

BS

Evolution of Computers

first gen 1945-55

electronic vacuum tubes or relays
patch bays

second gen 55-65

transistors + reliability

punch cards + batch processing

early languages: FORTRAN/COBOL

third gen 65-80

IC - Integrated Circuit

multi programming + multicasting
e.g. fairly between users
CPU

fourth gen 77 - now

personal computers

fifth gen 90 - now

mobile computers

OS

Hardware Abstraction + portability, reuse

VM for applications / execution environment

access to other resources Memory, IO, Devices, ...

Resource Management & Control

shared access

fair & efficient

optimize

caches

Concurrency & Sync

Isolation

memory protection

Access Control

to sensitive ops

System Design Concepts

Abstractions

processes, threads, sockets, files, disks, devices

Mechanism

open, create, fork, schedule, exec, close, send

Policies

how...?

Computers

Assembly Language Intel x86-32 x64-64 types

mnemonic, operands/instructions, opcode, param

Endianness

→ Big Endian MSB first 2.B. 0xAF7F

→ Little Endian LSB 0x0000AF7F

(CPU Registers tiny, fast)

(Segment r: identify segments in memory)

General Purpose r: immediate values, mem addresses

Program status and control r.

eip - extended instruction pointer next to be executed

other flags

8x 32 bit general purpose on x86

EAX - Accumulator

EBX - Base (mem addr.)

ECX - Counter (loops, strings)

EDX - Data (I/O Pointer)

ESI - Source Index for String op

EDI - Dest

EBP - Base Pointer (stack)

ESP - Stack Pointer

User Interface Program

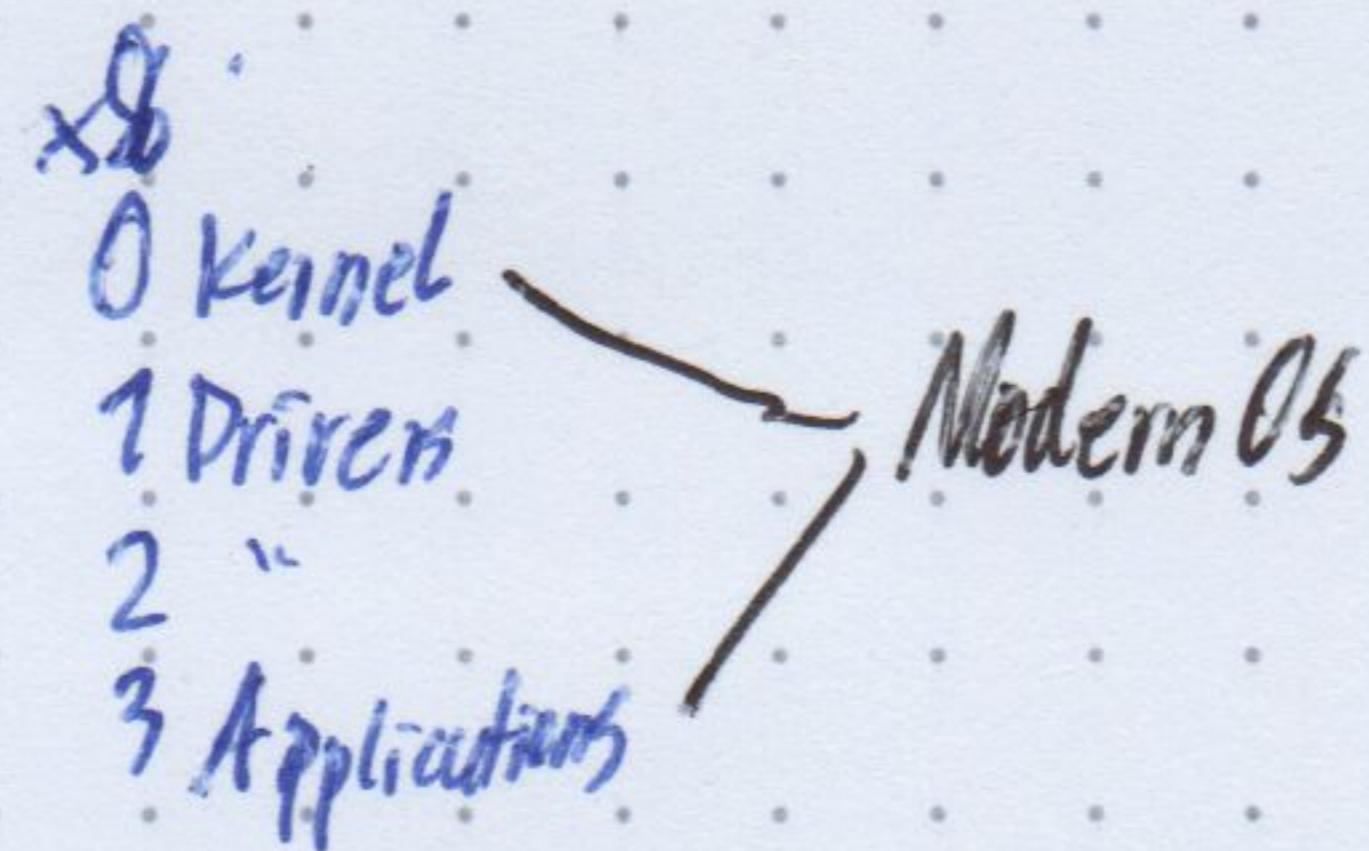
OS Interface

OS

Hardware Interface

OS Protection Boundary

User Mode by privilege bit
in processor
privileged ops in kernel m. else fault



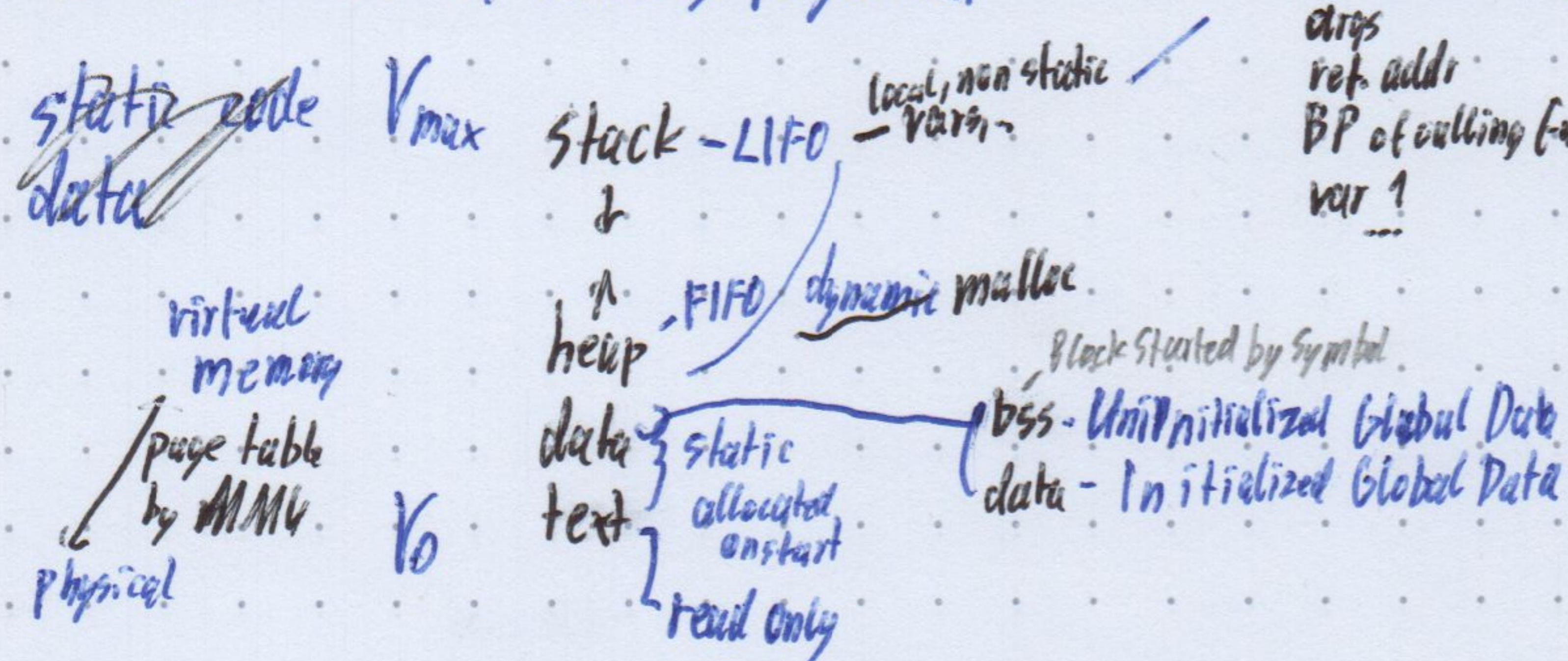
System Calls - request service from kernel

1. call library syscall wrapper
 2. accumulator := syscall number & param
 3. trap function \rightarrow software interrupt
 4. Processor switch kernel mode
invokes syscall dispatcher
 5. look up number \rightarrow handler
 6. processing
 7. return control in user mode
- hardware related
creation & exec processes
communicate with kernel services
e.g. scheduling

0. read
1. write
2. open
3. close

9. mmap map additional memory to process
10. mprotect change permission to mem area
11. munmap remove mem from virtual addr space

Process - abstraction of running program



Stack Frame per Function Call

args
ret. addr
BP of calling frame
var !

Iria x86 call instruction
e.g. call foo

automatically pushes ret. addrs
instruction after call
ret
pops ret addrs, loads it into EIP

Multiprogramming

multiple processes on system conceptually own virtual CPU

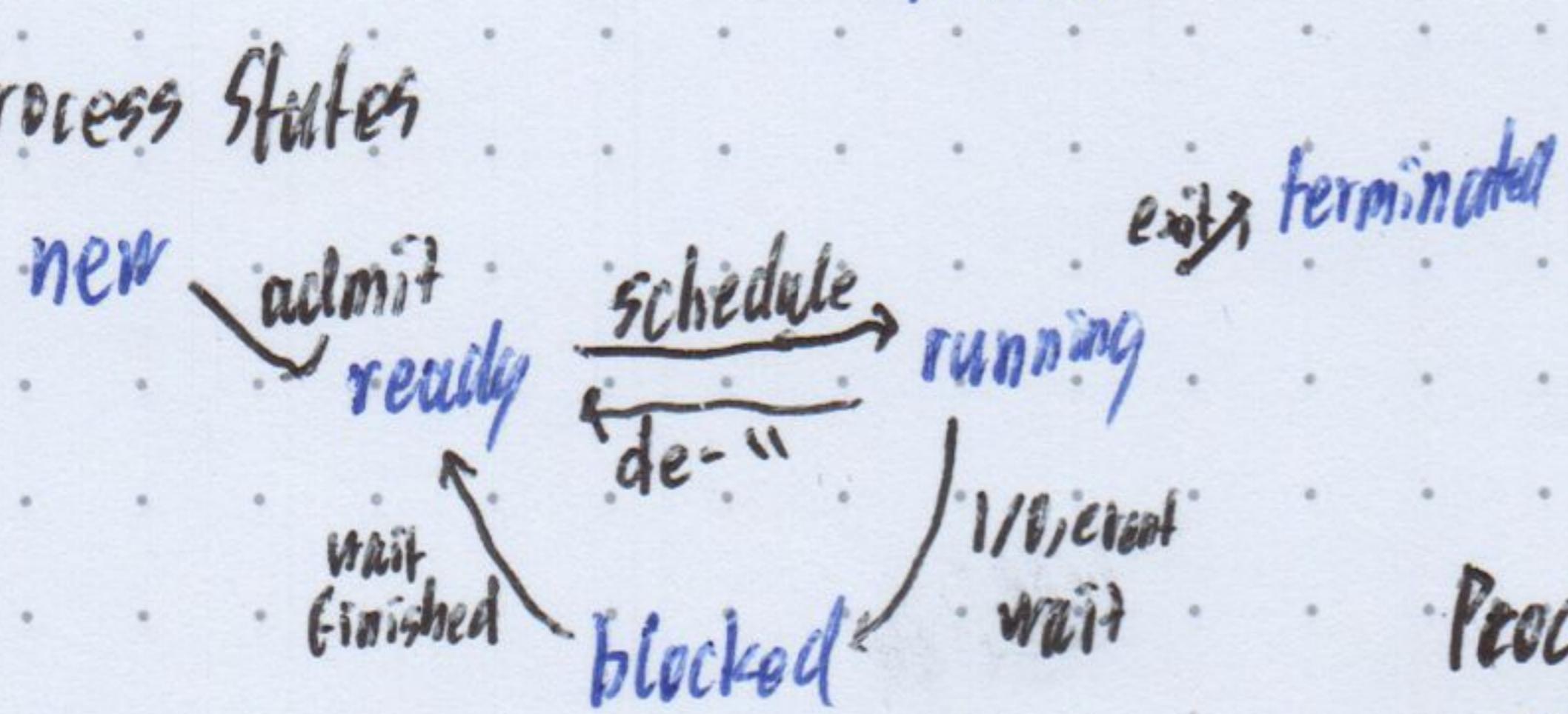
pseudoparallelism - through rapid context switches

Process Table

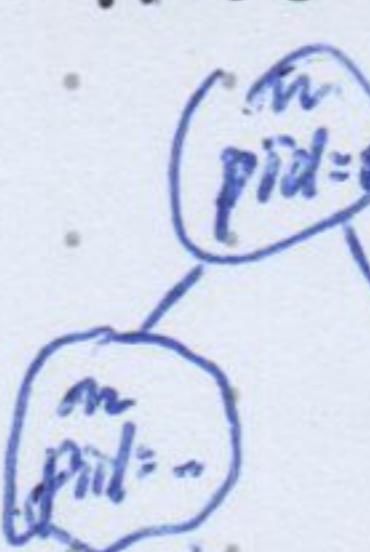
Non-active saved in PCB - Process Control Block

contains state
PC, Registers, Stack Pointer
Memory allocations, open files
accounting & scheduling info
ID, status

Process States



Process free



z.B.

POSIX - Portable Operating System Interface

Process Creation

pid = fork();

if (pid == 0) // child process

variant of exec() replace mem. space with new program

else

wait(NULL); // waits for child

exit(0); // terminates, deallocate

parent can kill() child processes

IPC - Inter Process Communication

Message Passing - based

Send() + Message Queue \rightarrow read()

- syscall overhead

BRUNNEN

Shared Memory - based

mem region mapped into both virtual mem

+ faster - complex

error prone

7

Memory Management / Organisation

Requirements

- Protection between processes
- Fast translation
- context switch

No Memory Abstraction

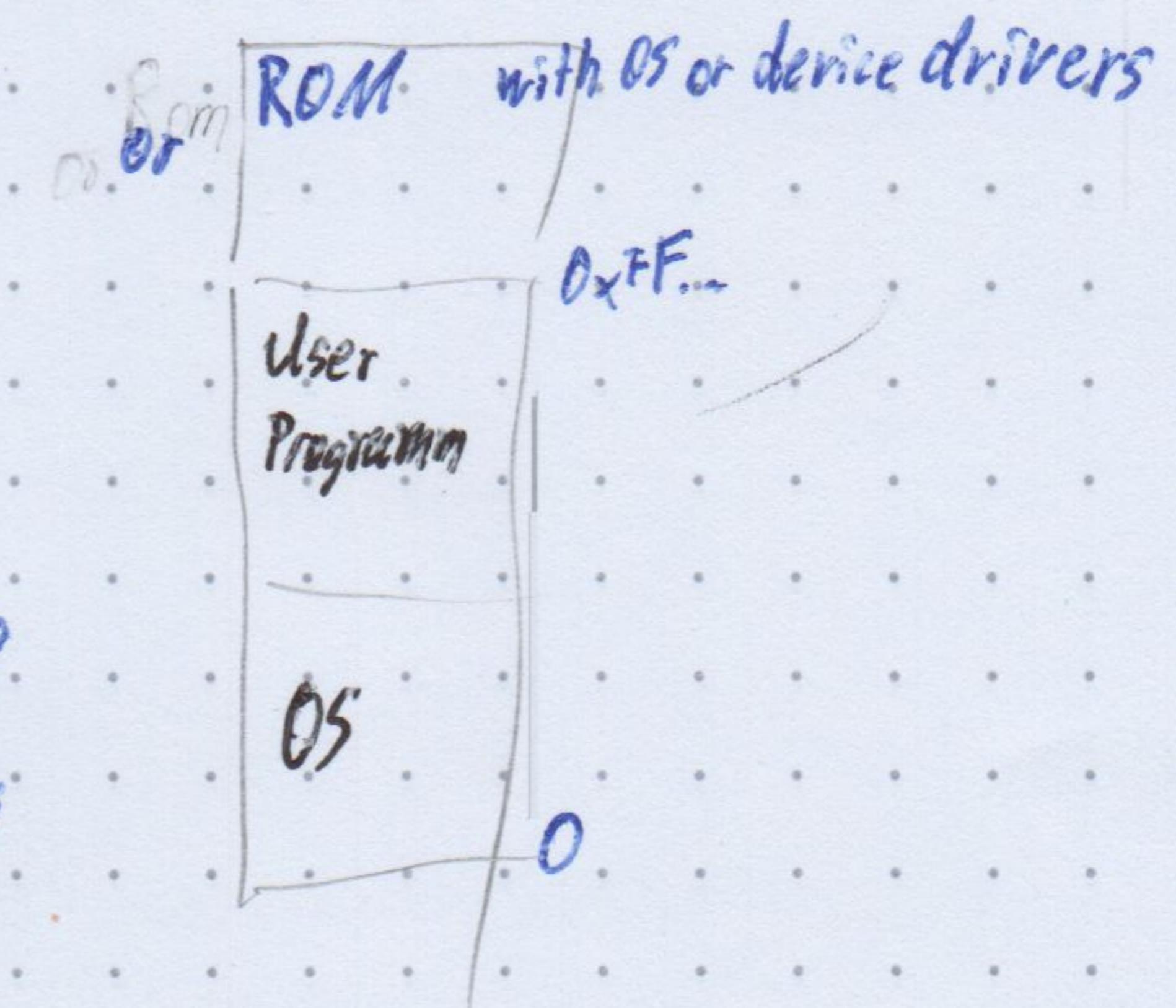
processes see physical memory

⇒ program needs to know where it lies ⇒ parallel only with swapping

or Static Relocation during loading

t_{slow}

needs additional info



Allocate Memory to Process + multiple processes in parallel

Fixed Partitioning

$$\text{physical addr} = \underline{\text{base addr}} + \text{logical addr.}$$

each process gets partition of same size

+ easy

+ fast context switch

- internal fragmentation: mem inside partition not used by process

zuria

zuwenig

S

+ memory protection

- external fragmentation: loading, unloading creates empty holes between partitions

Variable Partitions

variable size partitions

base and limit register

physical > base+limit ⇒ addressing error

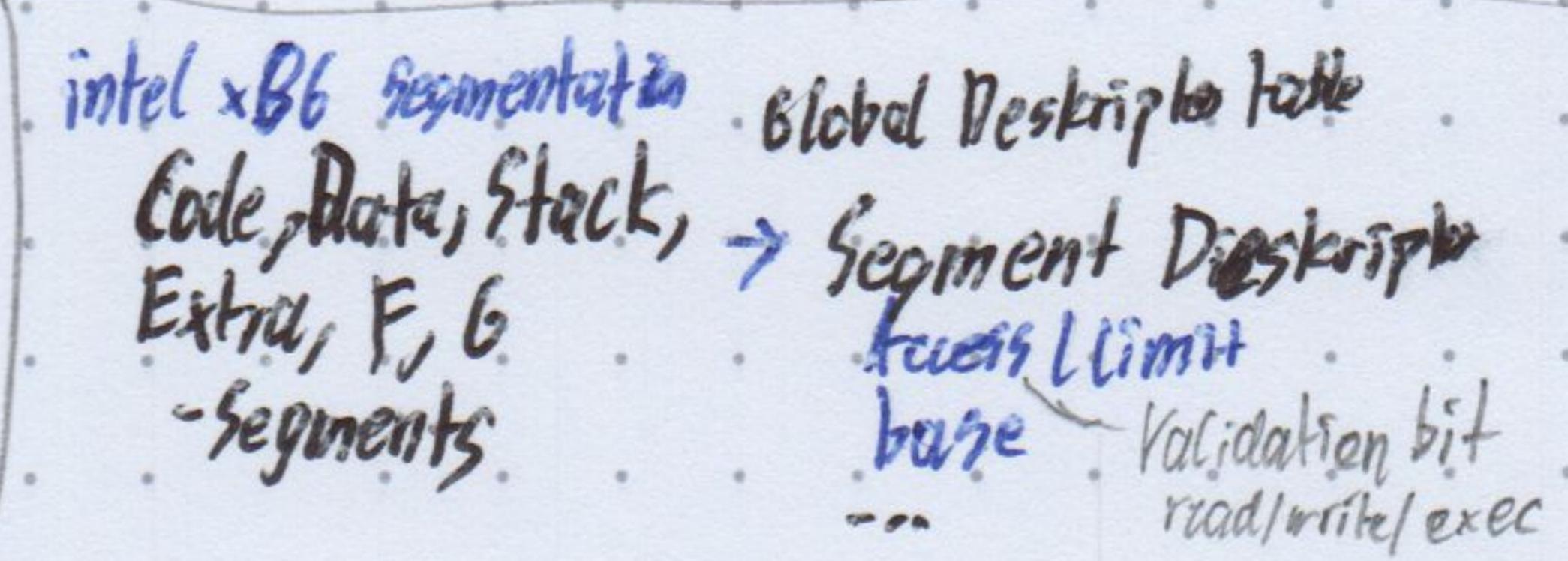
Virtual Memory

physical location can change

but not continuously

virtual addr must not be in physical memory

⇒ allows process to mem



Segmentation extension variable Partition

allows many segments per process

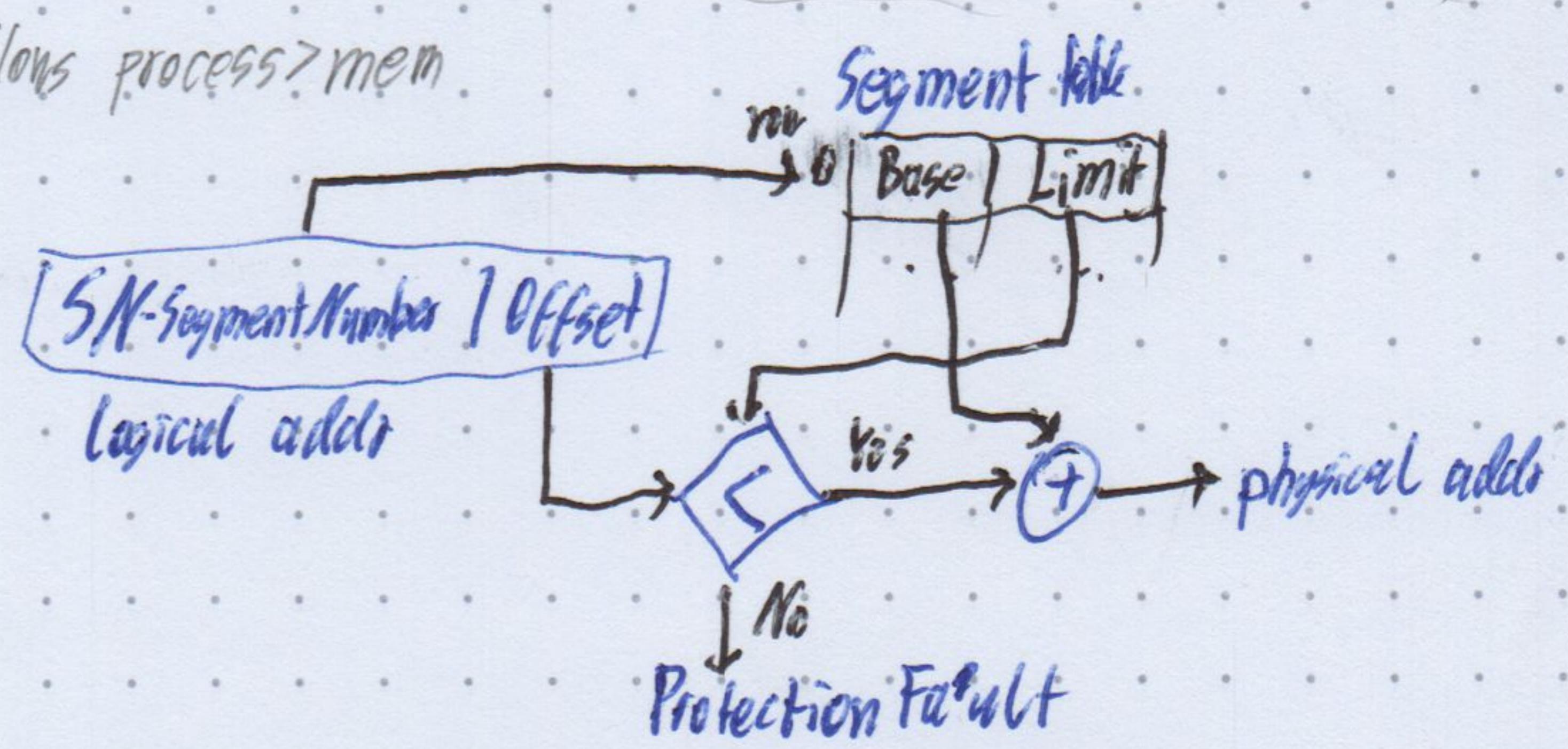
STBR - Segment Table Base Register

SLTR - Segment Length Register

+ only segments need continuous memory

+ no internal frag, + mem protection

- external fragmentation



Paging

Page - equally sized mem block (usually 4 kB)

Permissions
of mapping → managed independently
rwX
0.1 memory frame
physical page

Page Sharing

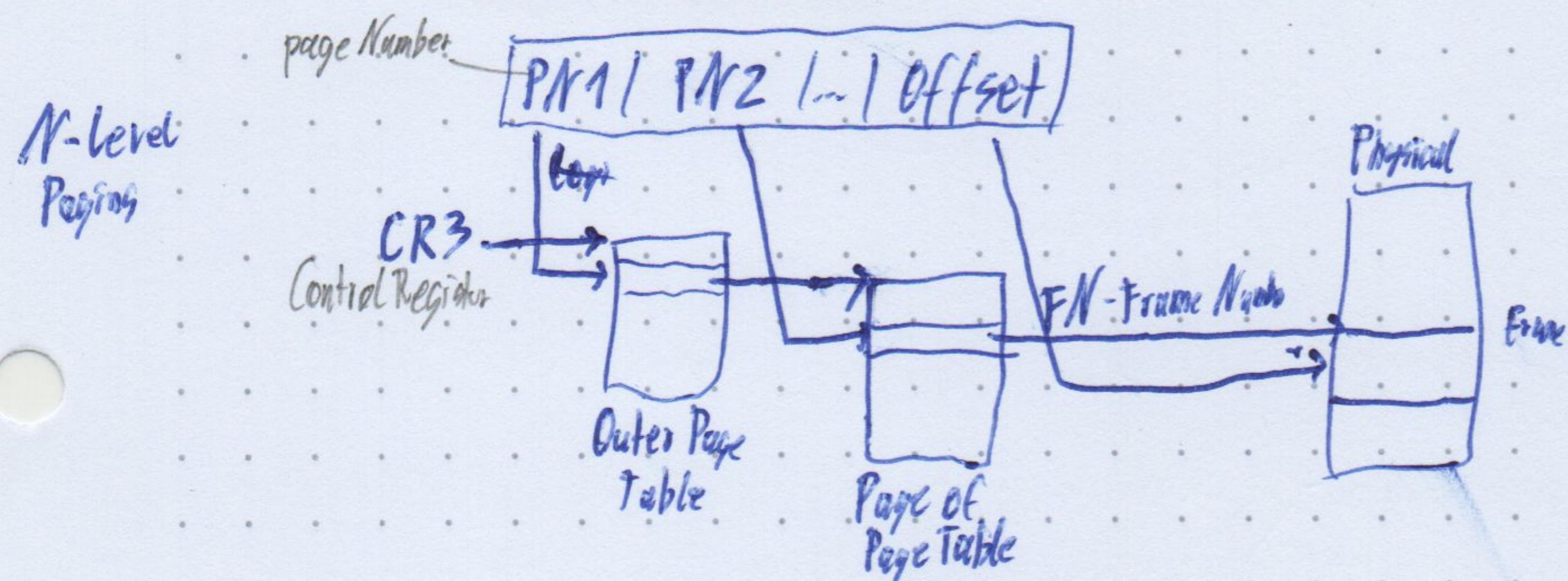
n virtual pages → 1 physical page

e.g. load common data only once

read-only

shared mem for IPC

Address Translation



large ⇒ not in Registers

→ TLB - Translation Look-aside Buffer
Cache for translations

× B6 4 Level Pages

PPML4E
Page-Directory-Pointer-table) can be skipped for
Page Directory
Page Table
larger pages

Page Protection

access control information in Page Table Entries

protection bit

Valid bits

↳ present in Memory

↳ page fault

Read/Write

User/Superuser

Write Through

Cache Disable

Accessed

Dirt

Global

Scheduling & Threads

$$\text{multiprocessing} \Rightarrow \text{CPU utilization} = 1 - P_L^n$$

degree of multiprogramming
fraction of time waiting for I/O

CPU - \rightarrow I/O-bound processes

Scheduler - lowest level of OS

handles interrupts & scheduling
↳ which first?

Interrupts - signals event to OS

e.g. input ready/output complete

exception (I/O, attempt protected mem.)

clock interrupts

privileged instr.

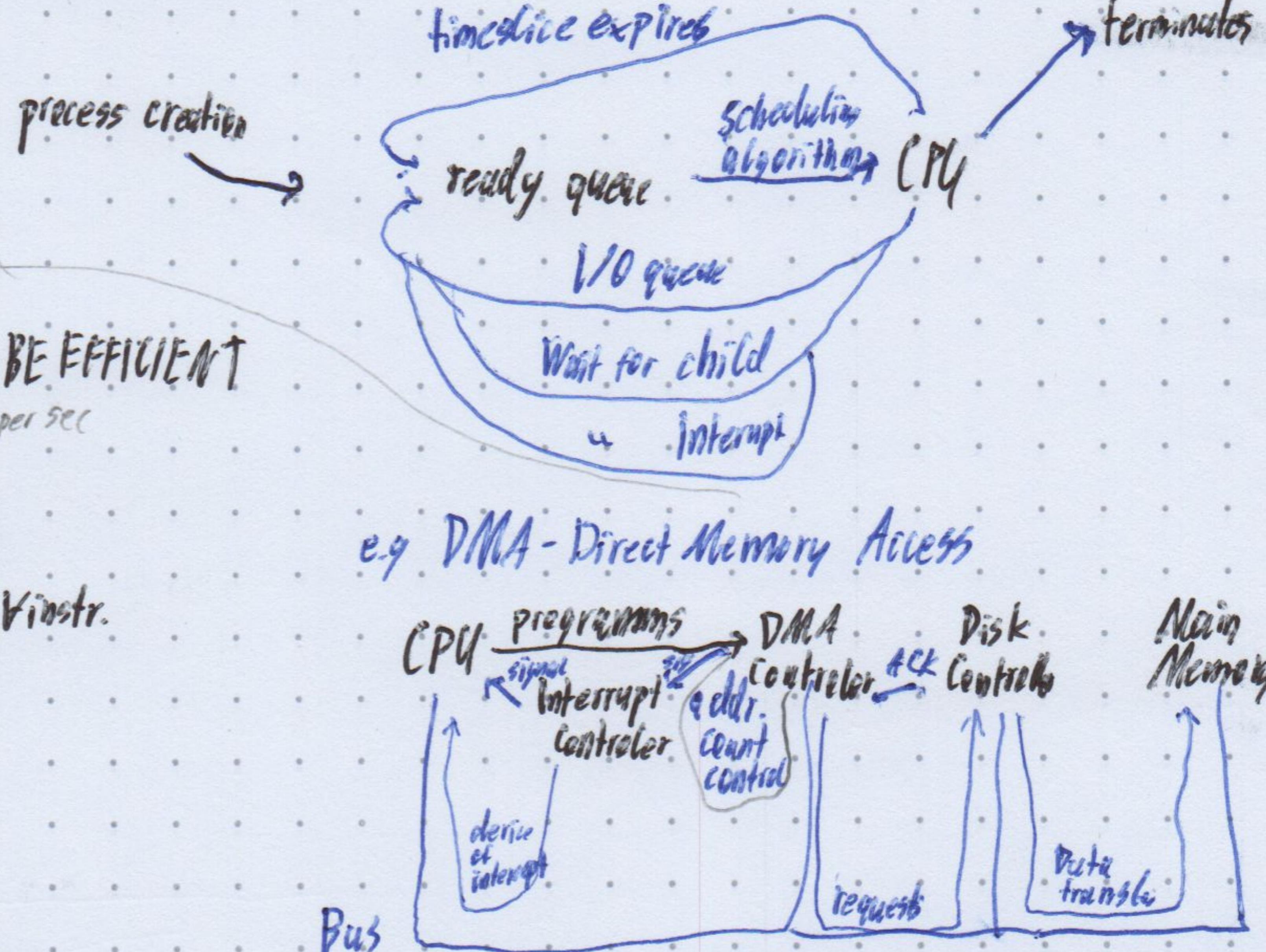
(CPU checks interrupt-request line after V instr.)

↳ dispatches to routine at specific addr.

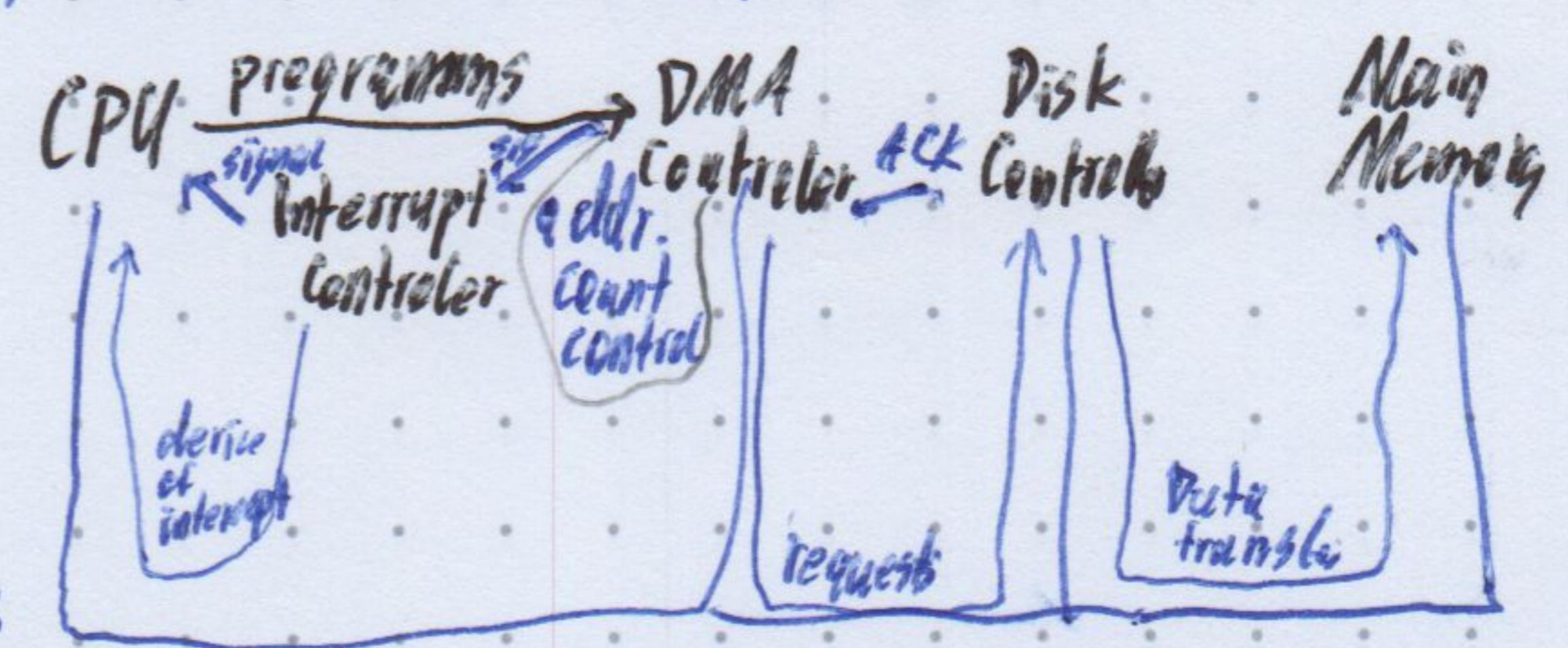
↳ context switch

MUST BE EFFICIENT

100s per sec



e.g. DMA - Direct Memory Access



When to schedule?

fork < Parent or child next?

process terminates blocked

interrupt incl. clock interrupt e.g. 250 Hz

How?

environment

batch - big tasks \Rightarrow throughput, turnaround time, utilization \Rightarrow non-preemptive

interactive - active users \Rightarrow response time (for interactive), enforce proportionality

e.g. Multimedia

Real-time - hard timing requirements

\Rightarrow + predictability

↳ else deterioration in media

general goals

fairness in same priority

Policy enforcement

Balance of computation/I/O

user

for utilization with

Preemptive \leftrightarrow Non-preemptive scheduling

↳ suspend, when

max timeslice

↳ until blocks or releases

↳ Strategies

First Come, First Served

resp. time

diff. types

SJF - Shortest Job First

- needs estimate

Shortest Remaining Time Next

+ optimal avg. turnaround time

preemptive SJF version

- long jobs wait longer

↳ Round-Robin

First come, first serve with quantum - max timeslice at a time

↳ too long - resp. time

too short - many context switches

↳ Priority

queue for each priority ~~higher scheduled before lower~~

higher before lower

↳ Multiple Queues

lower priority queues get \geq e.g. $\geq 2 \times$ mehr quanta when scheduled

↳ Shortest Process Next

estimates e.g. based on previous observations $T_h = \alpha T_{h-1} + (1-\alpha) T_{\text{proc}}$

+ if alternating pattern of comand-exec

running round

length param

↳ Guaranteed $\frac{1}{n}$ CPU time

↳ users next process with ~~with~~ min ~~entitled~~ $\frac{\text{received}}{\text{entitled}}$ time ratio

e.g. Linux CFS

+ fair Completely Fair Scheduling

↳ Lottery

winning \sim # tickets

↳ co-operating may share

↳ user A: A, B

user B: B, E

↳ Fair-Share user gets fraction CPU time regardless H processes \Rightarrow AEB-EAE...

Threads - lightweight processes

per thread PC, Reg, Stack, State
per process code, data, files,
accounting info, ...

+ thread scheduler
+ no IPC between
+ faster creation, destruction

User \leftrightarrow Kernel-Level implementation (where Runtime System, ~~Thread~~ Thread Table)

+ faster switching
if no blocking syscall
anyway
+ no OS support
needed

+ better blocking
syscall bundling

- replicate fork?

Hybrid Kernel-Level Threads having User-level threads

POSIX Threads

Pthread - create
exit
join wait for specific thread
yield release CPU
attr.init
attr.destroy

Swapping

no mem \Rightarrow swap process to disk

Manage free mem

Bitmaps
Free allocation unit
Occupied size

may waste space
large bitmap slower finding

Linked Lists

ordered (process / fl., start pos, length)

Hole

Allocation of segments

Allocation Strategies

First Fit

from last place of search

Best Fit

- searches all

smallest large enough - leaves small rest holes

Worst Fit

largest available

Quick Fit

list for each hole size + fast find

- slow combining on swap

Paging only pages 'swapped out' / 'enriched' to ~~mem~~ disk

Page Fault

minor in mem but not mapped in process add. \Rightarrow upd. page table

major not in mem \Rightarrow load

Segmentation fault invalid addr. \Rightarrow kill p

Page Replacement Algorithm

frame limit per process

Optimal

least likely referenced in near future

impossible, but after run + baseline toward other

Not recently used NRU

Class	R	M	N	random of lowest class
0	0	0		
1	0	1		
2	1	0		
3	1	1		

First-In, First-Out + easy
- problematic

Second-Cohesive

FIFO, but when R=1 requeued as R=0

Clock

second chance but circular list

(clock hand moves not pages)

@ every

Least Recently Used LRU - needs update @ reference
impl. with list last used in front, least in back

Not Frequently Used NFU

counter ~ page evict lowest counter
- @ clock interrupt $t=R$ - doesn't forget long past references

Ageing

counter ~ per

bit width limits remembering

- @ clock interrupt $t=R$, ch $\gg 13$

Working Set

set of pages process currently using utilizes Locality of reference

last time ~ page

@ page fault: R=1 update time R<0

R=0 age > T \Rightarrow evicted

if none random with oldest first

WSClock last time/age ~ page, clock hand

R=1 set age, advance

R>0 age > T \wedge M=1 schedule write, advance

M>0 evict

BRUNNEN



Page Table Entry

Caching Disabled bit - e.g. for memory mapped I/O

Referenced bit $\leftarrow 0 @ \text{clock interval} \leftarrow 1 @ \text{referenced}$

Modified/Dirty bit

Protection bits rwe

Present/absent

Thrashing if process \gg available frames

Seitenflattern

evicted soon needed \Rightarrow many page faults

Allocation Policies

Local \leftrightarrow Global page replacement

Selected among all

frames of P framed dynamic

- risk thrashing + if working set size varies
or wasting frames

Allocation size

periodic fixed Machine processes fixed fraction

- doesn't processes of diff. sizes

proportional to process size

adjusting based on Page Fault Frequency PFF

PFF > A \rightarrow increase alloc

PFF < B \rightarrow decrease alloc

r p tends to enlarges subset at a time

If no read no write scheduled, evict any clean else current

File Systems - controls how data stored/received

File - abstraction data blocks on disk

Types

regular files

directories - system files for structure

Character special files - model I/O devices

Block special files - model disk

Naming

allows other processes to reference it

conventions

Unix case sensitive, extension only convention

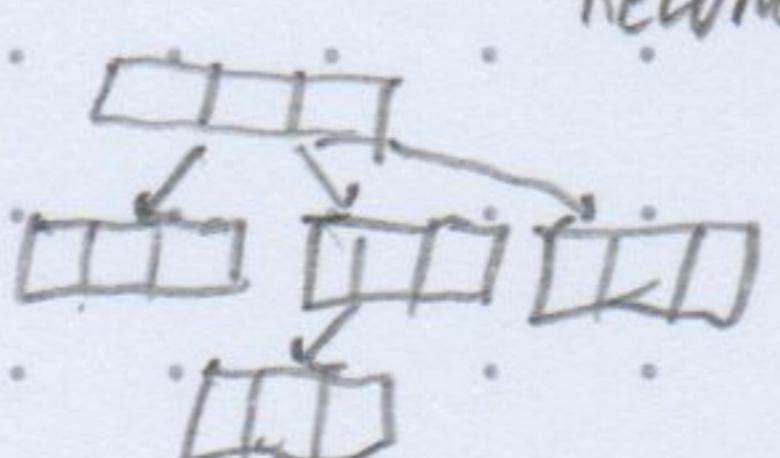
Windows not " , extension link to application

File Structure

byte oriented. Sequence of Bytes

record oriented Records

tree structure



Types of Access

e.g. sequential (magnetic tapes)

random out of order possible

Positioning File Read(), read() specifies start
or seek() moves position

Attributes

Protection

Password

Creator

Owner

Read-only flag

Hidden

System " System file?

Archive " Has been backed up

ASCII/binary " ASCII, binary

Random access " Sequences, random

Temporary " Delete on process exit

Lock " allows only one p access

Record length # bytes in record

Key position offset of key within each record

Length # bytes in key field

Creation time

Time of last access

Change

Current size # bytes

Maximum "

Directories

root dir

hierarchical:

userdir

Pathnames init win

sub-dir

absolute starts with root / , \

file ffd

relative to working dir of process

L syscall can change it

Special names : current

parent

operations

create ... only entries

delete

open dir for reading

close dir free mem

read dir next entry

rename

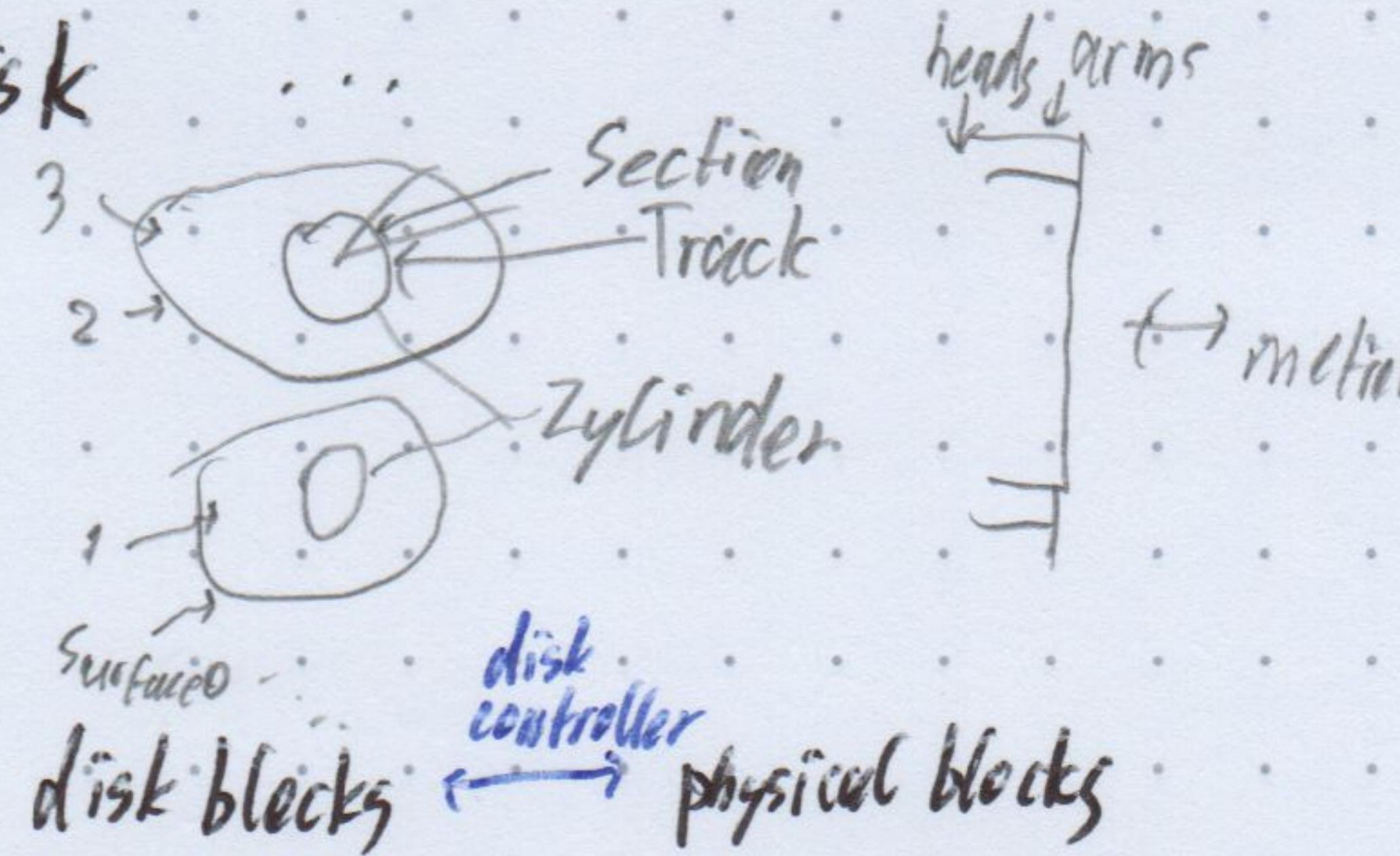
link new link file + dir

unlink removes "

filen dir

L removed if not referenced

Disk



Delays

Seek time to move arms to cylinder

rotational delay for sector

read/write time

e.g. executable

Header (Magic Number, Text size, Entry point, "", Flags)

Text

Data

Relocation bits

Symbol table

Data "

BSS "

Symbol table "

binary archive

Header (Module name, Date, owner, Protection, Size)

Object module

Header (...)

Object module

Disk Organisation

MBR - Master Boot Record in sector 0

used to boot

partition table at end

↳ always one marked active (contains OS)

one or more Disk Partitions

all start with boot block

typically: but file system dependent

superblock [magic num to identify type, #blocks, ...]

Free space map

I-nodes

Root dir

Files and dirs

File Allocation

Continuous Allocation

- + easy
- external fragmentation
- + efficient
- know max size beforehand

Linked-List Allocation

- first word of block pointer to next
- + no external frag.
- + dir entry only needs 0 for End
- + dir entry only needs slow random access
- start block
- block size not pow 2

File Allocation Table FAT

data structure in main mem

physical block → next block

{ 1. end offset
2. }

+ faster

+ Block per sf 2 - large

fully in mem

Index

I-Nodes

File attr.

Addr. of disk blocks 0

"

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

Remote Procedure Call RPC
 between processes on networked system

! may fail or duplicate
 timestamps, request history,
 ACKs, resend

stubs - proxy for actual procedure

client-side: locates server, marshalls call param

server-side: receives msg; unpacks param, performs procedure

Race condition
 result depends on order each processes ops on shared data

Critical region / section
 part of program modifying shared mem

Avoiding Race Conditions

1. No 2 p. simultaneously in crit. region
2. No assumption on CPU speed, number
3. No blocking other p. outside crit. regions
4. No waiting forever to enter crit. region

Disabling Interrupts after entering - guarantee enabling?
 enables just before leaving - multiprocessors

Lock Variables
 L shared var
 init 0
 @entering lock == 0 \Rightarrow 1
 lock == 1 \Rightarrow wait
 - descheduling right after lock test

Process Alternation / Spin lock
 shared var & turn
 p0:
 while(turn != 0) Xⁿloop;
 crit. region();
 turn = 1;
 non crit region;
 p1: analog
 - busy waiting
 : especially for high-prioritying
 only reasonable for short
 - p's of diff. speed

Peterson's Solution
 shared vars turn, flag array interested[N]
 interested[process] = TRUE
 turn = process;
 while(turn == process || interested[other]).
 / busy waiting.
 crit region
 interested[process] = FALSE;

Atomic Instructions by locking mem bus
 TSL RX, LOCK - loads LDR into RX
 TestSet Lock stores nonzero val in Lock
 XCHG RX, Lock - exchanges contents
 via TSL nach MOVE RX, #1
 Then check RX == 0
 More lock, NO after crit region

Monitors - programming language-level
 collection procedures, vars, data structures
 Only 1 p at a time in monitor
 impl by compiler
 monitor Name
 type name;
 procedure name(); begin ... end
 end monitor

ctr + t could be impl.

r=ctr
 r=r+1
 ctr=r

no busy waiting:

Sleep and Wakeup + processes blocked while waiting

e.g.
 producer:
 item = produce-item();
 if(count == 0) sleep();
 insert-item(item);
 count = count + 1;
 if(count == N) wakeup(consumer);
 consumer:
 if(count == 0) sleep();
 item = remove-item();
 count = count - 1;
 if(count == N-1) wakeup(producer);
 consume-item(item);

! deschedule right before sleep
 missing wake up
 go to sleep for good

Semaphores Dijkstra 1965

down(semaphore)
 value > 0 \Rightarrow decrement
 value == 0 \Rightarrow sleep, down op remains pending
 up(semaphore)
 increment
 if ≥ 1 p = H processes sleeping on semaphore
 \Rightarrow wake one up

E.g.
 Producer
 item = produce();
 down(&empty);
 down(&mutex);
 insert(item);
 up(&mutex);
 up(&full);
 Consumer
 down(&full);
 down(&mutex);
 item = remove();
 up(&mutex);
 up(&empty);
 consume(item);

Mutex simplified semaphore only locked/unlocked
 impl. with atomic inst. in user space

mutex-lock:
 TSL RX, MUTEX,
 CMP RX, #0
 JZ E ok
 CALL thread.yield
 JMP mutex-lock // try again
 ok: RET

mutex-unlock:
 MOVE MUTEX, #0
 RET

Condition Variables allows ps to block until condition
 wait(condition var) blocks p
 signal(condition var) wakes up random p sleeping on it

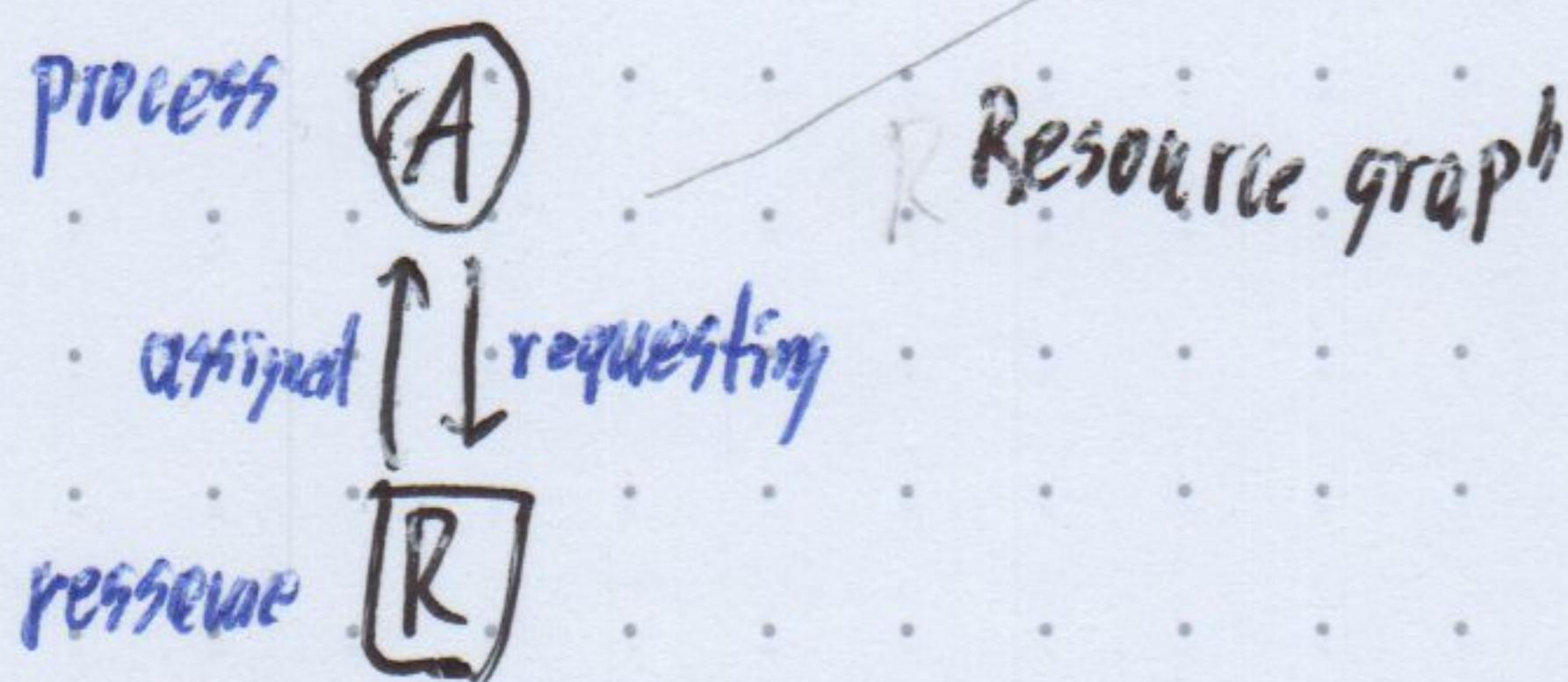
Deadlocks of set of p's

⇒ If p wait for event only another in set can cause

Conditions for Resource Deadlocks

1. Mutual Exclusion of Resources (assigned to 1 p or available)
2. Hold & Wait p holding Res. can request additional
3. No pre-emption
4. Circular wait condition circular list p: each waiting for R. held by next

Modelling



Handling

Ignoring

Likelihood?

more than other crash causes?

Detection in current situation

construct resource graph

single res:

DFS-list of cur. path, mark visited arcs
node twice in path → cycle

multiple res of res classes:

E = (E₁, ..., E_m) resource vector
#res Existing in class i

A = (A₁, ..., A_m)

#res Available in class i

C = Current Allocation Matrix

C_{ij}: P_i currently holds C_{ij} of class j

R: Request Matrix

P_i requests R_{ij} of class j

Recovery

Pre-emption

- not always gracefully possible

Killing processes

until cycle broken

Rollback

to periodically recorded checkpoints

- Work lost

- side effects?

- idempotent?

1. unmarked P_i with R_i ≤ A
No / found
all unmarked deadlocked
mark P_i
A + = C_i

Two-Phase-Locking

1. phase: req. res for next op

2. phase: exec, release

if 1. couldn't lock all → release, try again

- real-time systems
- process control systems

Communication Deadlocks

A, B waiting for msg from other

e.g. file system → timeout detection

resend/recovery

Live Lock hold by either p
because can't acquire res

2 p acquire and release res simultaneously

- not blocked but can't proceed

Starvation

e.g. shortest-job-first res allocation policy

If short jobs with need for r continue to come in

long job never gets r

avoided by first come, first served

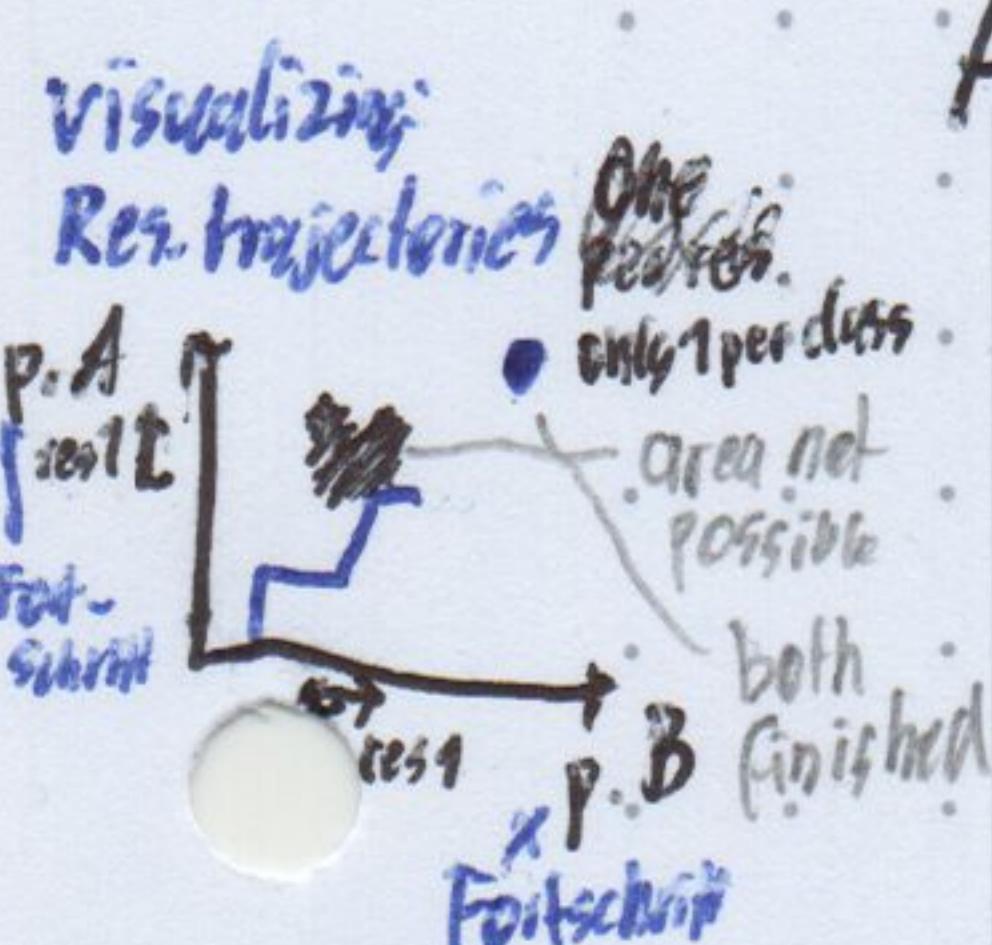
Avoidance only grant res. if safe

I scheduling order p complete even if p request max res. now

Banks Algorithm

like detection with mult. res.
marking as terminated

- needs rarely known
- # p's not fixed
- res may fail



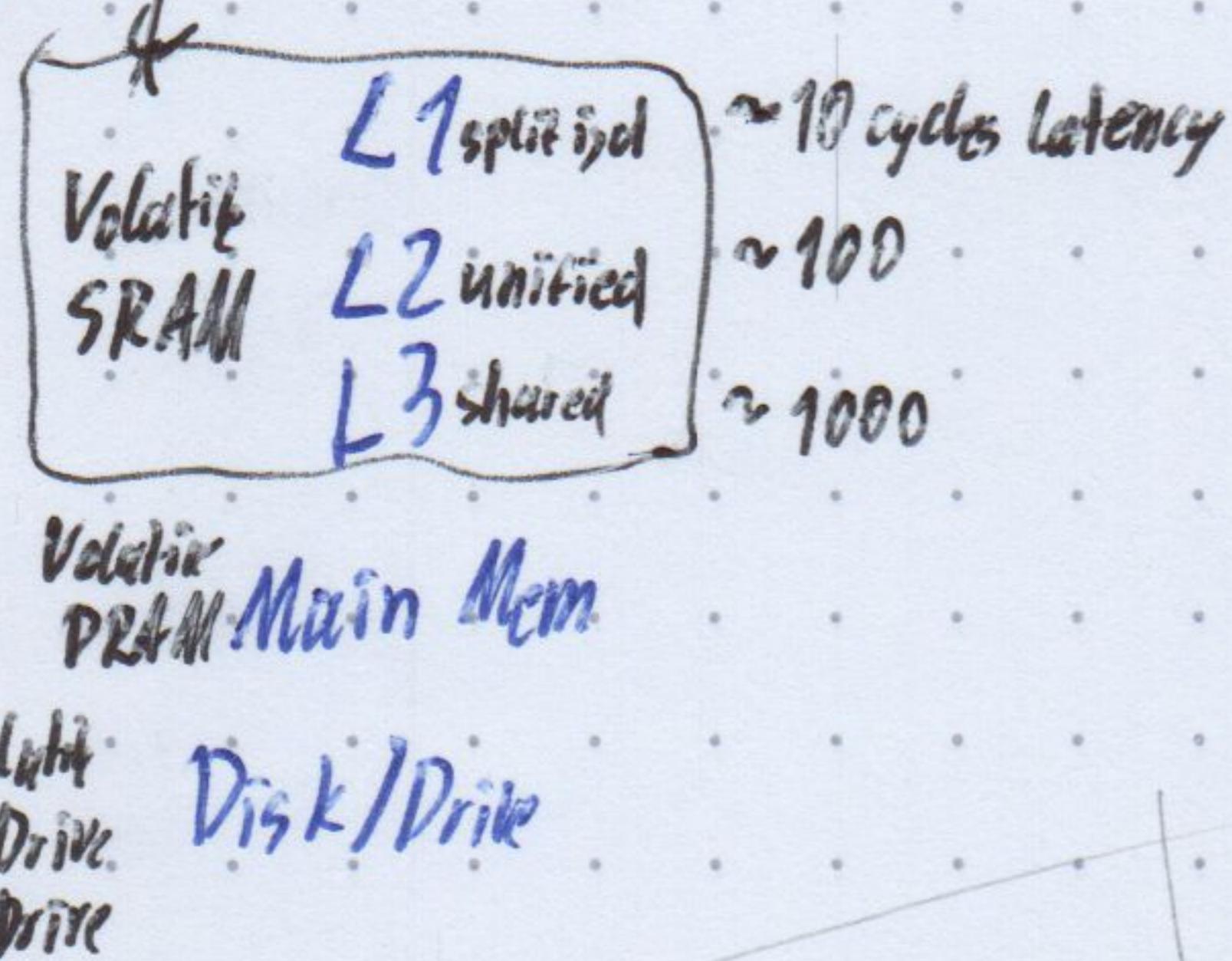
Caches & Memory Allocation

Motivation

main storage
shared between processes
cheap but slow
low frequency

Hierarchy

expensive
Small Latency



cheap
high capacity

Cache Addressing

PIPT - Physically Indexed
Physically Tagged

index, tag extracted from physical mem addrs.

- needs virt. → phy. translation before cache
⇒ latency

VVVT - Virtually Indexed
Virtually Tagged

index, tag extracted from virt. addrs.

- only unique per process

homonym problem: 1 virt. add. → diff. phy.

synonym ↪ 2 diff. virt. → same phy.

VIPT - Virtually Indexed index from virt.
Physically Tagged

addr. translation parallel to indexing step

physical tag prevents homonym

synonyms possible

Split → Unified Caches

↳ dedicated cache for instructions
data

Non-exclusive / Partially inclusive - can, but doesn't have to

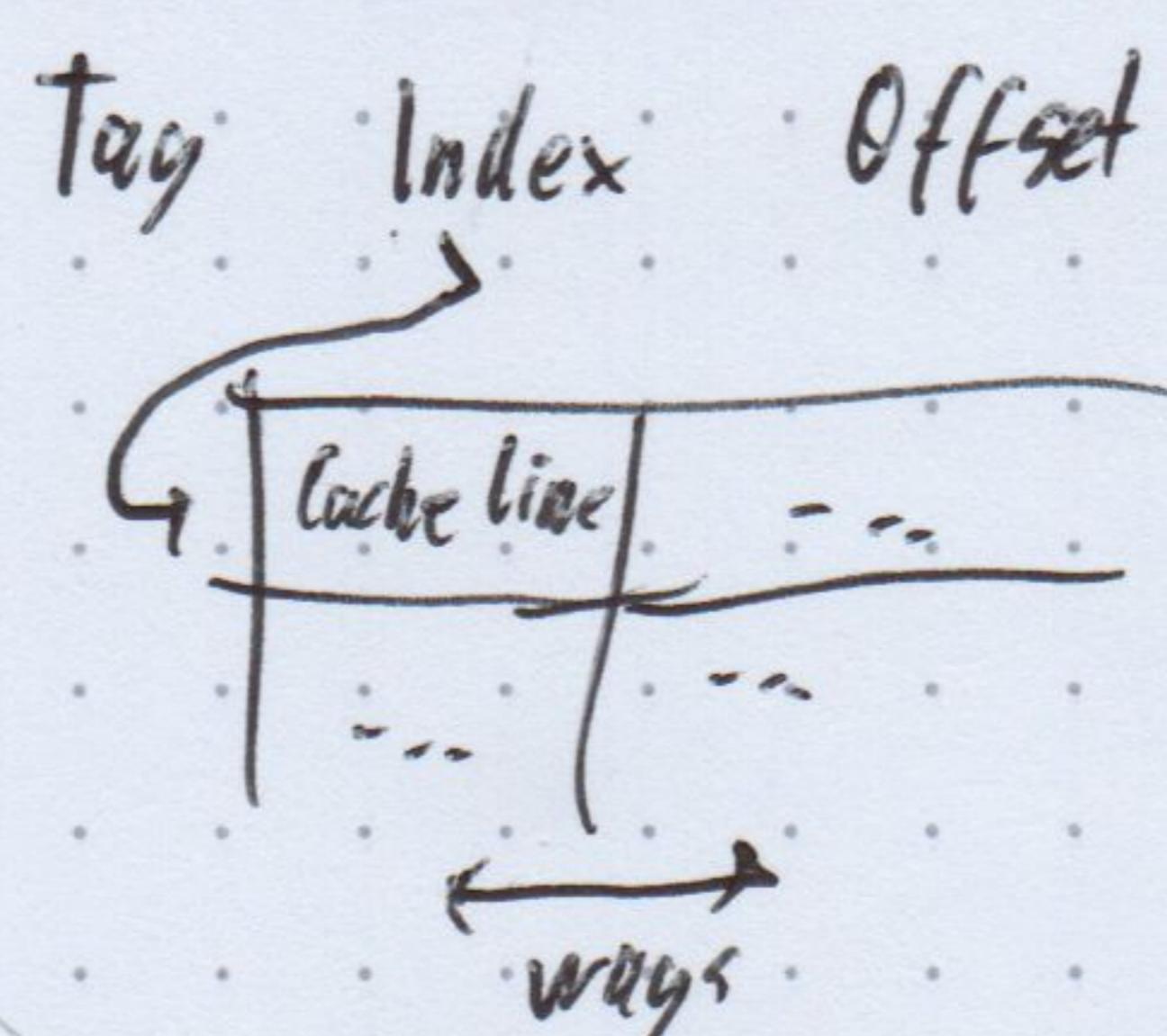
Inclusive ↔ Exclusive Caches

data stored in all cache levels
↳ either core exclusive
or shared

X-way set-associative

Fully associative ↔ Directly-mapped

Cache Architecture



Line:
V Tag DATA

Cache Control Flow

Write-hit Policies

Write-back ↔ Write-through

↳ changed in cache ↳ synchron change in cache
→ needs Dirty bit & mem forwarded to mem on invalidation

Write-miss Policies

Write-allocate ↔ No-write allocate

↳ write-miss leads to exception usually combined
triggering mem read req. with Write-back

Replacement Policies

LRU - Least Recently Used

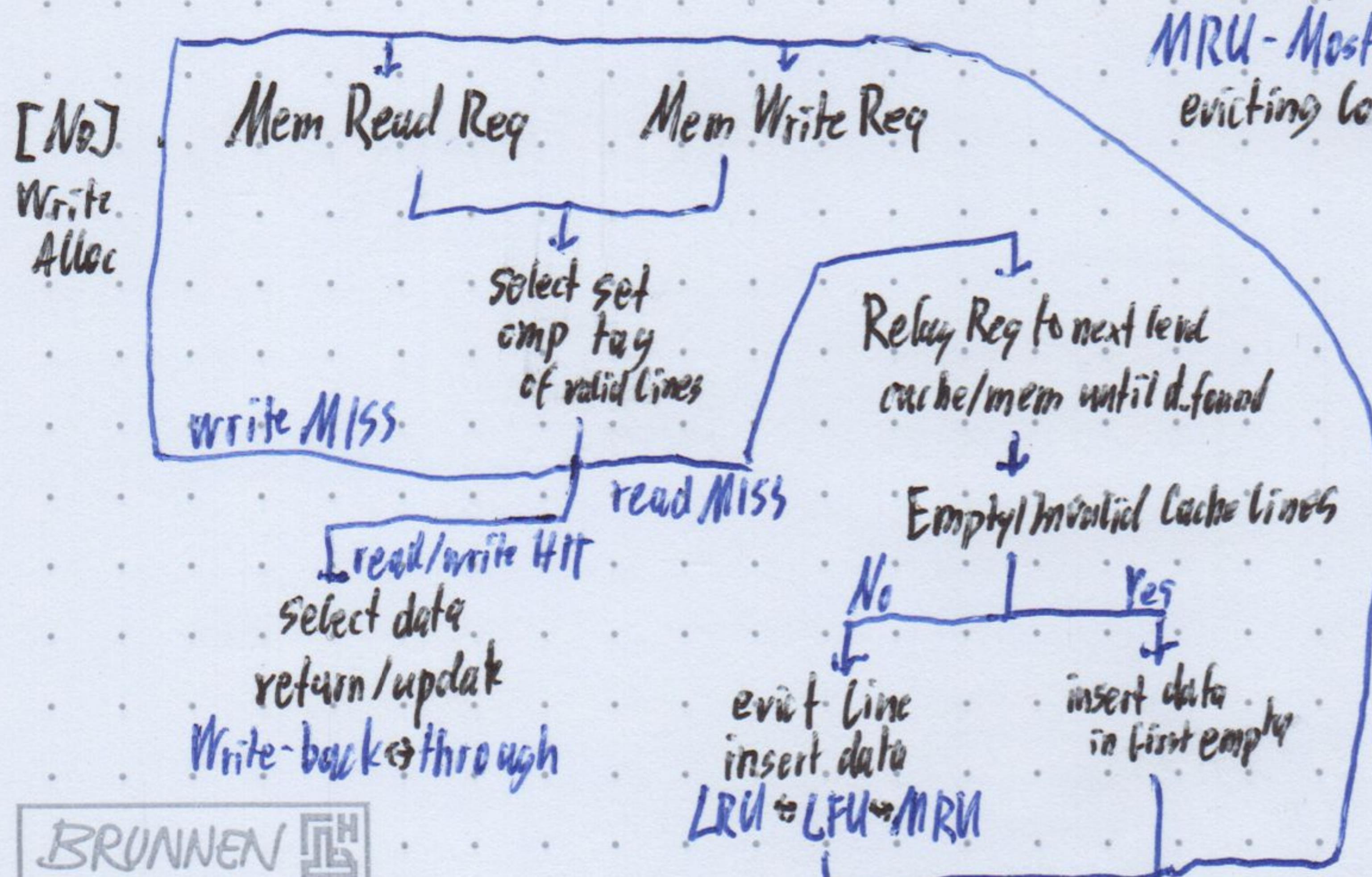
line with oldest age - last time hit performed ↳ pseudo-LRU estimate using less bits

LFU - Least Frequently Used

ctr. how often line used evicting lowest

MRU - Most Recently Used

evicting lowest age



Memory Mappings

Memory-Mapped Files

disk file mapped to phy. mem.

Copy-on-Write

@ change: apply on copied mem. page
+ same view on shared pages

I/O using I/O Ports ~ 8-/16-bit number

access to device control registers

special I/O instr.

IN Reg, Port

OUT Port

Port, Reg

Memory-Mapped I/O

control reg ~ unique mem. addr

+ op's directly on mem. addr (instead of loading in regs)

+ no extra instr.

+ protection via mem. mapping

- needs selective caching

- complicated with multiple buses
(intended device must see)

Hybrid

Separate Mem. Spaces

separately paged, dedicated page tables

e.g. I-space and D-space

DMA - Direct Memory Access

CPU programs DMA Controller

DMA repeatedly instructs device to transfer
byte-by-byte, device acks to addr in mem

DMA interrupts CPU when done

- if CPU much faster

- complexity (embedded systems)

Operating Modes

Word-at-a-time / cycle stealing

occasional short transfers 'stealing' bus

Block / Burst

series of transfers at once

→ blocks bus for long time

Fly-by

DMA instructs device

directly read/write to/from mem

DMA stores heads words itself

+ enables device-to-device

mem-to-mem

internal device buffer

+ checking checksums

+ data at steady rate

but bus volatile

Runtime Attacks

- launched @ exec
- exploit vulnerabilities
- affect benign programs
- 'binaries clean'
- Detection @ runtime expensive

Vulnerability

flaw enabling attacker using system in a manner not intended

Exploit

taking advantage of vulnerability
0-day not publicly disclosed

Denial of Service Attack

[Remote] Arbitrary Code Execution (e.g. system("/bin/sh"), creating user)

privilege escalation

e.g. Buffer-Overflow vulnerability

risk: e.g. override ret.addr. to injected code \Rightarrow arbitrary code execution

fix: bounds check of untrusted inputs

before COPY(buffer[8], *user-input)

Malware = malicious software

e.g. install backdoors

turn into zombie

botnet

controlled with command & control server

Ransomware encrypt disk

extort victim for money

Spyware steal sensitive info

e.g. identity theft

Trojan Horses appearing good, voluntarily downloaded

Viruses reproduces by attaching itself to other programs

exec payload potentially time-delayed

to give time to propagate

Worms spread independently over network

using exploits

bootstrap/loader to download actual worm

on victim system

IOT Malware targeting IOT devices

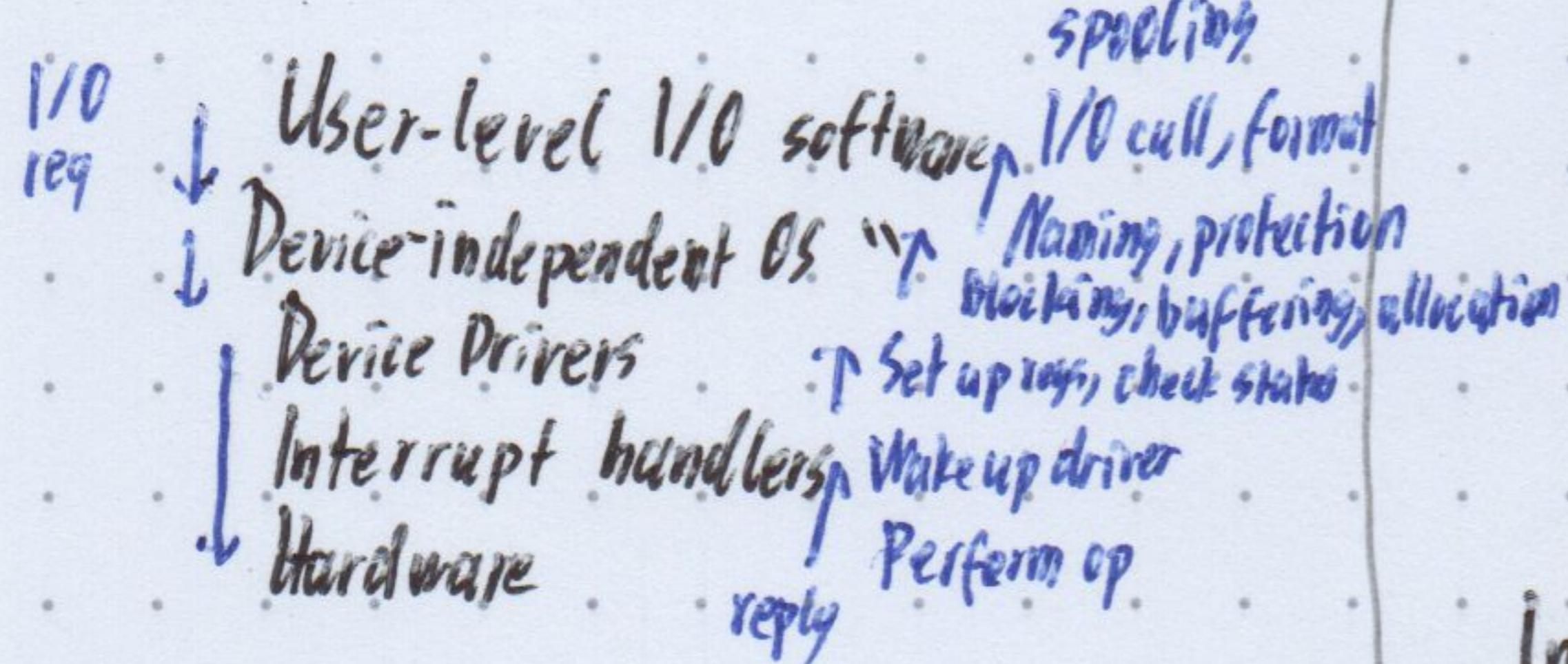
often bad impls/

e.g. Mirai

weak passwords

Device I/O

Layers



Device types

Block devices

1 block at a time

blocks independent of other
(character devices)

stream of characters
not addressable, no seek

Device Controller

electronic component of device

controlled via control reg's

controls device on low level

e.g. serial bit stream → blocks
error-correction

Interrupts

signals on bus

disabled until enabled again

Interrupt Controller

issues prioritized interrupt to CPU

places number on addrs. bus

interrupt-table contains

routing addrs. for if.

CPU acts to allow next

through interrupt

Interrupt Handling

Where to store in interrupted

Internal Regs

can only store fixed amount

Process Stack

Stack pointer no legal value?

page fault?

where to start for page fault routine?

Kernel Stack

switching changes MMU context

invalidates cache and TLB

Processors

pipelined

superscalar

1 op on n instr.

in parallel

Precise vs Imprecise Interrupts

well-defined state

more state

PC saved

instr. before

completed

no instr. beyond

state of instr. of PC

known

may differ per type

Intel x86 provides precise for backward compatibility

Programmed I/O - busy waiting

system call:

```

copy_from_user(buffer, p, count)
while(i < count; i++)
    while(*printer_status_reg != READY),
        *printer_data_reg = p[i];
    return to_user();
    
```

Interrupt-Driven I/O

Interrupt service procedure:

```

if(count == 0) unblock_user();
else { *printer_data_reg = p[i];
    count--;
    i++;
}
ack_interrupt();
return from_interrupt();
    
```

I/O using DMA

Aspects

Device independence

Uniform naming using identifiers

Error handling as low as possible transparent recovery

Synchronous ↔ Asynchronous

Buffering data rate of processing

Handling (including Program Status Word)

1. sure reg's not already by hardware
2. setup context for interrupt service procedure
3. " stack "
4. ack interrupt controller / readable interrupts
5. copy reg's from where saved to process table
6. run interrupt service procedure extract info from device control reg's
7. schedule new process
8. Set up MMU context, load reg's
9. run

Device-Independent I/O Software

typical tasks:

Buffering

Error Reporting

Allocation, Releasing Dedicated devices

Prioritizing device independent block size

Uniform Interfacing for Device Drivers

+ allows new devices without changing OS

Set of functions per class of devices

Driver provides table of pointers

Uniform Naming

UNIX: device's special file i-node contains

major number to locate driver

minor number passed to driver to identify unit

I/O buffering

Unbuffered → interrupt for every character

User-level + reduces interrupts

→ page fault of buffer?

Kernel-level then copied to user space

→ no page fault - buffer full before copied?

Double buffering in kernel for use if:

interrupt while copying filled up buffer to userspace

User Space I/O Software

library provide routines to issue requests

Spooling System for shared access

daemon has exclusive access

p's generate files in spooling directory

Security

Goals

Confidentiality - no exposure of data to unauthorized parties

Integrity - no unauthorized modification

Availability - nobody can make system unusable
denial of service attack

authenticity source of msg, sender's identity

accountability e.g. logs

non-repudiability - can't deny transaction,
privacy validity, msg

Protection Domains - names e.g. user,
set of objects mapped into it processes
e.g. File [ERWX]
Plotter [W]
can be mapped into n. Domains

Protection Matrix

Objects ————— domain
level of access.
(RWX)

typically sparse \Rightarrow stored as

ACL - Access Control Lists

File \rightarrow (domain, level of access)
or Capabilities set of
domain \rightarrow set of capabilities
 $L \sim (\text{object}, \text{level of access})$

Authentication

user knows e.g. password
has smartphone, smartcard
is biometric

UNIX Passwords

password file

username \rightarrow encrypted password

L da von allen usern
einsehbar

Encrypted

da z.B. von Admin einsehbar

Salted

user, salt, e (passwd, salt)

to hinder precomputation Attacks

L list of common passwords and
encrypted

Adversary Model - capabilities of an adversary

in general Turing Machine able to perform arbitrary ops

e.g. for Networks: Dolev-Yao adversary can

intercept any msg

authorized participant (send/receive)

spoof identity

DAC - Discretionary Access Control

Users decide access rights

MAC - Mandatory Access Control + high security (military, hospitals)
system enforces rules

Bell-LaPadula Security Model read down \Rightarrow Confidentiality
write up

Users \sim security levels
objects

simple security/integrity properties

Biba Model no write up
no read down \Rightarrow Integrity

Covert Channels

P passes info to other

by influencing observable system property

e.g. Modulating CPU usage

paging rate

locking files

acquiring devices

Side Channel - info extracted from impl. of computer system

e.g. Cache leakage in shared physical system

Timing info for computation

Power consumption

electromagnetic Radiation

TCB - Trusted Computing Base

hardware & software components necessary to enforce security rules.

typically

Security-relevant hardware

Parts of OS-kernel

Vor programs with superuser

privileges

necessary functionality

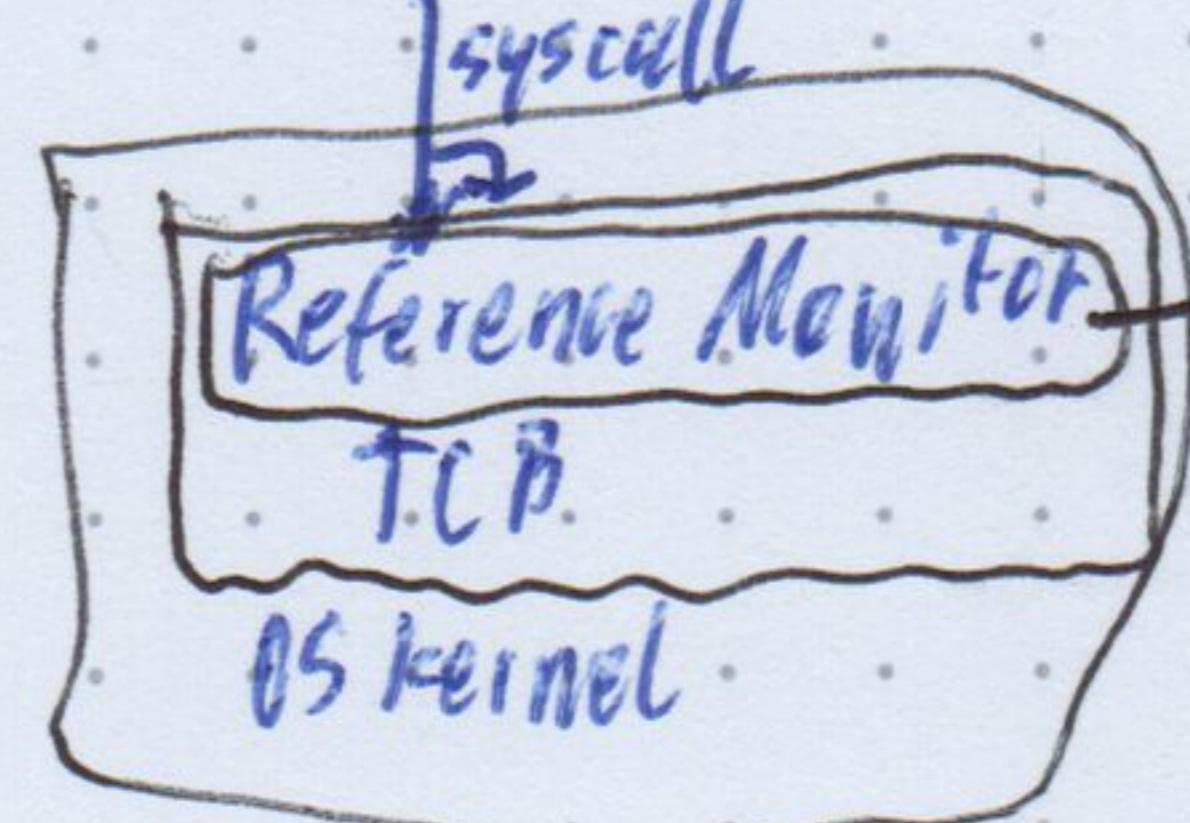
P. creation, switching

Mem. Management

Part of File, I/O Management

User Process

Systemcall



TCB as small as possible

call allowed?

CFG - Control Flow Graph

node ~ Basic block
code sequence
with 1 branch in/out

Code Reuse - adding path to CFG

- limited to available nodes
- needs static analysis of attacked code

Hybrid - e.g. use library code to copy
Malware into area, then making it
executable

Code Injektion - adding new node to CFG arbitrary code (remote console, exploit BS)

but

DEP - Data Execution Prevention

Return-to-libc Attacks Solar Designer 1997

- relies on library (no changes)
- not directly arbitrary exec
- redirect execution to security critical func in shared libraries
 - e.g. libc with open(), malloc(), printf(), system(), execve()
 - e.g. system("sh"), exit()

Cache-based Side-channel Attacks

Primer + Probe - measures if victim overwrote cache line

- Overrides all cache lines

- Noise e.g. durch scheduling

column: least significant nibble of address
row: most

e.g. AES uses substitution box

- key ~ S-box access
- know code ES-box (open)
- know memory mappings
- trigger single encryption

Virtualisation

Pros:

Kernel Kernel	Fault tolerance if system crashes Isolation of sensitive data
VM1 VM2	No dedicated hardware
virtual machine manager	Deployability: all dependencies, right versions allows freezing and starting from new location
hardware	Reliability hypervisor way smaller than OS

Requirements

Safety - hypervisor full control of virtualised resources

Fidelity - behavior VM = as if hardware

Efficiency - little hypervisor intervention

Sensitive instr \Leftarrow Privileged instr. (Popek & Goldberg 1972)

exec diff in user mode
kernel \Downarrow trap in User mode

Trap-and-Emulate

Guest OS thinks it's in kernel mode (virtual kernel mode)

privileged op \Rightarrow trap to hypervisor

First by guestOS: carry out instr
 \Downarrow process: emulate HW response

Binary Translation for sensitive ops. that don't support trapping

rewrite Basic Code Blocks before exec

\Downarrow sensitive op \Rightarrow call to hypervisor

final branch instr.

Guest OS in ring 1 to protect kernel Mem from user processes

Para virtualisation \leftrightarrow full virtualisation

not transparent to guests

Machine-like interface: hypercalls for sensitive ops

+ simpler, faster impl.

Type 1 \leftrightarrow Type 2 hypervisor

\Downarrow process on host OS

on hardware

Older Processors: e.g. popf instr. loads diff flags from stack
 \Downarrow designed for calculator in kernel+user mode

backwards compatibility
Technology: Intel VT, AMD SVM, Secure

Current: 2005

Container for guest OSs.

hardware bitmap controls trapped ops

Trap-and-