any file consists of a number of blocks + possibly a number of fragments (never more than  $k - 1$ ) at the end.

For example, assume that the file length is 11,000 bytes. The file takes two blocks and three fragments.



$$11,000 = 2 \times 4096 + 2 \times 1024 + 760$$

Winter 2019

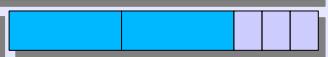
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

63

### The trick

There is no difference between a large file and a small one (in its representation on the disk). Any file consists of a number of blocks + possibly a number of fragments (never more than  $k - 1$ ) at the end.

For example, assume that the file length is 11,000 bytes. The file takes two blocks and three fragments.



$$11,000 = 2 \times 4096 + 2 \times 1024 + 760$$

The important fact is that ... just by looking at the file size we know which pointers point to blocks and which point to fragments.

Winter 2019

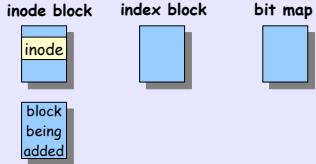
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

71

## JOURNALING

### Journaling

Suppose that a block is to be added to a file. These are the steps involved in the operation:



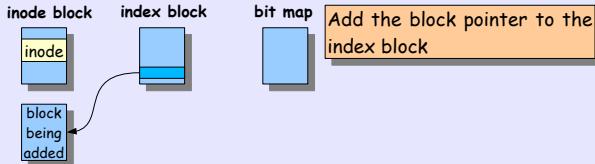
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

72

### Journaling

Suppose that a block is to be added to a file. These are the steps involved in the operation:



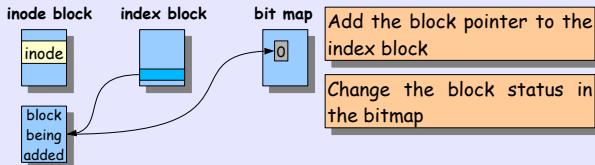
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

73

### Journaling

Suppose that a block is to be added to a file. These are the steps involved in the operation:



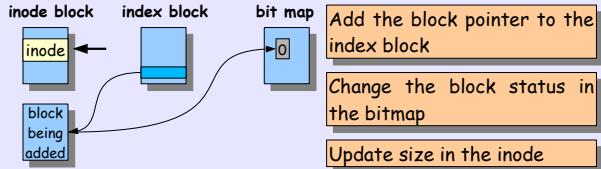
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

74

### Journaling

Suppose that a block is to be added to a file. These are the steps involved in the operation:



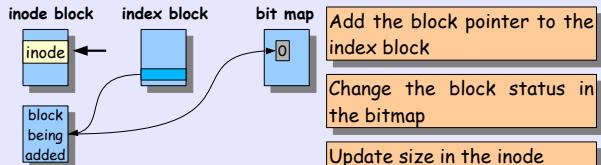
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

75

### Journaling

Suppose that a block is to be added to a file. These are the steps involved in the operation:



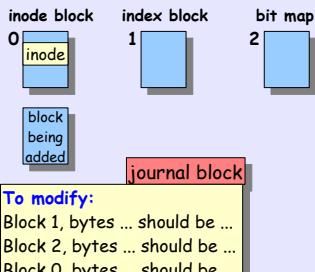
Three blocks must be written to disk to complete the operation, i.e., to make the disk structure fully consistent with the changes. If the system crashes after the first write but before the completion of the last one, we may be in trouble.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

76

### Similar to transaction roll-back in DB

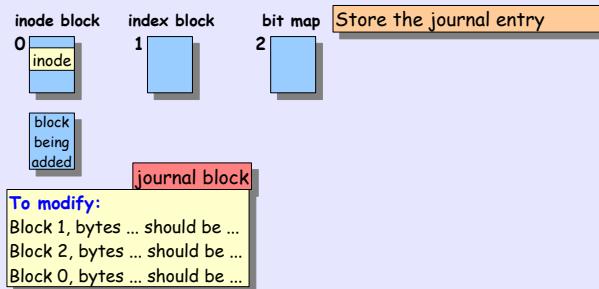


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

77

## Similar to transaction roll-back in DB

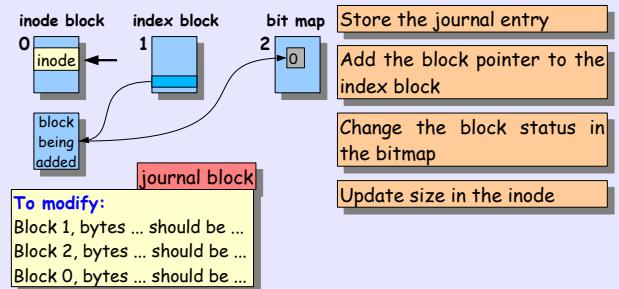


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

78

## Similar to transaction roll-back in DB

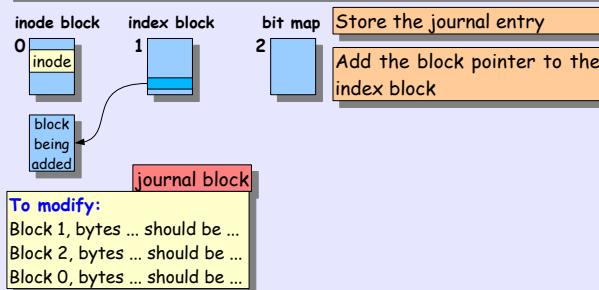


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

81

## Similar to transaction roll-back in DB

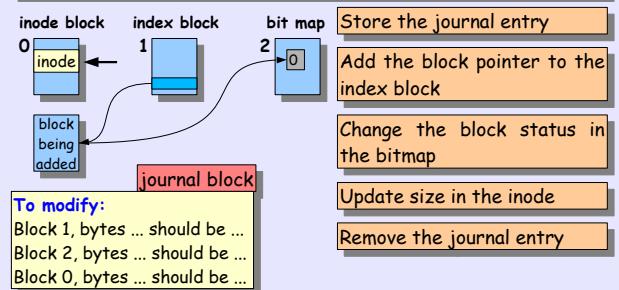


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

79

## Similar to transaction roll-back in DB

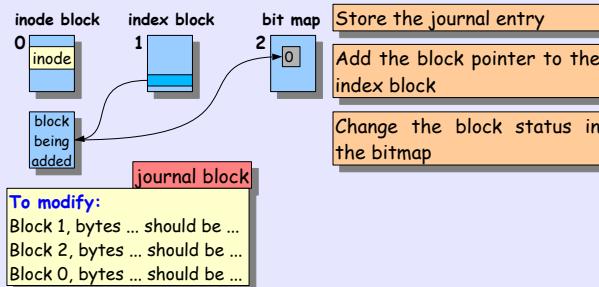


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

82

## Similar to transaction roll-back in DB

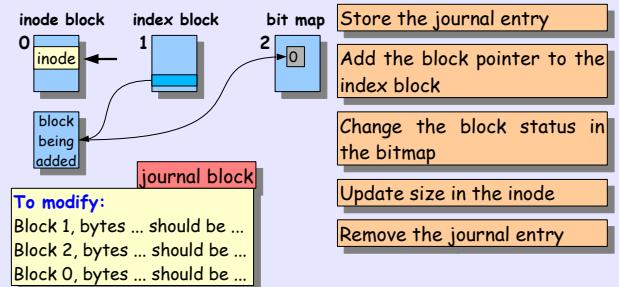


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

80

## Similar to transaction roll-back in DB



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

83

System consistency does not depend upon a success of multiple writes to the disk. After every single write, the system either is consistent, or is able to complete the unfinished transaction.

# EXTENTS

## Extents

I told you that the simple idea of continuous allocation is too valuable to be discarded as naive. Perhaps we shouldn't use it directly, but it may make sense to build files of continuous chunks (block ranges) called **extents**.

Statistically, even a large file may need just a few extents.

If long extents tend to be frequent, a file's description will tend to be short.

This means that access time to the file will tend to improve - both for sequential and random operation.

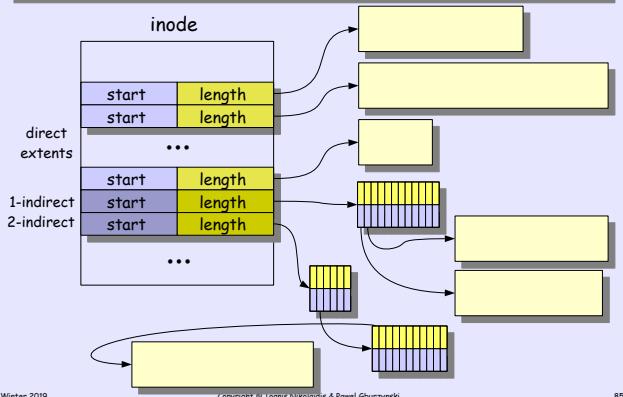
Also, describing free storage with extents may be more efficient (fewer things to keep track of, faster location of free space).

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

84

## For example



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

85

## Comments

Holes cause no problem: they are taken care of by **empty extents** with `start == NULL`.

One problem is the apparent complication of random access to distant blocks.

But we hope that the extents will tend to be large, so there will be few of them, so traversal of index extents will be infrequent.

Caching will also help, assuming that the "random references" are not completely random.

Besides, there are other (more efficient from the viewpoint of random access) ways of representing series of extents than the straightforward extrapolation of UNIX indexes and index blocks.

Winter 2019

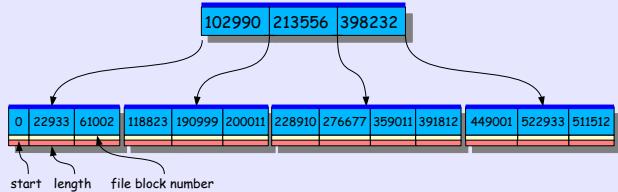
Copyright © Ioannis Nikolaidis & Paweł Gburzyński

86

# B-TREES

## B-Trees

You know what they are if you have taken a database course. It used to be a prerequisite for 379.



A file can be viewed as a database of extents with block numbers used as keys.

Winter 2019

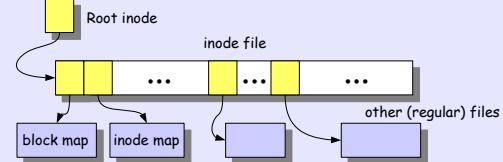
Copyright © Ioannis Nikolaidis & Paweł Gburzyński

87

# WAFL

## Other ideas: WAFL

**Write Anywhere File Layout**. All metadata is kept in files:



The entire file system forms a tree of blocks rooted at the single (Root) inode. This means that:

No blocks are dedicated to inodes, and there is no need to draw any lines.

An arbitrary number of filesystems can coexist on the same partition.

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

88

## All pointers in inode are the same level



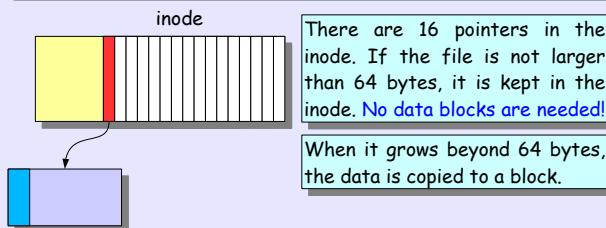
There are 16 pointers in the inode. If the file is not larger than 64 bytes, it is kept in the inode. **No data blocks are needed!**

Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

89

## All pointers in inode are the same level

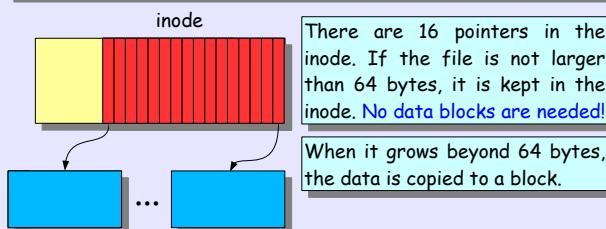


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

90

## All pointers in inode are the same level

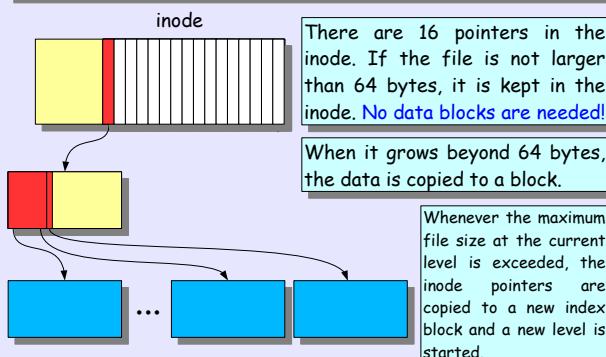


Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

91

## All pointers in inode are the same level



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzyński

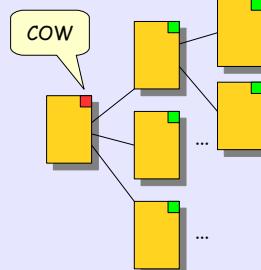
92

## SNAPSHOTS

### Snapshots

One neat feature of WAFL is the possibility to create quick "snapshots" of the entire filesystem. To create a snapshot, you:

1. duplicate the root inode
2. enter a COW mode



Winter 2019

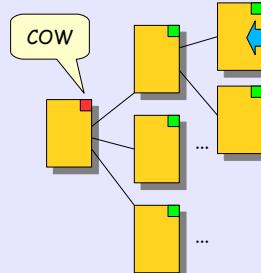
Copyright © Ioannis Nikolaidis & Paweł Gburzyński

93

### Snapshots

One neat feature of WAFL is the possibility to create quick "snapshots" of the entire filesystem. To create a snapshot, you:

1. duplicate the root inode
2. enter a COW mode



Winter 2019

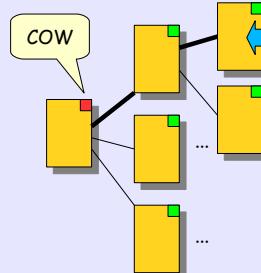
Copyright © Ioannis Nikolaidis & Paweł Gburzyński

94

### Snapshots

One neat feature of WAFL is the possibility to create quick "snapshots" of the entire filesystem. To create a snapshot, you:

1. duplicate the root inode
2. enter a COW mode



Winter 2019

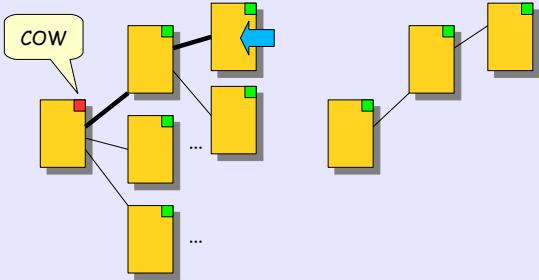
Copyright © Ioannis Nikolaidis & Paweł Gburzyński

95

## Snapshots

One neat feature of WAFL is the possibility to create quick "snapshots" of the entire filesystem. To create a snapshot, you:

1. duplicate the root inode
2. enter a COW mode



Winter 2019

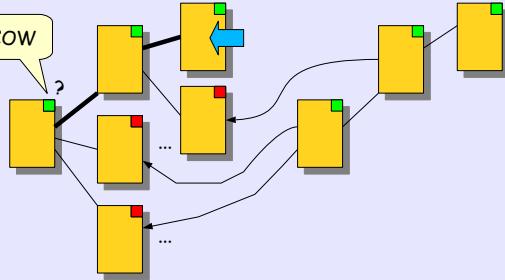
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

96

## Snapshots

One neat feature of WAFL is the possibility to create quick "snapshots" of the entire filesystem. To create a snapshot, you:

1. duplicate the root inode
2. enter a COW mode



Winter 2019

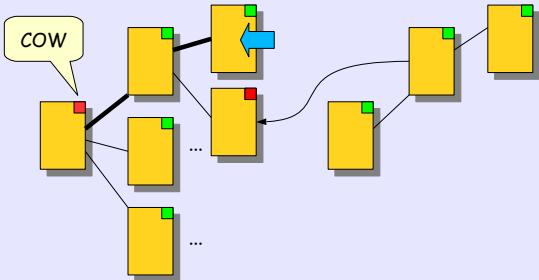
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

99

## Snapshots

One neat feature of WAFL is the possibility to create quick "snapshots" of the entire filesystem. To create a snapshot, you:

1. duplicate the root inode
2. enter a COW mode



Winter 2019

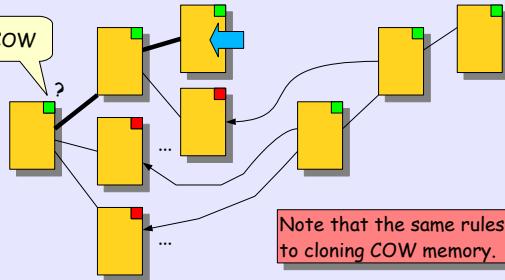
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

97

## Snapshots

One neat feature of WAFL is the possibility to create quick "snapshots" of the entire filesystem. To create a snapshot, you:

1. duplicate the root inode
2. enter a COW mode



Winter 2019

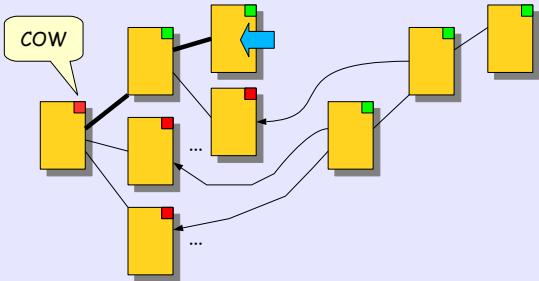
Copyright © Ioannis Nikolaidis & Paweł Gburzynski

100

## Snapshots

One neat feature of WAFL is the possibility to create quick "snapshots" of the entire filesystem. To create a snapshot, you:

1. duplicate the root inode
2. enter a COW mode



Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Gburzynski

98

## Snapshots

Note that the same rules apply to cloning COW memory.

# 12: DISK SCHEDULING

I/O request to disks may arrive from several places at pretty much the same time (in bursts), occasionally faster than they can be handled by the disk. This is because

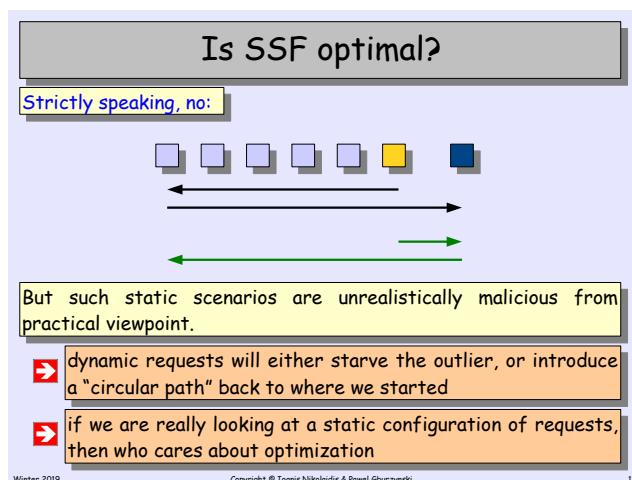
- **Multi-programming:** multiple programs doing I/O concurrently
- **Paging:** programs don't have to do explicit I/O to trigger I/O
- **caching:** I/O can happen independently of programs' intentions
- Moreover, various performance enhancing tricks like read-ahead contributes to this picture.

**Seek time** is the absolute time/distance it takes for the head to move from the current cylinder to the desired cylinder. The **total seek time** is the summation of these seek times.

## SSF: SHORTEST SEEK FIRST

Having completed processing a request, select the next request with the shortest seek distance from your current position. This is similar to SJF, and it looks like it would be optimal but its not. It still have starvation issues and weird edge cases where it is sub optimal.

### IS SSF OPTIMAL



Lets assume that the head starts reading from the yellow cylinder and moves leftward as those cylinders in the queue have the shortest SSF. However, after the most leftward cylinder finishes processing it then needs to go all the way right to process the next cylinder in the queue. If the device waited a little longer it would've seen the blue cylinder first, processed it then move leftward, decreasing the seek distance. However, this scenario are unrealistically malicious from a practical viewpoint.

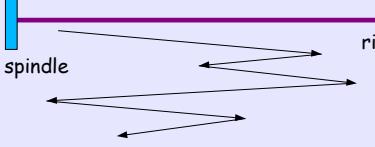
# LOOK STRATEGY

**LOOK (aka Elevator) Strategy**



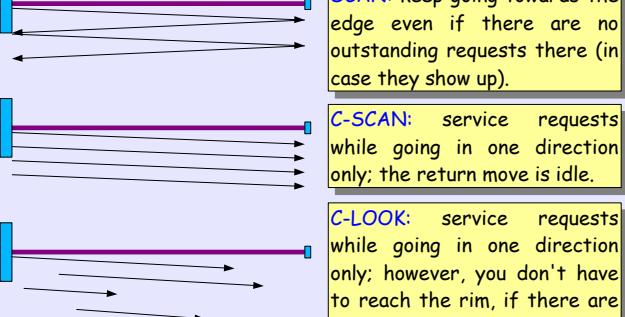
Once you start going in a given direction, **do not reverse** unless there are no more requests to be serviced along the way.

LOOK is the most popular and widely used strategy. It is not optimal (as SSF), but decent and reasonably fair. There are still some rudimentary fairness issues remaining with LOOK, which lead fairness nuts towards more fair strategies.



The middle range of cylinders is still favored. If you are located there, your waiting time will likely be shorter than if you are closer to an edge.

**Some (more fair) variations**



- SCAN:** keep going towards the edge even if there are no outstanding requests there (in case they show up).
- C-SCAN:** service requests while going in one direction only; the return move is idle.
- C-LOOK:** service requests while going in one direction only; however, you don't have to reach the rim, if there are no more requests towards it.

Winter 2019      Copyright © Ioannis Nikolaidis & Paweł Gburzynski      13      Winter 2019      Copyright © Ioannis Nikolaidis & Paweł Gburzynski      14

There are also other issues involved. Like optimizing other components of the disk access cost:

- The file system will try to allocate continuous ranges of blocks to files
- The hopscotch approach to numbering consecutive sectors/blocks
- Reading ahead.
- Merging multiple requests

Multiple disks in high performance systems:

- Disk arrays: for performance and reliability

# MULTIPLE DISKS

**Multiple disks, perhaps lots of them**



How to best take advantage of them?

What are the options for a single file system spanning them all?

**JBOD:** Just a Bunch Of Disks. View them all as a single large virtual disk (a virtual partition?).

Disk 0	Disk 1	Disk 2
0 1 2 3 4 5 6 7 0 1 2 3 4 5 6 7 0 1 2 3 4 5 6 7	0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23	0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23

This may be a better approach from the viewpoint of access time:

0 0 0 1 1 1 2 2 2 3 3 3 4 4 4 5 5 5 6 6 6 7 7 7
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23

Winter 2019      Copyright © Ioannis Nikolaidis & Paweł Gburzynski      19

# RAID

With JBOD, if a single disk fails, the whole file system becomes essentially useless. It's like if a hole suddenly popped up in your partition.

## RAID 1

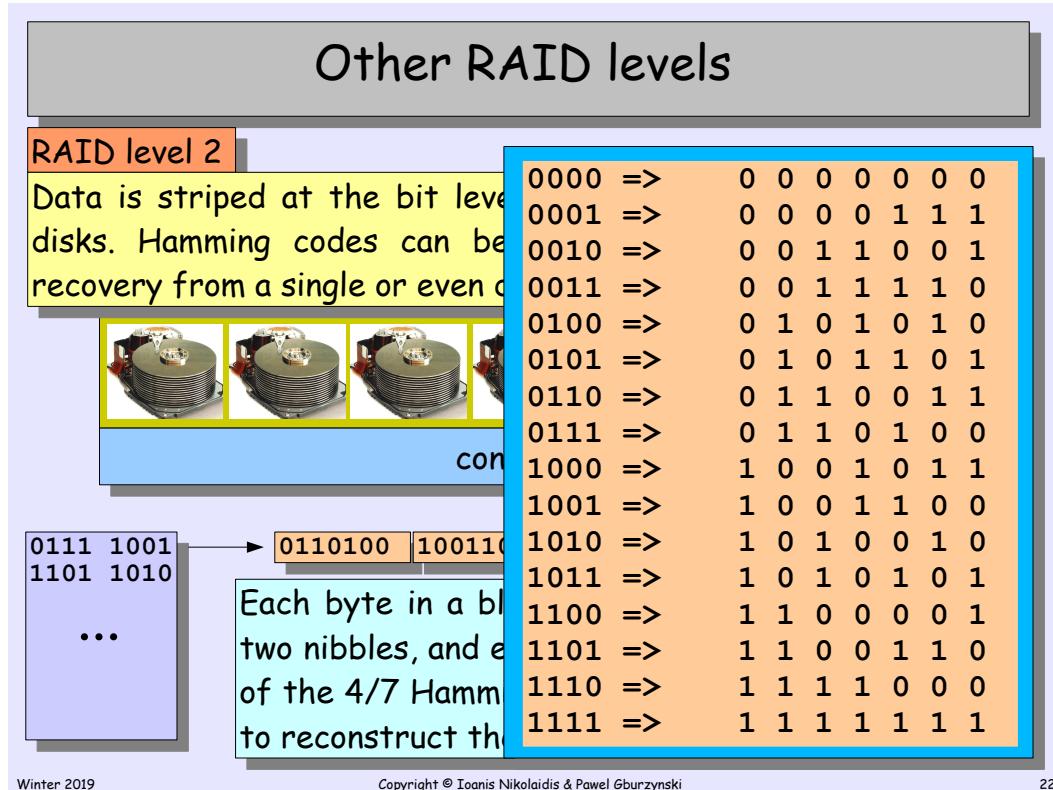
Creates a mirror on two or more disks. The multiple disks are for redundancy and reliability rather than larger size. If you have a lot of disks, you can group them into a JBOD and mirror them into a RAID 1 array.

Note that either set can be used for reading (alternately) to reduce average access time. Writing goes to the basic set, and then the mirror is updated.

## RAID 2

Data is striped at the bit level. The controller synchronizes all disks. Hamming codes can be used for error correction and recovery from a single or even dual disk failure.

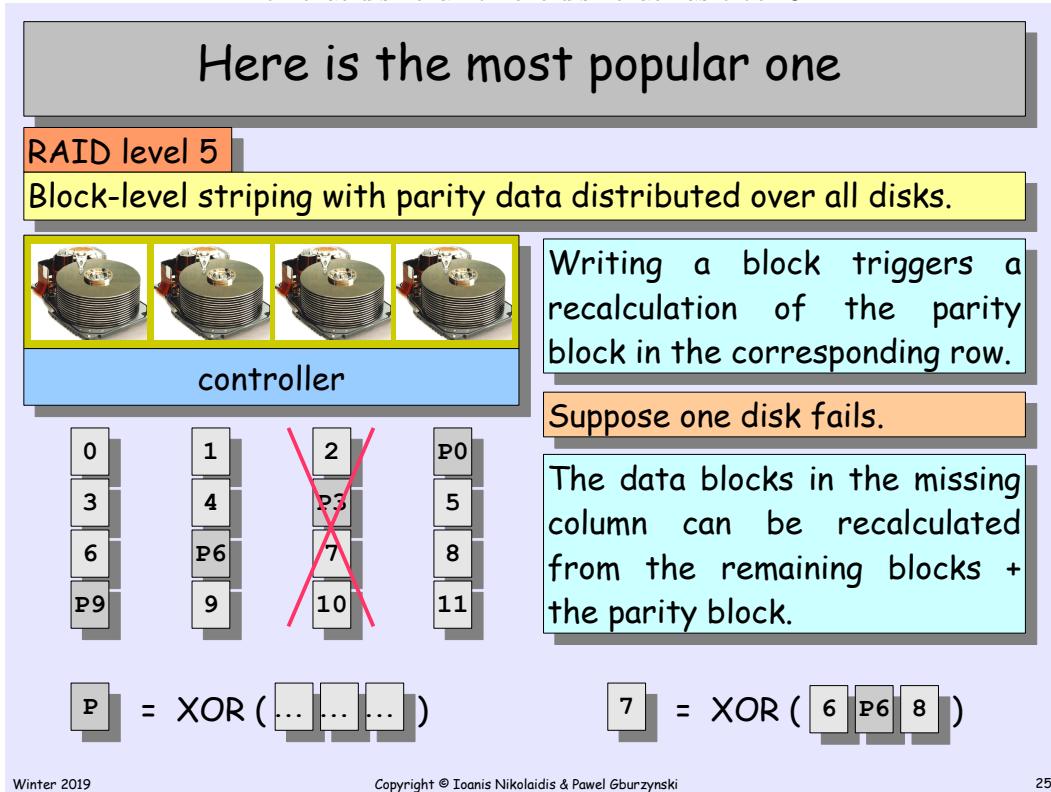
Each byte in a block visible by the FS is split into two half bytes, each half byte is transformed into the 7 bits of the 4/7 Hamming code. This allows the controller to reconstruct the half byte from any 5 of the 7 disks.



## RAID 5

Block-level striping with a parity data distributed over all disks. When a disk fails, the entire array can be reconstructed using the parity data and XORing the blocks. When you are re-silvering the array however, there is a performance hit as not only you are reading parity data and doing reconstruction operations but any data that is written onto the disk will force the parity to change. Causing a slowdown in the reconstruction process.

Moreover, if the data that is being written is block 1 and block 8 in the figure below then it would take much longer than if the data was written on block 0 and 4 as you are just updating the parity on the disk an nothing else. If block 1 and 8 were updated then you are updating block 1 and P8 which involves more write operations on that disk than on the disk that has block 8.



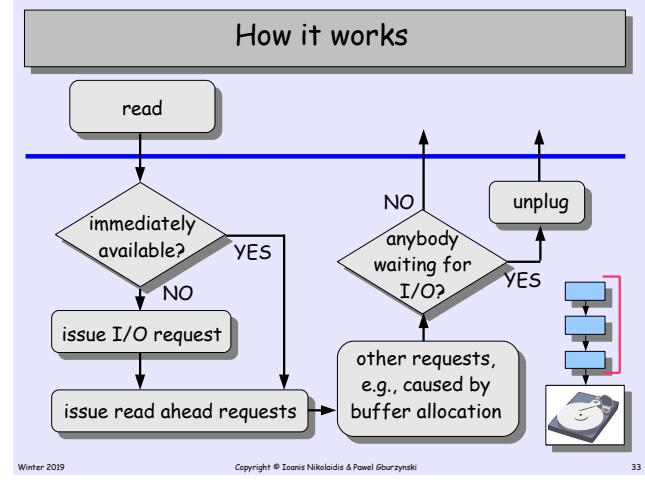
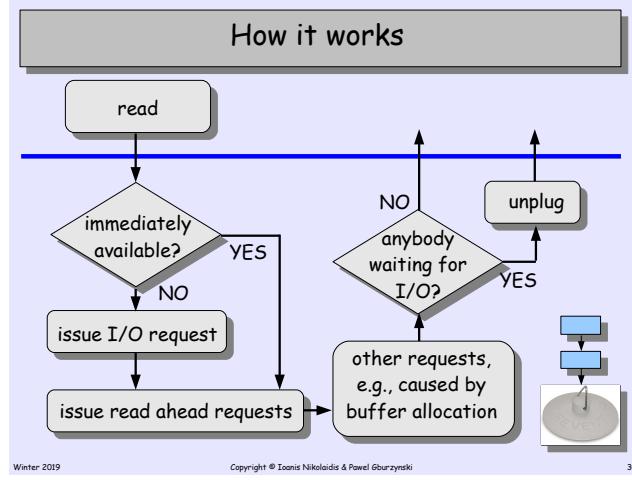
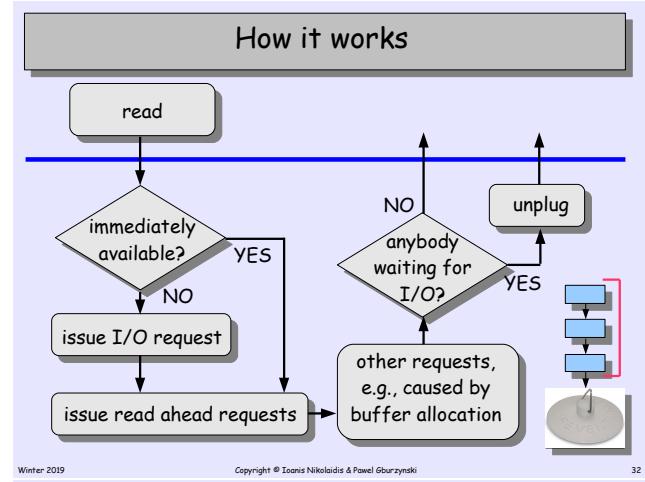
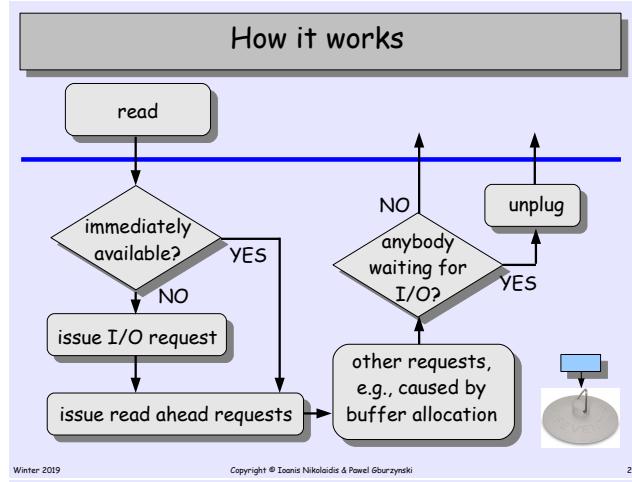
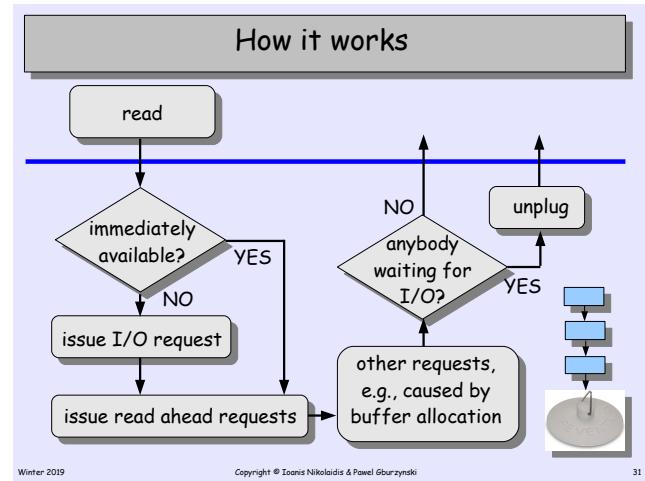
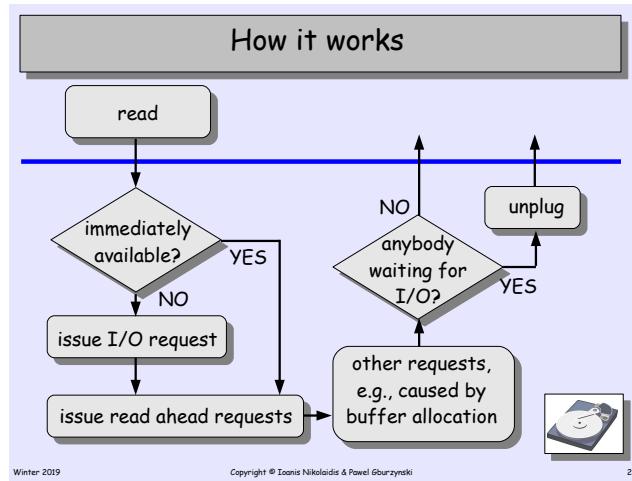
## READING AHEAD (FOR FILE OPERATIONS)

When a file is opened, the system tentatively sets the read ahead option for the file and assumes some window size specifying how much extra stuff should be read ahead after each actual read. If the subsequent read operation falls within the window, the system assumes that the read ahead option makes sense and it should remain on. Otherwise, the system assumes that the file is being read at random, and the read-ahead option is switched off. Unfortunately we can't have write-ahead.

# PLUGGING

Whenever a new I/O request is added to the disk queue, the driver tries to merge it with another request already queued there. What you are able to do depends on the capabilities of the controller. Sometimes it makes sense to hold a request at the front of the queue waiting for merging opportunities

## HOW IT WORKS



# 13: VIRTUAL MACHINES

Virtualization is about independence from mundane aspects of reality: A program language makes programs independent of hardware, but still depend on libraries and syscalls.

You can have multiple OSes running on a machine either by dual booting or by running within a host OS.

## SOFT VIRTUALIZATION

Java is a great example of a soft virtualization as you convert the source code from Java to its bytecode that can be executed on a JVM. The Java bytecode does not need to be recompiled again as the Java byte code is system independent from the OS and machine. This is all handled by the JVM.

In fact interpreters are usually quite portable so it is easy to "rapidly" deploy the VM on different HW/OS platform.

Of course you can compile the Java byte code into machine code that can be executed naively on a machine for optimal performance.

## CYGWIN

Cygwin is a UNIX/Posix function and wrapper for system calls. Its not a perfect map as there are some syscalls that does behaves differently or does not exist at all.

### NOTE

That the UNIX executables (compiled for the same ISA) cannot run under Cygwin (The program must be recompiled)

## OS VIRTUAL MACHINES

The guest OS goes through a hardware emulator that goes through the virtual machine monitor (a meta OS) that talks to the hardware.

### NOTE

There are different types of virtualization architecture.

OS virtualization is useful as we can:

- Better utilize expensive resources more effectively, especially for multi-user systems.
- Better development as you don't need to hog the system and underutilize it
- Multiple OS as different OS have different syscalls and behaviors.
- Server consolidation and service isolation (increased security or sand-boxing)

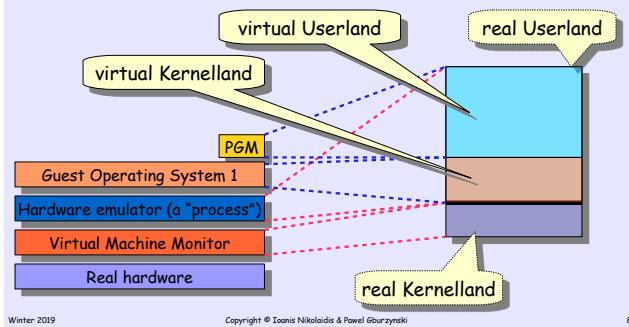
Virtualization is not interpretation in this case as we want the guest OS to execute everything directly as much as possible with a few "sensitive" events (instructions) going towards the host OS.

Thus, the virtual computer set up by the VM monitor executes as a "process" in user space. First, note that it **has to** execute in non-privileged mode as it **cannot directly execute certain machine instructions**.

- What the guest OS perceives as "physical" memory is in fact a chunk of virtual memory setup by the monitor.
- When the guest OS thinks it is running privileged mode, it is in fact running in (real) user mode, with the monitor knowing that the mode is "privileged" (**virtual privileged**)
- So when the guest OS executes a privileged instruction (in its fake privileged mode), it will cause a fault. This is exactly what should happen.

## This is not interpretation!

The trick is to execute everything directly as much as possible, only intercepting certain "sensitive" events (instructions).



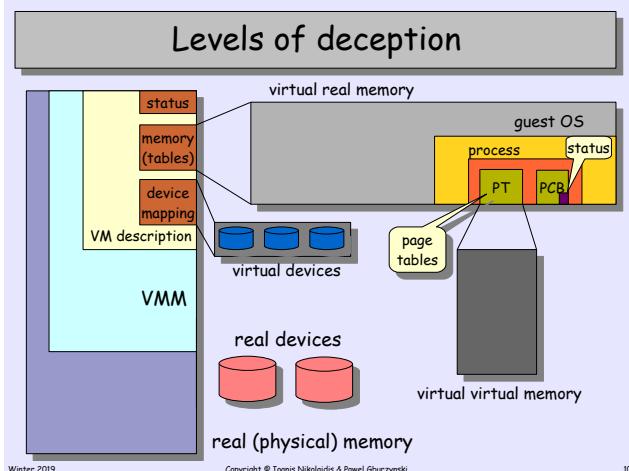
Winter 2019

Copyright © Ioannis Nikolaidis & Paweł Górzynski

## EXAMPLE: SYSTEM CALL

If a process in the guest OS executes a system call:

1. it is captured by the host OS.
2. The host OS (the virtual machine monitor) will see the syscall coming from a process which it knows to be an OS
3. bounce back the system call back to the kernel space of the guest OS. Where it will take care of it.



The idea is very simple. When the CPU is running regular instructions, there is no need to do anything (regardless of the level). This assumes that anything that must be emulated requires an interrupt. E.g., **suppose that the guest OS wants to switch the privilege level (to user)**

1. It executes a privilege instruction (thinking it has a right to do so, as it is running in "privileged" mode)
2. But this is an illegal instruction (the OS is in fact running in non-privileged mode at the host OS level), so this causes an interrupt
3. The monitor receives it and changes the virtual privilege for the emulated machine (and emulates the necessary effects)
4. Then, it gives back the CPU to the virtual machine.

## I/O REQUEST

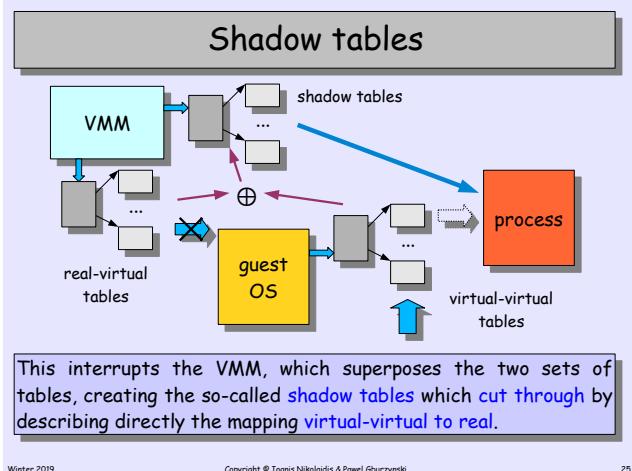
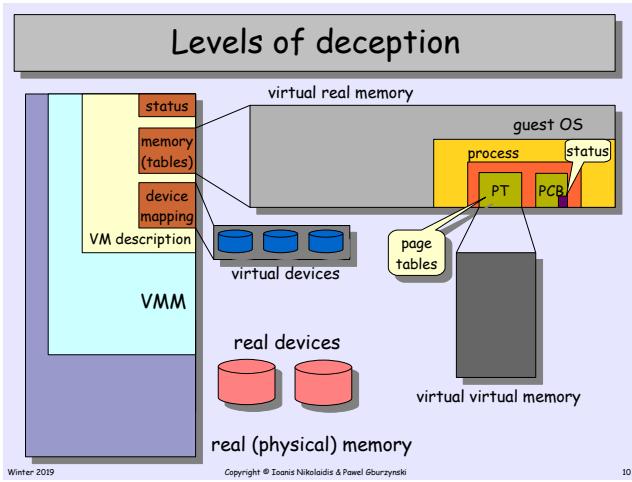
If a process in the guest OS is trying to do an I/O operation then the guest OS

1. is trying to do something "sensitive" like referencing the fake device registers, which are not formally mapped.
2. This will interrupt the virtual machine monitor which will emulate the I/O operations
3. and passes the outcome to the guest OS.
4. Eventually, the guest OS will pass the results to the process (all the time thinking it is doing the real job)

The real job is in fact being done by the virtual machine monitor, which, in particular, will emulate all the requisite interrupts in the guest OS, and so on.

## MEMORY MAPPING

When the guest OS tries to create its own page table as it wants to run a new program it thinks it has access to physical memory. In reality it does not, in reality the guest OS is mapping the its logical address to the logical address of the host OS, created by the virtual machine monitor. In order to aid this translation the VMM will create and maintain a shadow table that combines the table of the two.



## TABLE SWITCHING/UPDATING

Note that when the virtual-virtual process leaves its context, there is always an interrupt which (necessarily) goes to the VMM, so it can replace the tables as required - to keep the guest OS fooled. A tricky bit is suppose that the guest OS changes something in the V-V tables. How does the VMM learn about it?

It involves the update to the respective shadow tables.

Unfortunately, updating tables involves just writing to memory (nothing "sensitive"), so the VMM is not interrupted.

One way out would be to keep the V-V tables in read-only area (the VMM knows where they are). But it isn't really needed! Remember the TLB? The guest OS has no right to assume that a change in the tables becomes effective, until it invalidates the TLB!

So long as the instruction to do that is privileged, we are in the clear.

## THE GOAL OF VIRTUALIZING OSES

Smooth virtualization is possible if:

- All "sensitive" instructions are privileged.
- There are no instructions that silently behave in a different manner depending on the state (privileged/user). In other words: you either execute the same way in both states, or are privileged.

If you want to be formal, then you can (theoretically) emulate everything by declaring all memory "absent" and intercepting ALL instructions (at a tremendous performance hit).

## FUCKING INTEL

### The Intel (x86 and friends) is broken!

Why? Because at that time, most people (but not your instructor :-)) thought that virtual machines were gone for good (so nobody cared).

**pushf, popf** Push/pop **EFLAGS** onto/from the current stack. In non-privileged mode (any mode != 0), ignore **IF**.

This is a small and irrelevant detail, but it seriously damages the otherwise impeccable concept.

Note: on the 370, the virtualization was perfect!  
E.g., you could have this:



### Paravirtualization

The **pushf/popf** issue only affects (guest) OSes (which you have to be able to trick into the **virtual-privileged** mode).

So you can persuade their authors not to use those instructions (easy, say, with Linux - macros + recompilation).

Generally, an approach where the virtual machine (somewhat) differs from the real machine is called **paravirtualization**.

→ The VMM can provide hooks for guest OSes (sometimes optional), e.g., protocol shortcuts for virtual devices

→ The source code of a guest OS can be modified to take advantage of such hooks

→ That was done even on the 370 (VM-CP/CMS), even though no tricks were formally required there

## BINARY REWRITING

---

Sometime the source code may not be available; then you can try to modify the binaries:

- initial preprocessing. Generally, you cannot grow the code, but you may be able to replace the opcode (with one of the unused/illegal opcodes).
- On demand. recall lazy loading. Generally, you know when a code page id demanded (and about to be executed).

Note that you need this for the OS code only. The translated code can then be cached (VMWare).

Generally, the VMM will try to give preference to the pages belonging to the guest OS image, so they are not victimized (at least not too often)

## HARDWARE ASSIST (INTEL'S VT-X)

---

- A new state (root) added - to be used by the VMM
- Hardware data structure to keep track of the state of VMs
- Sensitive instructions cause VM exits, i.e., automatic transitions to the VMM Device virtualization hooks. i.e, devices like NICs can be "context switched" to a particular guest OS, with DMA and interrupt redirection.

# PROFESSOR'S VILLAINS AND CREDITS

## C Minion

*Medium Website psychic, neutral evil*

**Armor Class** 12-15

**Hit Points** 100

**Speed** 1 Month

STR	DEX	CON	INT	WIS	CHA
29 (+9)	10 (+0)	17 (+3)	12 (+1)	11 (+0)	15 (+2)

**Senses** —

**Languages** C

**Challenge** 1 (200 XP)

## COVERAGE

Everything that is being discussed in class during the existence of the minion. Including the readings

## COMPOSITION

**Programming 1.** You create an client to client messenger using signals.

**Programming 2.** Creating a file deduplicator.

**Programming 3.** Simulating a TLB and analyzing the performance.

## Midterm Champion

*Large paper psychic, neutral evil*

**Armor Class** 20

**Hit Points** ?

**Speed** 70 minutes

STR	DEX	CON	INT	WIS	CHA
23 (+6)	14 (+2)	21 (+5)	16 (+3)	13 (+1)	20 (+5)

**Senses** —

**Languages** C, Unix

**Challenge** 8 (3,900 XP)

## COVERAGE

Everything from the beginning to Lecture 10

## COMPOSITION

**Multiple Choice.** Most of the exam is MC

**Short Answer.** Some short answer in a specific format

## Grade Slayer Tiamat

*Gargantuan paper psychic, Lathrain of grades, chaotic evil*

**Armor Class** 38

**Hit Points** ?

**Speed** 180 minutes

STR	DEX	CON	INT	WIS	CHA
30 (+10)	10 (+0)	30 (+10)	26 (+8)	26 (+8)	28 (+9)

**Senses** —

**Languages** C, Unix

**Challenge** 30 (155,000 XP)

## COVERAGE

- Compiling, linking, and loading
- Processes
- Process scheduling
- Synchronization
- Deadlock
- Memory/Virtual Memory
- Disk scheduling
- File system
- VM

## COMPOSITION

It going to be like the midterm but with 19 questions

# CREDITS

## GENERAL

---

- Created by u/DnD\_Notes, April 26, 2019
- Lecture Slides by Ioannis Nikolaidis & Paweł Gburzynski
- Typesetting engine: [LATEX](#)
- Dungeon and Dragon (5e) LaTeX [Template](#)

## ART

---

- Figures for the cover art is from [D&D Beyond](#)
- Andy the D&D Ampersand is from [Dungeon and Dragons](#)
- Cover art formatting and design done in Photoshop CC 2019

## DISCLAIMER

---

This document is completely unofficial and in no way endorsed by Wizards of the Coast. All associated marks, names, races, race insignia characters, locations, illustrations and images from Dungeon and Dragon world are either ®, ©, TM and/or Copyright Wizards of the Coast Ltd 2012-2018. All used without permission. No challenge to their status intended. All Rights Reserved to their respective owners.