# ~~Paging~~
## Introduction

OS → usually takes one of two approaches when solving any space-management problem

chop things up into variable-sized pieces

**Segmentation**

problems:
→ fragmented
↓
allocation challenging over time

chop-up page into fixed-size pieces

**Paging**

the Atlas

Instead of splitting up a
process's address space into
some number of variable-sized
logical segments (code, heap, stack)
ⓧ
divide it into fixed-size units
↓
page

Physical memory
becomes an
array of fixed-size slots
called page-frames

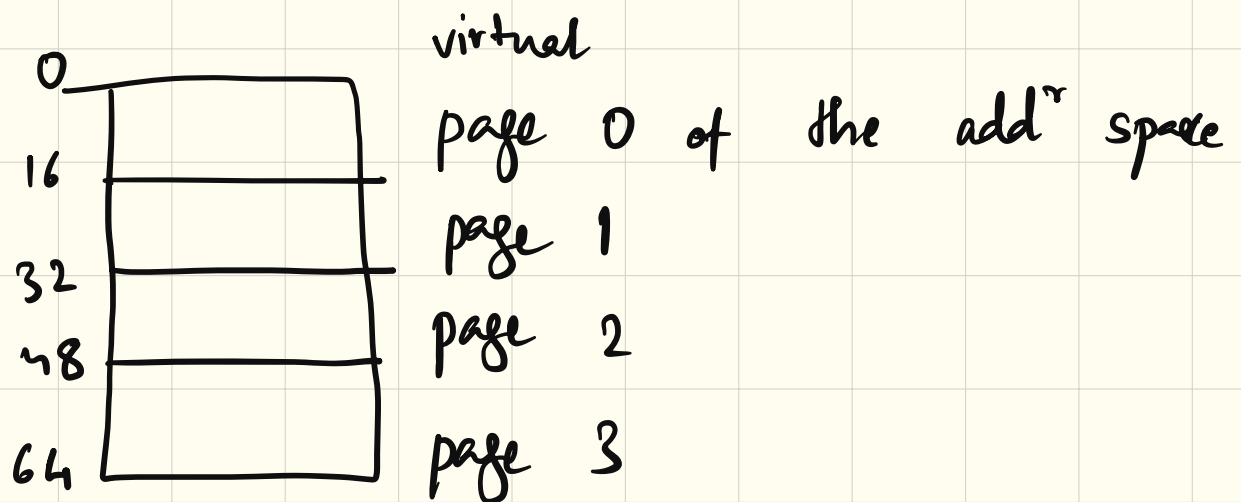VA                    ⊘A
Process               Physical Memory
↓                     ↓

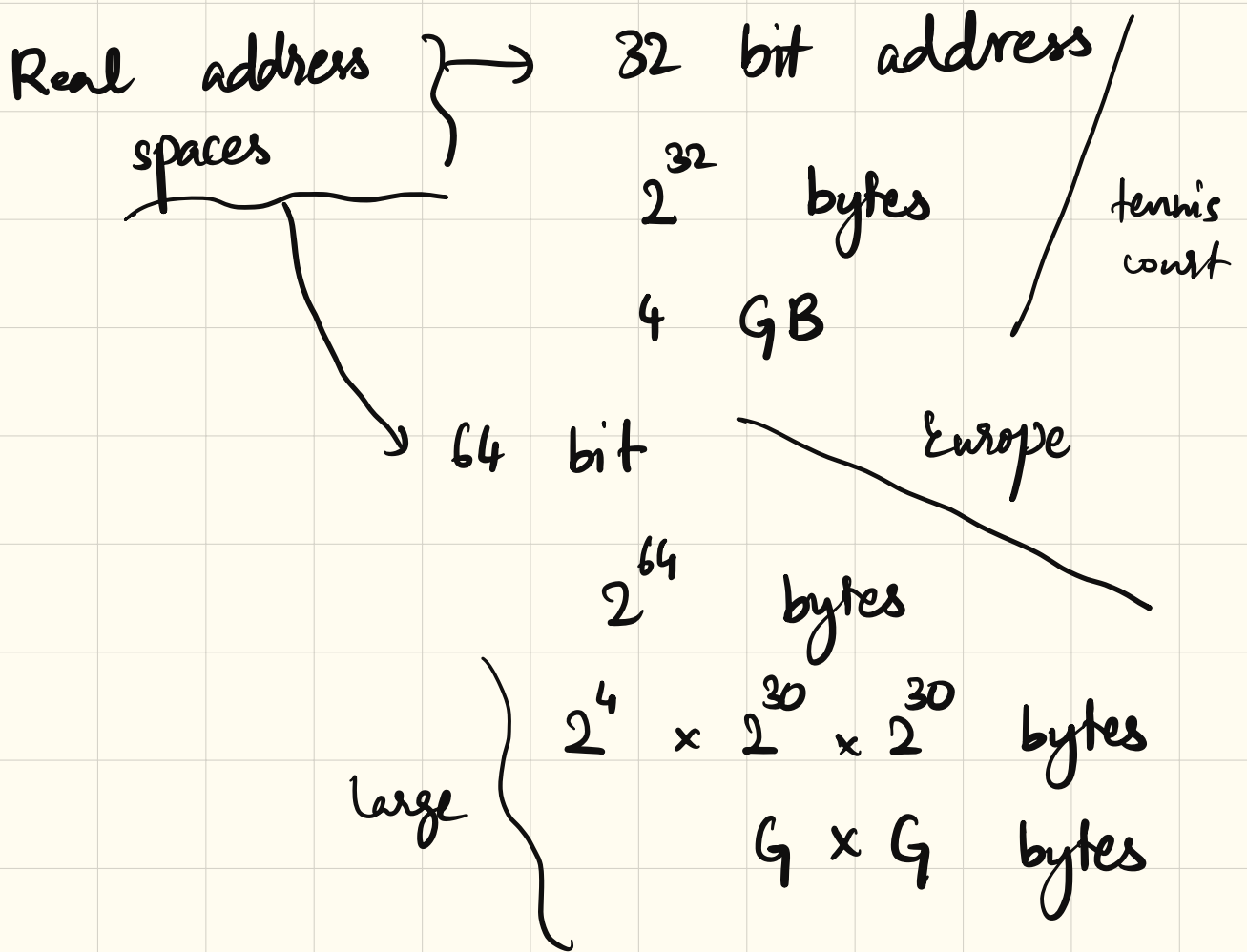Page                  Page frames.

Challenge

How to virtualize memory
 with pages:
  — avoid segmentation problems
  — minimal space and time
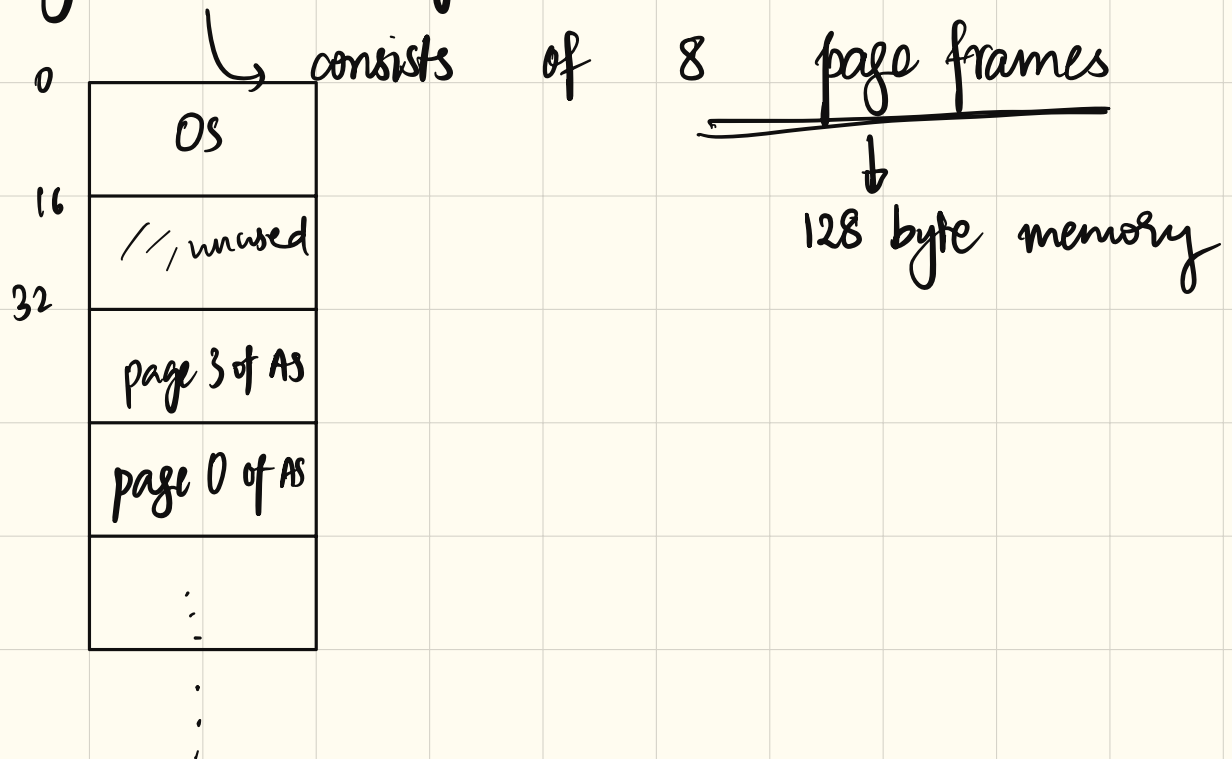          overheads

Simple example and overview.



virtual
page 0 of the add$^r$ space
page 1
page 2
page 3

Simple 64-byte add$^r$ space
 — 4  16-byte pages

Real address $\}\rightarrow$ 32 bit address

spaces

/ tennis
court

$2^{32}$ bytes

4 GB

$\searrow$ 64 bit

Europe

$2^{64}$ bytes

large $\{$ $2^4 \times 2^{30} \times 2^{30}$ bytes

$G \times G$ bytes

Physical memory

$\searrow$ consists of 8 page frames

$\downarrow$

128 byte memory

| | |
|---|---|
| 0 | OS |
| 16 | /// unused |
| 32 | page 3 of AS |
| | page 0 of AS |
| | |
| | ⋮ |

⋮

Most important advantage of paging

↓
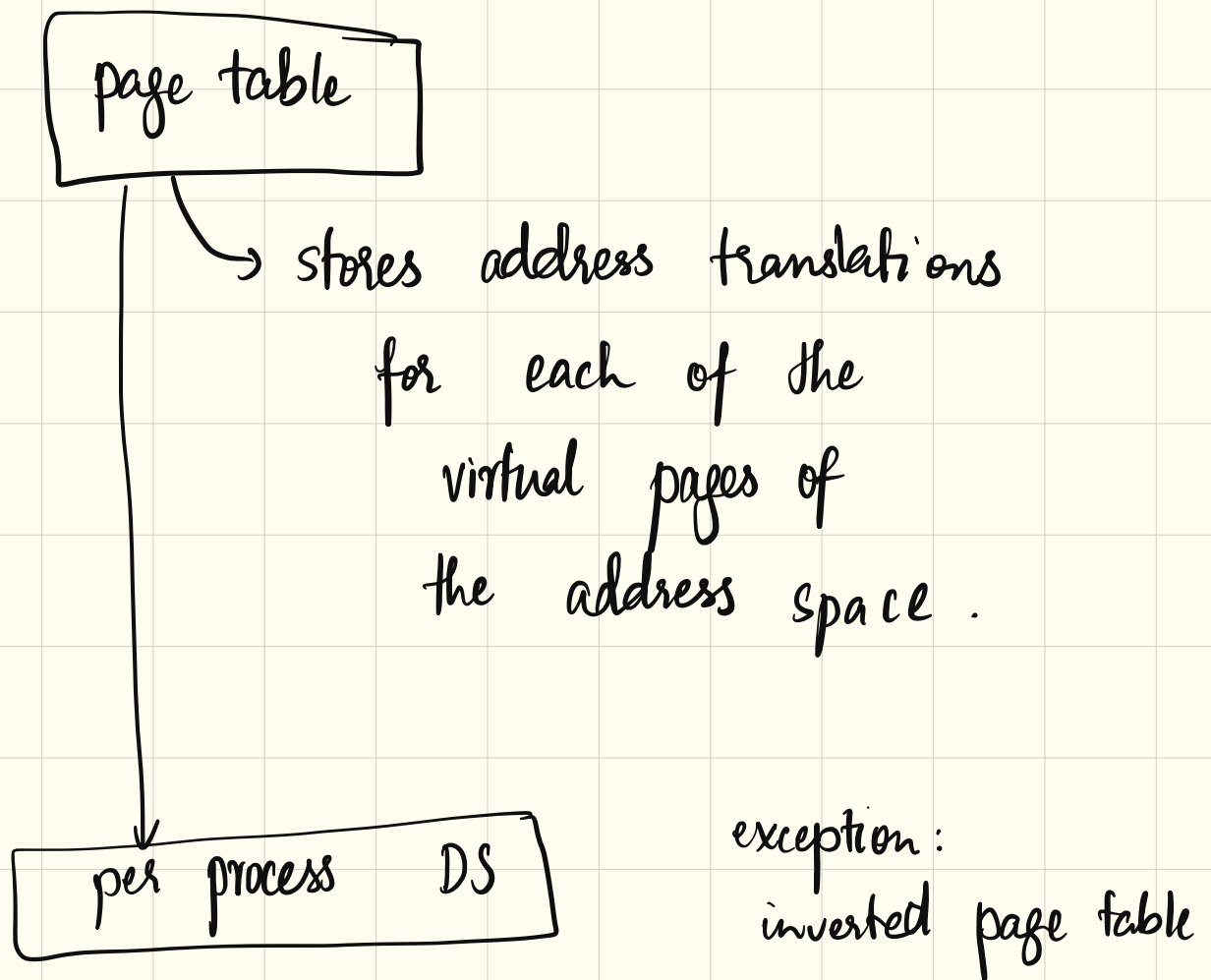
① flexibility

→ support the abstraction effectively

→ we won't make assumptions about direction of stack and heap.

② Simplicity

(maybe)

→ OS keeps a list of free pages to allocate memory.

→ To record where each virtual page of the address space is placed in memory, the OS usually keeps a per-process DS known as page table

page table

→ stores address translations
  for each of the
  virtual pages of
  the address space.

per process   DS

exception:
inverted page table

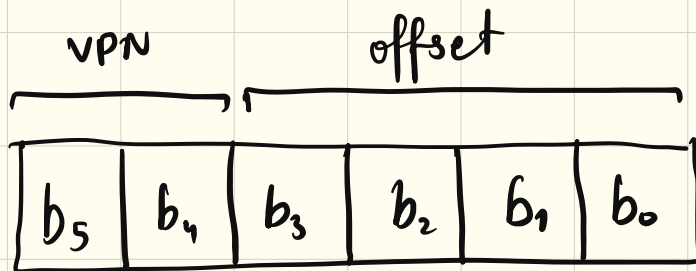OS will manage many such tables.

movl <virtual add$^r$> , % eax.

→ ignore instruction fetch

to translate VA

split it into
① Virtual Page No.
② offset

e.g. VA Size = 64 bytes

$\Rightarrow$ 6 bits for VA

VPN      offset

| $b_5$ | $b_4$ | $b_3$ | $b_2$ | $b_1$ | $b_0$ |
|-------|-------|-------|-------|-------|-------|

page size = 16 bytes = 4 bits = offset

2 bits = select one of 4 pages

movl    21    , %eax

$$\begin{cases} \\ \end{cases} \Downarrow$$

0 1 0 1 0 1

$\downarrow$

page 1 ⸺ use page table $\rightarrow$ find which physical frame

v. page no. 1 resides in.

PFN (Physical Frame #)

| 0 | 1 | 0 | 1 | 0 | 1 |

VPN

↓

Add$^r$ Trans lation

e.g | 1 | 1 | 1 | 0 | 1 | 0 | 1 |

## Questions

① Where are the page tables stored

② What are the typical contents of the page table, and how big are the tables ?

③ Can it be slow?

# Where are the page tables stored

Page tables can get terribly large.

e.g.   32 - bit   address space

4 KB   page size

$4 KB = 2^{12}$ bytes

offset = 12 bits

VPN  = 20 bits

$2^{20}$ translations  ~  million

OS would

have to

manage

per process

assuming   we   need   4 bytes   per

PTE   to   hold   the   physical translation

+ some other stuff

4MB   memory

per process

100  processes $\implies$  400 MB
                           just for
                               translations
                     ———————————————
                        Not  good

For  64 - bit
              $\leadsto$ even  larger
                    gruesome

Page  tables  are  $\longrightarrow$  $\therefore$  not  stored
      so  big                      in  any
                                    special  hardware
                                        $\downarrow$

Page  tabe  $\leadsto$  older        store  in
data  structure          systems
                      ( hardware )    memory
  $\}$                   ( determined )

  $\downarrow$

( modern )
  systems
flexibly  managed
   by  the  OS.

# What's actually in the Page Table?

→ DS to map VPN to PFN

→ simplest : | linear page table |

↓

just an array

→ indexed by VPN

Page table [ VPN ] ← PFN

→ more advanced DS in later chapters

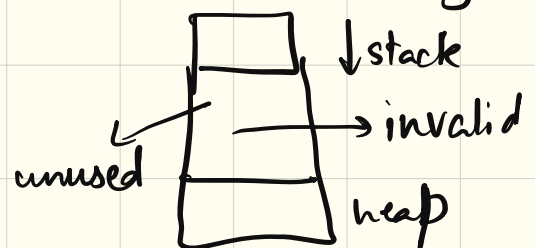## Contents of each PTE (page-table entry)

① Valid bit ————————→ e.g when a program starts running

if process tries to access such address (invalid)

↓

generate trap to the OS

↳ terminates the process

↓ stack

→ invalid

unused

heap

Valid bit $\longrightarrow$ crucial for supporting
sparse address space

→ remove the need
to allocate physical
frames for unused
memory.

② Protection bits

→ read write execute

traps if different permission
than
expected.

③ Present bit

↳ physical memory or disk?
(swapped out?)

}

Swapping — address spaces
larger than
physical memory.

Swapping allows the OS to
free up physical memory
by moving rarely used
pages to disk.

④ Dirty bit

⑤ Reference / accessed bit
↳ whether the page has been
accessed

useful
during
page replacement  → to know which pages are
popular.

⑥

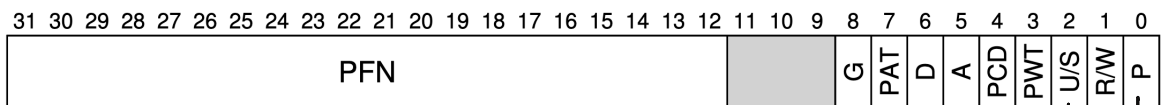| 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 | 11 10 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| PFN | | G | PAT | D | A | PCD | PWT | U/S | R/W | P |

Figure 18.5: **An x86 Page Table Entry (PTE)**

present

hardware cachig    user / supervisor

## Why no valid bit?

P ✓

P = 1 ⟹ both present and valid

P = 0 ⟹ may not be present
      but is valid

         OR

   not present
   not valid

↓

triggers trap to the OS

↓

use add$^n$ structure
it keeps to determine
if the page is valid,
and thus perhaps
should be swapped
back in.

or not
illegal memory

access

# Paging: <u>Also Too Slow</u>

→ Page Tables might be too big

→ They can slow things down too.

movl 21, %eax

→ physical address
117

first fetch the
proper PTE

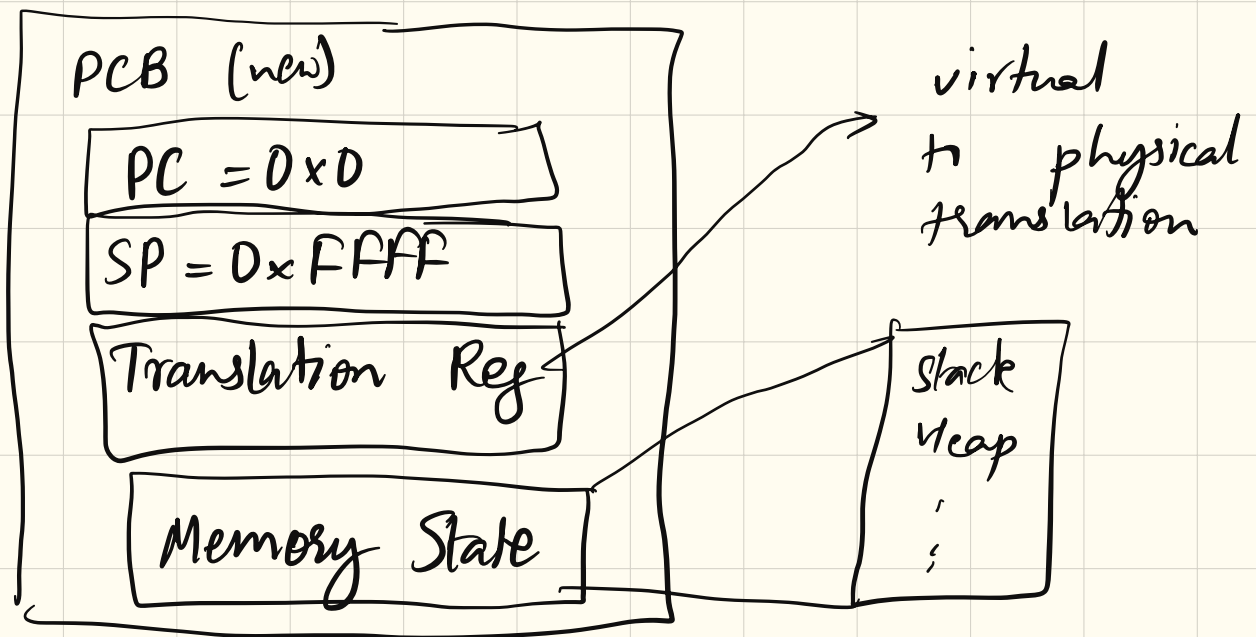---

Slides: important stuff

* Executing program

→ File to process memory view.

* All processes have same add$^r$ space.

→ OS performs memory virtualisation
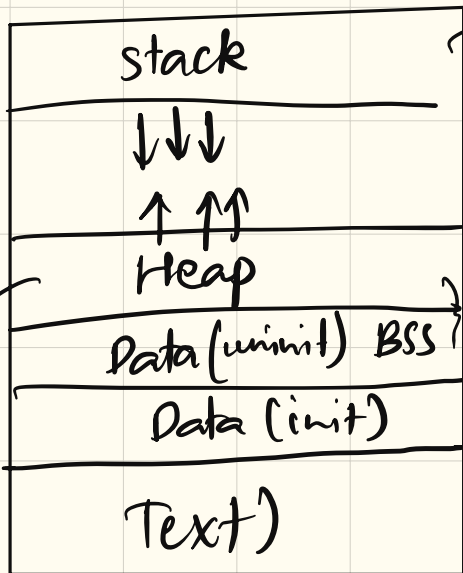with the help of hardware.

\* Responsibilities of OS during program load :

    → Create add$^r$ space →

    → Load binary

    → Update PCB register state

```
┌─────────┐
│  stack  │
├─────────┤
│  Heap   │
├─────────┤
│  unun   │
├─────────┤
│  Data   │
└─────────┘
```

```
┌──────────────────────────────────┐
│  PCB (new)                        │
│  ┌────────────────────────────┐   │
│  │  PC = 0×0                   │   │
│  ├────────────────────────────┤   │
│  │  SP = 0× FFFF               │   │
│  ├────────────────────────────┤   │
│  │  Translation Reg           │   │
│  ├────────────────────────────┤   │
│  │                            │   │
│  │  Memory State              │   │
│  └────────────────────────────┘   │
└──────────────────────────────────┘
```

virtual to physical translation

```
┌──────────┐
│  stack   │
│  Heap    │
│    ;     │
│    ;     │
└──────────┘
```

## Memory APIs

→ User has no direct control on physical memory.

```
┌────────────────────┐
│       stack        │
├────────────────────┤
│       ↓↓↓          │
│                    │
│       ↑ ↑↑         │
├────────────────────┤
│       Heap         │
├────────────────────┤
│ Data (uninit) BSS  │
├────────────────────┤
│    Data (init)     │
├────────────────────┤
│                    │
│      Text)         │
└────────────────────┘
```

no special system
calls

BSS → brk, sbrk

input
void* address
↓
make end
of BSS
= addr

input
long
size
make
end
= end + size

return old
addr

sbrk(0)
↳ returns the
current
location of
BSS

* MAP_FIXED
NULL ⇒ remap things

mmap: discontiguous
allocation

e.g: allocate 4096 bytes
w/ read + write
                  len   prot
prt = mmap (NULL, 4096, PROT_READ)
       start addr      PROT_WRITE,

MAP_ANONY | ... , -1, 0 }
   flags           fd (??) offset

→ etext → end of Text
  edata → end of data (initialized)
  end → end of BSS

} at
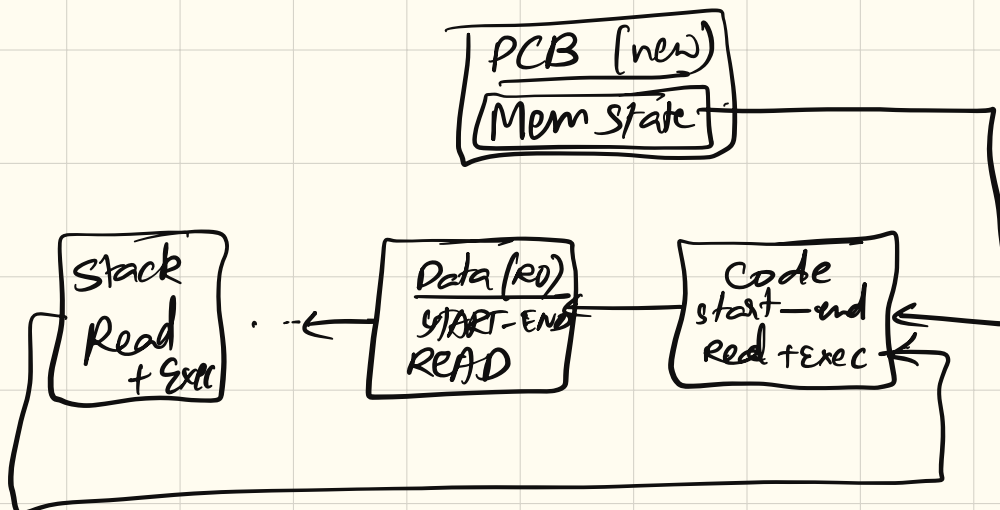program
load
time

→ printing adds of $f^n$ and
    vars:

    Linux :    /proc/ \<pid\> / maps


## Example: Memory State of PCB

→ Circular list

   START and END never overlap b/w
   two segment areas.

→ Can be merged/extended if
  permission matched.

```
        ┌──────────────┐
        │ PCB (new)    │
        │ Mem State ───┼──────────────────┐
        └──────────────┘                  │
                                          │
┌─────────┐     ┌──────────┐     ┌─────────────┐
│ Stack   │     │ Data (RO)│     │  Code       │
│ Read    │◄-- -│ START-END│     │ start—end   │
│ + Exec  │     │ READ     │◄────│ Read +Exec  │◄─┘
└─────────┘     └──────────┘     └─────────────┘
     └──────────────────────────────────┘
```

\* **Inheriting Address Space through fork ( )**

→ child inherit the <u>memory state</u> of the parent

↳ data structure is copied into child PCB

→ Any change through mmap( ) or brk( ) is per-process

\* **Overriding Address Space through exec'**

→ Add$^r$ is reinitialized using the new executable

→ Changes to newly created AS depends on the logic of new process
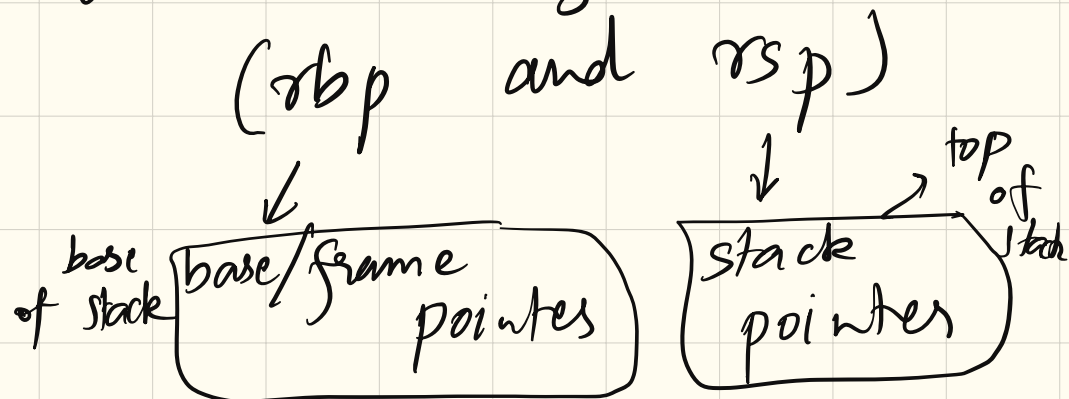
# Address Space Granularity

* Translation at address space granularity

→ ISA, x86

| | |
|---|---|
| register | mov %rcx, %rax |
| immediate | mov $5, %rax |
| absolute | mov 8000001, %rax |
| indirect | mov (%rcx), %rax |
| displacement | mov -16(%rbp), %rax |

→ examples

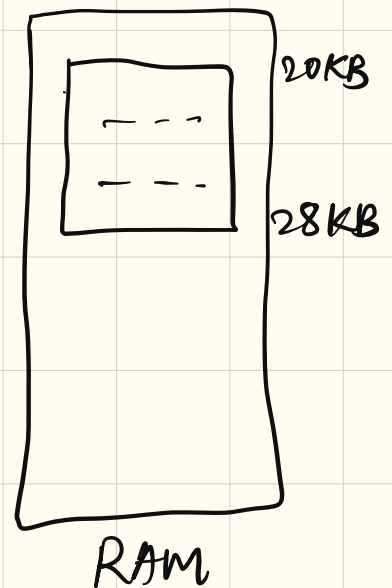→ compiler: does not know stack address, hence blindly uses the register (rbp and rsp)

base of stack → base/frame pointer

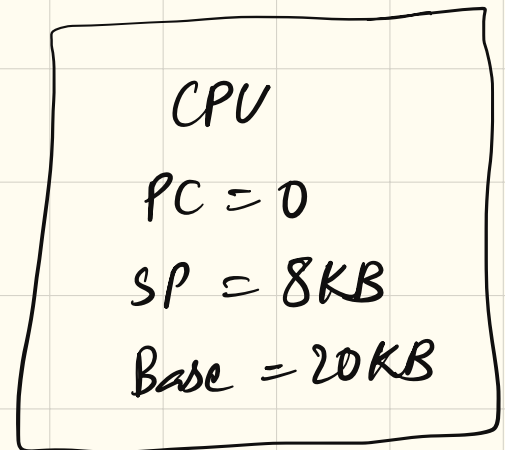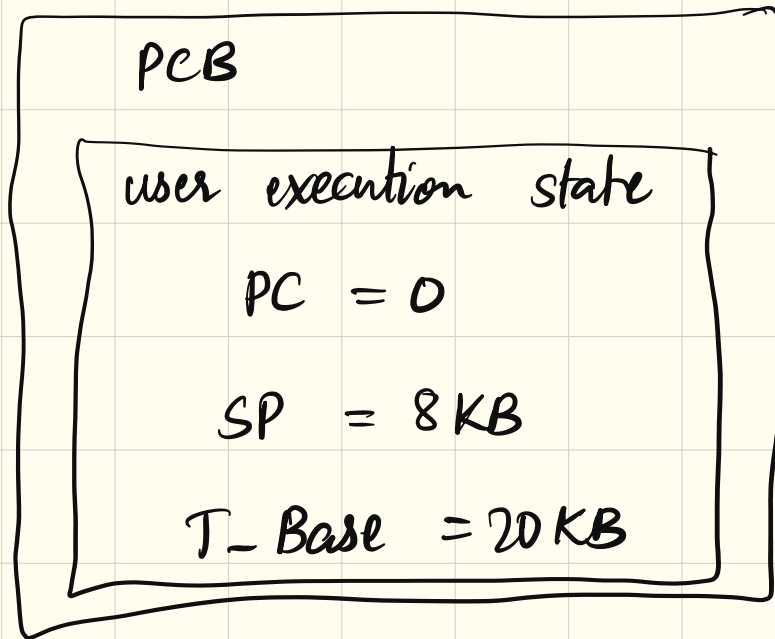stack pointer → top of stack

OS during binary load

```
load_new_executable ( PCB *current, File *exe){
    verify (exe )
    reinit_addr_space (current → mm_state);
    allocate_phys_mem (current )
    load_exe_to_phys_mem (current, exe);
```
VAS { 
```
    set_user_sp ( current → mm_state →
                            stack_start);

    set_user_pc (current → mm_state →
                            code_start)
```

```
    return to user

}
```

Process State After exec ( )

→ OS configures the base register
    depending on the physical loc^n

PCB

user execution state

PC = 0

SP = 8 KB

T_Base = 20 KB

CPU

PC = 0

SP = 8KB

Base = 20KB

20KB

28KB

RAM

instr fetch ( vaddr = 10)

→ instr fetch (paddr = 20KB + 10)

" push %rbp "

→ Assuming RSP = 8KB

" push %rbp " → results in a

memory store at

addr (8KB − 8)

CPU translates the addr to
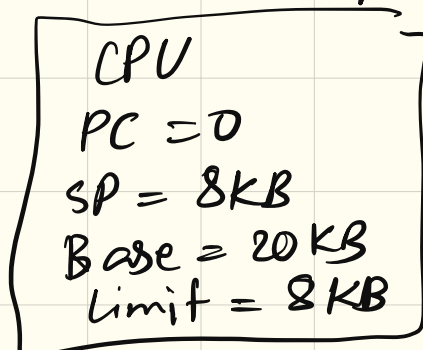
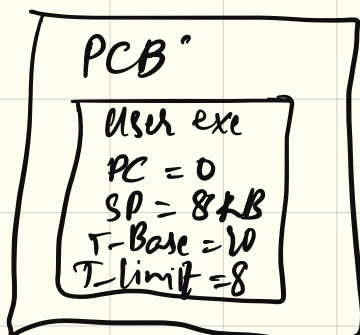(28 KB − 8)

$\rightarrow$ How to stop illegal accesses?

VA size = 8 KB

Access = 20 KB

$\downarrow$

Phy Add$^r$ = 40 KB ⊗

$\rightarrow$ Limit Register

↳ can be changed
from privileged
mode

$\rightarrow$ Hardware raises a fault if the program violates the limit

OS fault handler may kill the process.

PCB°

| |
|---|
| User exe |
| PC = 0 |
| SP = 8kB |
| r-Base = 20 |
| T-limit = 8 |

CPU

| |
|---|
| PC = 0 |
| SP = 8KB |
| Base = 20 KB |
| Limit = 8 KB |

∴ How is memory isolation achieved?
→ limit register.

## Context switch
→ Base and limit register
values saved in the outgoing
process PCB during context
switch.

→ Loaded from PCB to the CPU
when a process is scheduled

## Disadvantages of Translation at AS Granularity
→ Physical memory must be
greater than AS size.
→ against abstraction philosophy

→ Memory inefficient

→ Physical memory size is same as addr space size irrespective of actual usage

→ Memory wastage

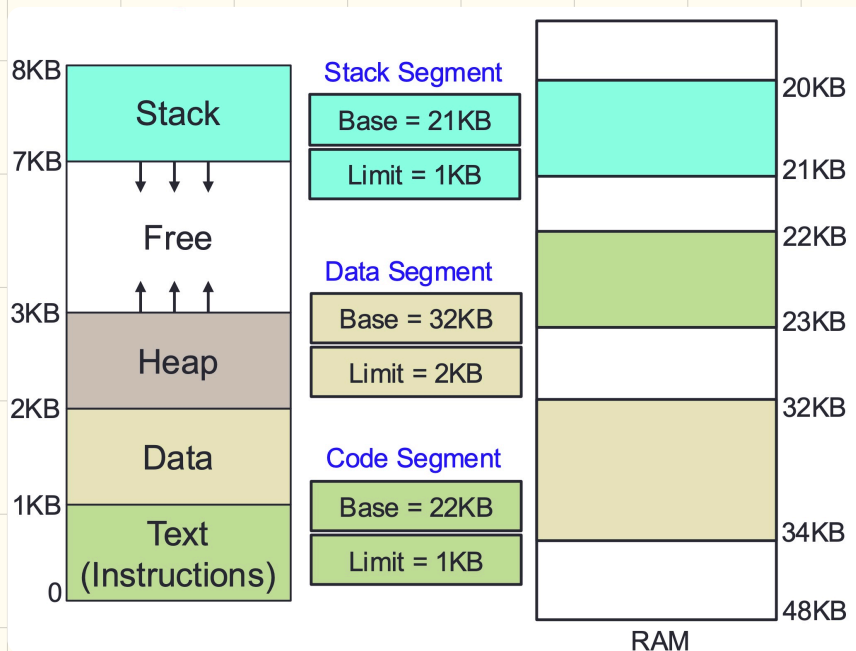→ Degree of multiprogramming is very less.

## Segmentation

→ Extension of the basic sheme w/ more base – limit register pairs.

→ Ex: code add$^r$
data add$^r$
stack add$^r$

# Segmentation : Explicit Addressing

→ Part of the code is used to explicitly specify segments.

→



| 8KB | | Stack Segment | | 20KB |
|---|---|---|---|---|
| | Stack | Base = 21KB | | |
| 7KB | | Limit = 1KB | | 21KB |
| | ↓ ↓ ↓ | | | 22KB |
| | Free | Data Segment | | |
| 3KB | ↑ ↑ ↑ | Base = 32KB | | 23KB |
| | Heap | Limit = 2KB | | |
| 2KB | | | | 32KB |
| | Data | Code Segment | | |
| 1KB | | Base = 22KB | | 34KB |
| | Text (Instructions) | Limit = 1KB | | |
| 0 | | | | 48KB |

RAM

$VA = 8\,KB$ ; $\text{add}^n$ length $= 13$ bits ; 3 segments

→ Two MSBs used to specify the segment :

$00 \rightarrow$ code

$01 \rightarrow$ data

$11 \rightarrow$ stack

→ hardware selects the segment register based
on the value of 2 MSB
bits and the rest of the bits
are used as offset.

→ Max size of each segment
$= 2KB$

→ Physical allocation is still done on an on-demand
basis

## Disadvantages with explicity addressing

→ Inflexible

→ Data and Stack cannot
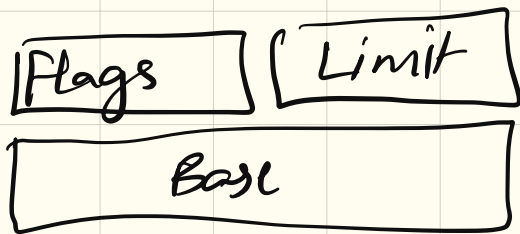be sized dynamically

→ Wastage of VA space,
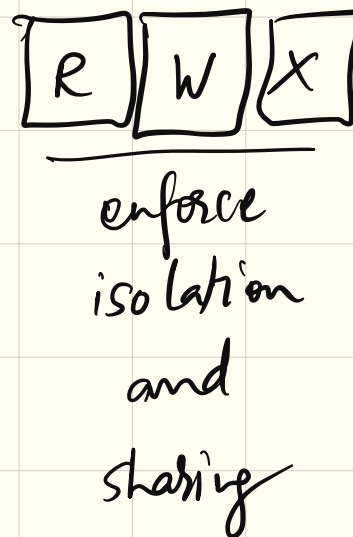
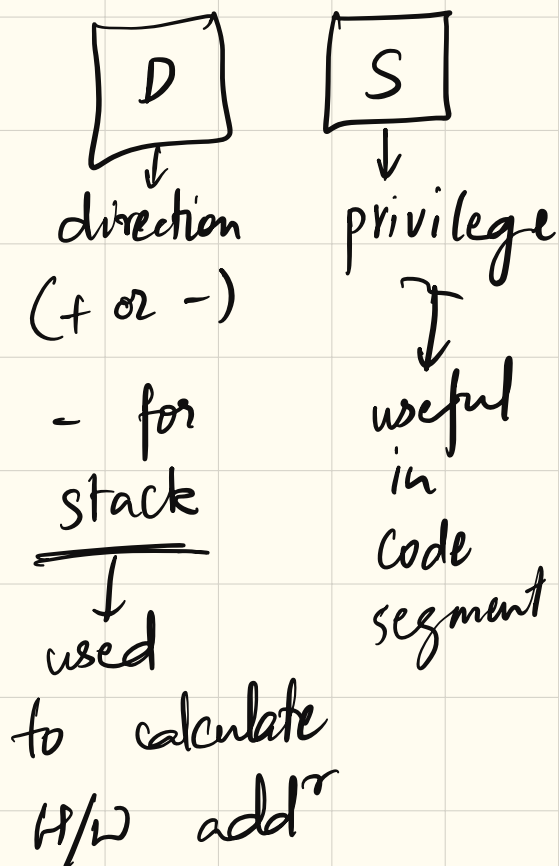→ in our example, 2KB VA
is
unusable

# Segmentation : Implicit Addressing

→ Hardware selects the segment register based on the operation

→ Code segment for instruction access
  → fetch $addr^r$, jump target, call $addr^r$

→ Stack segment for stack operations
  → push pop, indirect addressing with SP, BP

→ Data segment for other addr

# Segmentation: (Protection and Direction

```
| Flags |  | Limit |
|     Base       |
```

Flags

```
[ D ]   [ S ]       [ R ][ W ][ X ]
```

direction   privilege        enforce

(+ or -)                     isolation

  - for      useful          and

  stack       in            sharing

   |          code

  used        segment

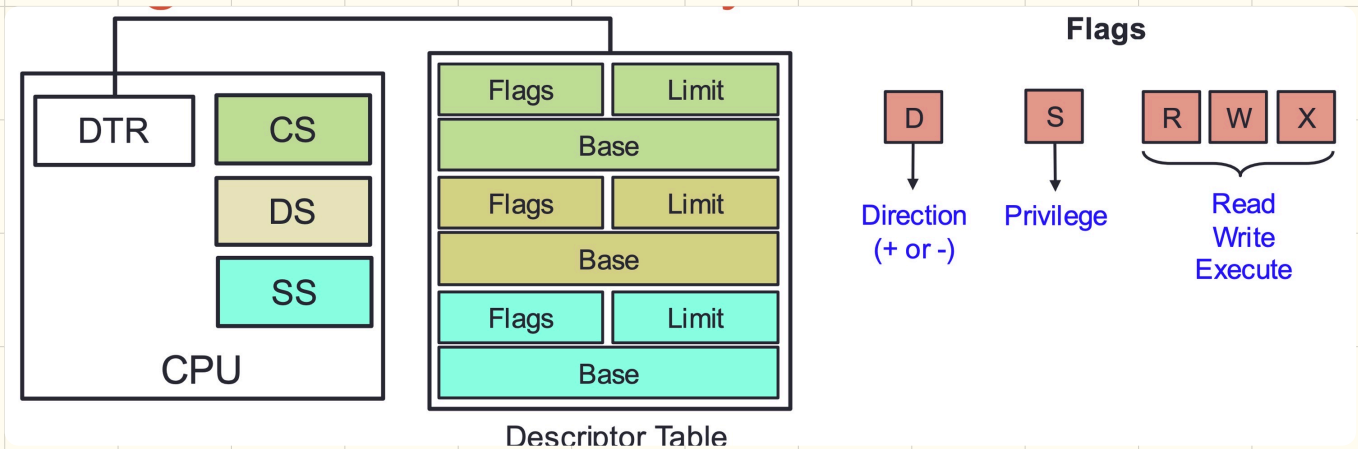to calculate

H/W addⁿ

# Segmentation in Reality

→ DTR : descriptor table register
            is   used  to   access
                       descriptor table

→ # descriptors depend on architecture
→ Separate descriptors for user and
                              kernel mode



| | | |
|---|---|---|
| DTR | CS | Flags / Limit |
| | DS | Base |
| | SS | Flags / Limit |
| CPU | | Base |

Descriptor Table

**Flags**

| D | S | R | W | X |

Direction (+ or -)    Privilege    Read Write Execute

# Advantages of Segmentation

→ Easy and efficient addr translation
→ Save memory wastage for unused
    addresses

## Disadvantages

→ External fragmentation

→ Cannot support discontiguous
  sparse mapping

---

## Paging

→ addresses external fragmentation
     due  to  variable
       size  segments

→ allows discontiguous sparse mapping


→ Basic idea :

① Partition the address space into
   fixed size blocks (call it pages)

② Partition physical memory in a
   similar way (call it page frames)

③ OS :  page ⇄ page frames
              mapping

④ H/W uses the mapping to
   translate VA to PA.

→ Paging example (Pages)

$$32 \times 1024$$
$$2^5 \quad 2^{10}$$

VA size = 32 KB,

Page size = 256 bytes ↝ $2^8$

Add^r length = 15 bits
$$\{0 \times 0 - 0 \times 7FFF\}$$

\# of pages = $\dfrac{32 \text{ KB}}{\frac{1}{4} \text{ KB}}$

= 128

Example: For VA = 0x0510

page number = 5

offset = 16

| 7 bits | 8 bits |
|--------|--------|
| Page # | Offset |

# Paging example ( Page frames )

Physical Addr size = 64 KB
   Addr length = 16 bits {0x0 − 0xFFFF}
   # of PFNS = 256

e.g: For PA = 0x 1F51

PFN = 31         offset = 81

Page Table Mapping:
→ each entry in page table is
      called a Page - Table Entry
                ( PTE )

# Page Table Walk
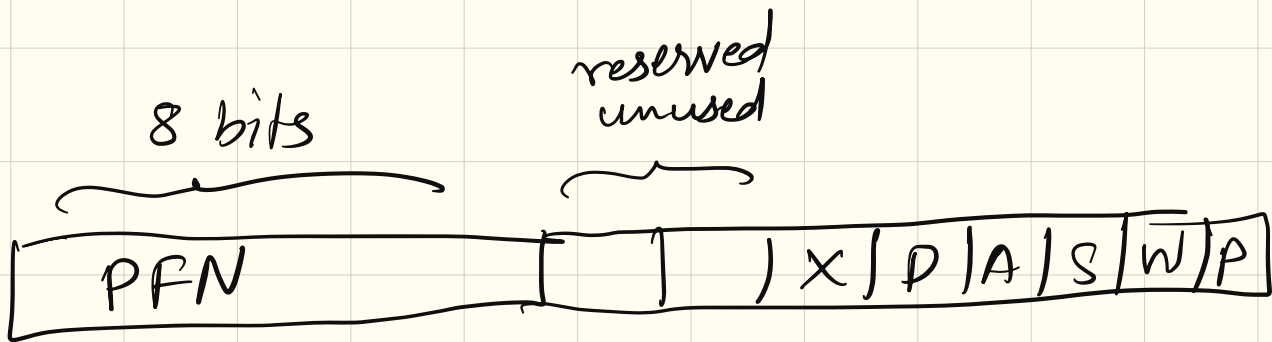
```
PTW ( vaddr V, PTable P)
  // returns physical address
{
    Entry  = P[ V >> 8 ]
    if ( Entry. present )
       return ( Entry. PFN << 8 )
                    + ( V & 0x FF ) ;
    Raise   Page Fault;
}
```

Where is page table stored?

→ Page table is stored in the RAM.
   Page table register (CR 3 in x86) contains the addr.

# Structure of PTE

PFN | 8 bits

reserved unused

| PFN | | | ) X | D | A | S | W̄ | P |

P → present bit → entry is valid

W → write bit → write allowed

S → privilege bit, 0 → only kernel mode access

A → accessed bit, → set by H/W during walk

D → dirty bit → set by H/W during walk

X → execute bit (instruction fetch)

---

Maximum size of physical memory

→ PFN → 8 bits

page size = 256 bytes = $\frac{1}{4}$ KB

$$\therefore \quad \text{Max RAM size} = 2^8 \times \frac{1}{4} \text{ KB}$$

$$= 64 \text{ KB}$$

## Address Translation: Multilevel Paging and TLB

→ With increased addr space size, single level page table entry is not feasible:

→ ↑ page size → fragmentation ↑

→ Small pages may not be suitable to hold all mapping entries.

\* Example: Paging with 32-bit AS
→ Sol$^n$: Multi-level page table.

# Two-level Page Tables (32-bit VA)

→ Level-1 page table contains entries
   pointing to level-2 page table
                              structures

→ L2 entry contrains PFN + flags

10 bits      10 bits           12 bits   4KB block
                                                   size

L1-offset                              Page offset

L2 entry        PFN                Physical
                                   frame
                                   4KB

CR3

→ 4 level page table: 48-bit VA (x86_64

# 4-Level Page Tables: 48-bit VA (x86_64)

Virtual Address (48-bits)

| 9 bits | 9 bits | 9 bits | 9 bits | 12 bits |
|--------|--------|--------|--------|---------|

pgd_offset   pud_offset   pmd_offset   pte_offset   Page offset

pgd_t   pud_t   pmd_t   pte_t   Physical frame (4KB)

CR3

❑ Virtual address size = $2^{48}$, Page size = PF size = 4KB

❑ 4-levels of page table, entry size = 64 bits

# 4-Level Page Tables: Example Translation

| 9 bits | 9 bits | 9 bits | 9 bits | 12 bits |
|--------|--------|--------|--------|---------|
| 000000000 | 000000110 | 000000000 | 000000001 | 000000001000 |

0x2007000      0x2008000      0x200B000      0x200C000      0x640E000

$0^{th}$ 0x2008027

$6^{th}$ 0x200B027

$0^{th}$ 0x200C027

$1^{st}$ 0x640E007

$8^{th}$ User data

0x2007000

CR3

❑ Virtual address size = 0x180001008

❑ H/W translation by repeated access of page table stored in physical memory

❑ Page table entry: 12 bits LSB is used for access flags

# Translation efficiency

Consider 4-level page table, how many memory accesses required for translation?

sum = 0;

for (ctr = 0; ctr < 10; ctr++)

    sum += ctr

Assume no cache

```
0x20100: mov $0 %rax;
0x20102: mov %rax, (%rbp);          //sum=0
0x20104: mov $0, %rcx;              //ctr=0
0x20106: cmp $10, %rcx;←            //ctr<10
0x20109: jge 0x2011f;←              //jump if>=
0x2010f: add %rcx, %rax;←
0x20111: mov %rax, (%rbp);←         //sum+=ctr
0x20113: inc %rcx;←                 //ctr++
0x20115: jmp 0x20106;←              //loop
0x2011f: ret;
```

Instruction execution:  loop = 10 × 6

               others = 5

Memory accesses during translation $= 65 \times 4 = 260$

Data / stack access : initialization = 1

loop = 10

Memory access
during translation   = 11 * 4 = 44

Distinct accesses ?

Assume : stack address range

0x 7FFF000 - 0x 8000000

One code page (0x 20)

and one stack page 0x 7FFF

→ cache these : TLB
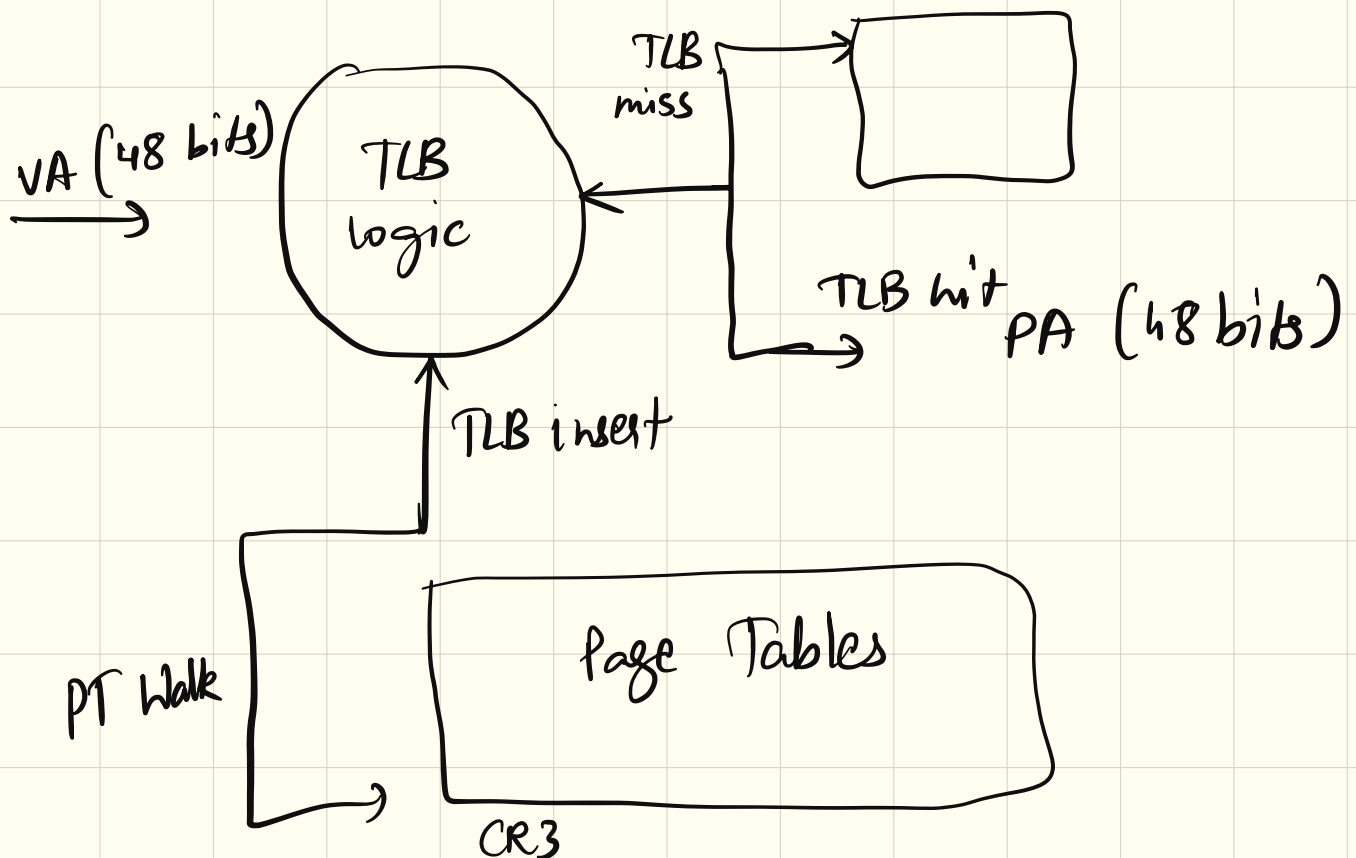
# Paging w/ TLB: Translation efficiency

TLB

| Page | PTE |
|------|-----|
| 0x20 | 0x750 |
| 0x7FF | 0x890 |

```
Translate ( vaddr V ) {
    PageAddress   P  =  V >> 12;
    TLBEntry   =   lookup (P);
    if (entry. valid)
        return   entry. pte;
    entry  =  PageTableWalk (v);
    MakeEntry (entry);
    return   entry. pte;
}
```

TLB → Translation Lookaside Buffer

     → hardware cache

     → stores Page to PFN
                         mapping

→ after first miss for instr fetch, all
                              accesses hit

## Address Translation (TLB + PTW)

VA (48 bits) → **TLB logic** → TLB miss → [ ]

TLB hit → PA (48 bits)

TLB insert ← Page Tables

PT Walk → Page Tables

CR3

→ Seperate TLBs for instruction an
data , multilevel TLBs

→ In x86 , OS cannot make
entries into the TLB directly,
it must flush entries

* How is TLB shared across
multiple processes?
Option 1: flush whole TLB ⟶ performance
problem

Option 2: Address Space Identifier
(ASID) along with each
TLB entry to identify
the process.

| Process (A) | Process (B) |
|---|---|

| ASID | Page | PTE |
|:---:|:---:|:---:|
| A | 0x200 | 0x200007 |
| B | 0x200 | 0x300007 |
| A | 0x201 | 0x205007 |
| B | 0x201 | 0x305007 |

**TLB**

Why is page fault necessary?
→ Page fault is required to support memory over-comitment through lazy allocation and swapping.

# Page fault Handling in x86

```
if ( ! pte. valid ||
            (access == write && !pte. write)
         || (cpl != 0   && pte. priv==0))
         {  CR2 = address;
            errorcode = pte. valid
                          | access << 1
                          | cpl << 2;
         Raise Page Fault ;
      } // simplified
```

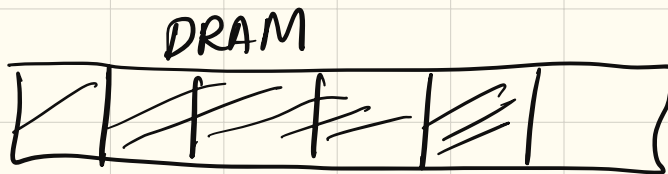cpl = current priviledge level

pte. priv = priv level of pte.


→ error code is pushed into the
   kernel stack by the hardware (x86)

## OS Fault Handler

```
Handle Page Fault (u64 address, u64 error_code)
{
    entry = P[V >> 8];
    if (AddressExists (current -> mm_state,
                            address) &&

            AccessPermitted (current -> mm_state,
                                error_code)) {

            PFN = allocate_pfn();
            install_pte (address, PFN);
            return;
    }
    Raise Signal (SIGSEGV);
}
```

virtual (annotation above "address")

# Page Fault and Swapping

## Swapping (swap_out)

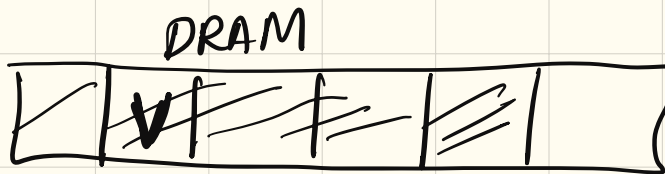DRAM



SWAP ( Hard disk)

Allocate PFN():

free PFNs ↓

cannot break promise to applications,

∴ swap out, but which one?

Decide using a page replacement
policy

DRAM

```
┌──┬──┬──┬──┬──┬──┐
│ │V │ │ │ │ │
└──┴──┴──┴──┴──┴──┘
```
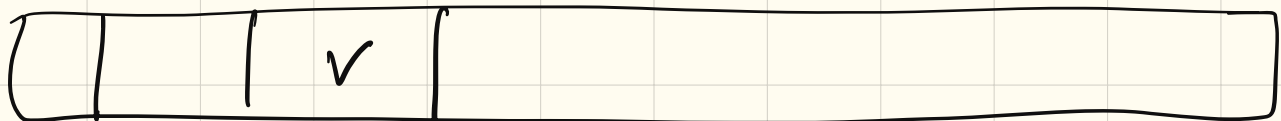
┌─────────────────────┐
│ Page Replacement    │
│           Policy    │
└─────────────────────┘

```
┌────┬────┬────┬────┬────┬────┬────┬────┬────┐
│    │    │    │    │    │    │    │    │    │
└────┴────┴────┴────┴────┴────┴────┴────┴────┘
```

SWAP ( Hard disk)

→ Update the present bit of
   victim entry to 0
        → page fault
Replace PFN with swap Address

Befor

Present

| PAN (v) | - flags | 1 |

After

↓

| Swap Address (v) | same flags | 0 |

| | | ✓ | |

Hard disk

Any future translation will result in
page fault.
        Page Fault Handler will copy
                it back from the
        swap device

## Swap - in :

```
Handle Page Fault ( u64  address  ,  u64  error_code)
                            virtual
{
      entry  =  P [ V >> 8 ];
      if  (AddressExists (current -> mm_state,
                                     address)  &&

           Access Permitted (current -> mm_state,
                                     error_code) ) {

           PFN =   allocate_pfn ();
           if  ( is_swapped_pte (address))
                  swap in (getPTE (address), PFN);

           install_pte (address,  PFN);
           return ;

      }
      Raise Signal (SIGSEGV );
}
```

# Page Replacement

Objective: Minimize # of page faults
(due to swapping)

→ 3 parameters
    → A given sequence of accesses
        to virtual pages
    → # of memory pages (frames)
    → page replacement policies

→ Metrics to measure effectiveness
    → # of page faults
    → Page fault rate ⟶ fraction of
    → AMAT              memory
                            accesses that
                            result in
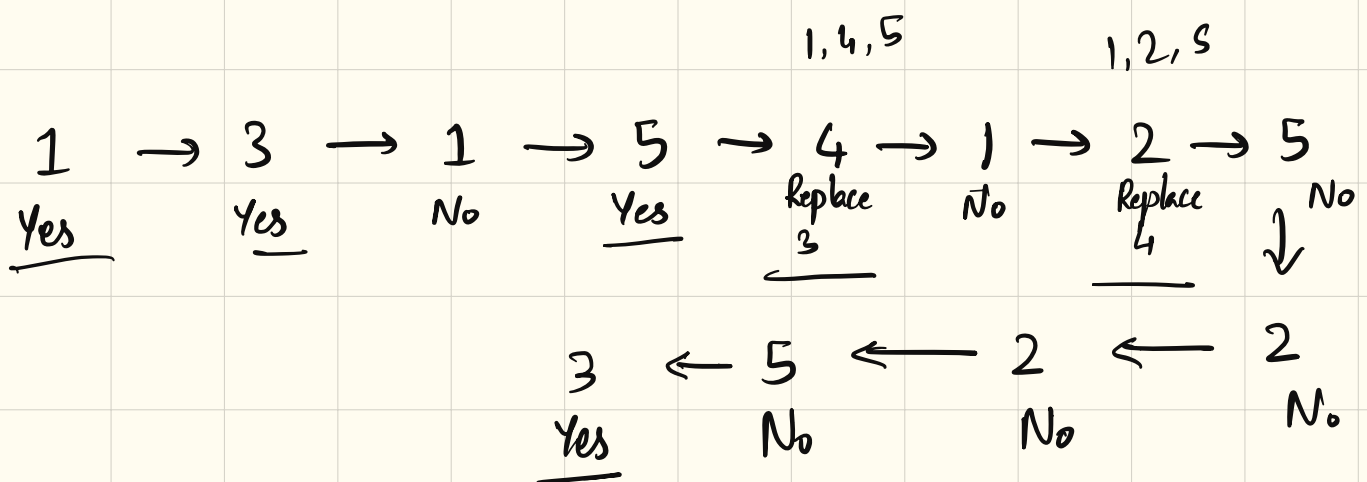                               page fault

# Belady's Optimal Algorithm (MIN)

→ Strategy : Replace the page that will
be referenced after the
longest time.

→ Example :

→ # of frames = 3

→ Reference sequence (in temporal
order )

$$1 \rightarrow 3 \rightarrow 1 \rightarrow 5 \rightarrow \overset{1,4,5}{4} \rightarrow 1 \rightarrow \overset{1,2,5}{2} \rightarrow 5$$

Yes  Yes  No  Yes  Replace  No  Replace  No
                       3              4     ↓

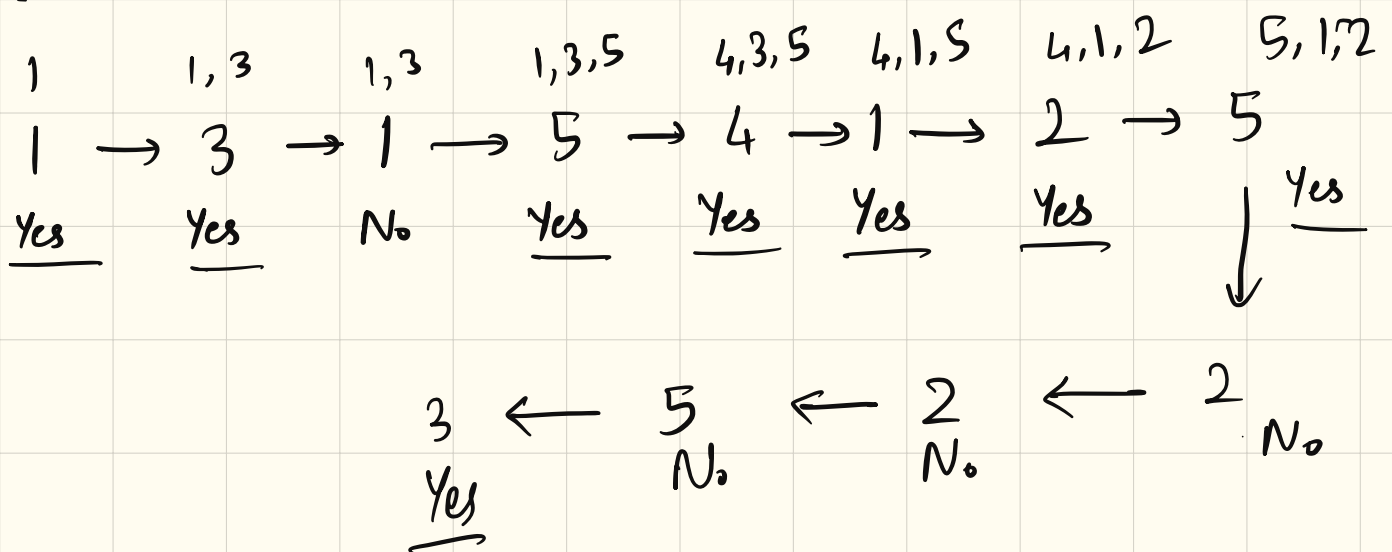$$3 \leftarrow 5 \leftarrow 2 \leftarrow 2$$

Yes  No  No  No.

# of page faults = 6

Belady's MIN algorithm is proven
to be optimal, but impractical
as it requires knowledge of
future accesses.


## FIFO

Strategy: Replace the page that is
in the memory for the longest
time


Ex:      # of frame = 3

| 1 | 1,3 | 1,3 | 1,3,5 | 4,3,5 | 4,1,5 | 4,1,2 | 5,1,2 |

$1 \rightarrow 3 \rightarrow 1 \rightarrow 5 \rightarrow 4 \rightarrow 1 \rightarrow 2 \rightarrow 5$

Yes    Yes    No    Yes    Yes    Yes    Yes    $\downarrow$ Yes

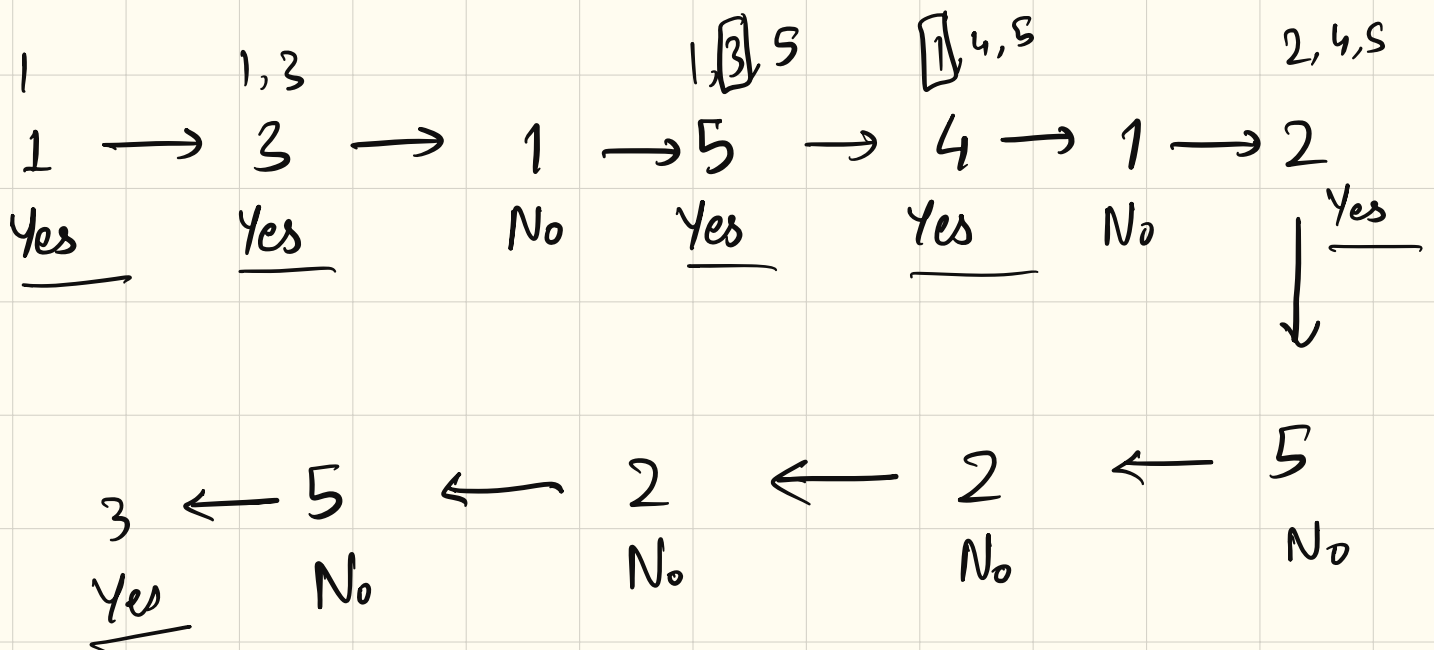$3 \leftarrow 5 \leftarrow 2 \leftarrow 2$
Yes      No       No        No


8 page faults — 3 cold start

# Least Recently Used (LRU)

Replace the page that is not referenced for the longest time.

Eg: # of frame = 3

1      1,3         1,3 5     1,4,5        2,4,5

$$1 \rightarrow 3 \rightarrow 1 \rightarrow 5 \rightarrow 4 \rightarrow 1 \rightarrow 2$$

Yes     Yes     No    Yes    Yes    No     Yes ↓

$$3 \leftarrow 5 \leftarrow 2 \leftarrow 2 \leftarrow 5$$

Yes    No     No     No     No

→ LRU is shown to be useful for workloads with access locality

→ approximate using CLOCK