



# last time (section 001)

deadlock: circular waiting  
not just locks

preventing deadlock:  
stop sharing resources  
stop waiting  
stop holding other resources when waiting  
**consistent order**

deadlock detection:  
dependency graph  
single resource: look for cycles  
multi-resource: simulate threads finishing

# last time (section 002)

deadlock: circular waiting

not just locks

preventing deadlock:

stop sharing resources

stop waiting

stop holding other resources when waiting

consistent order

# deadlock prevention techniques

**infinite resources**

or at least enough that never run out

*no mutual exclusion*

**no shared resources**

*no mutual exclusion*

**no waiting** (e.g. abort and retry)  
“busy signal”

*no hold and wait/  
preemption*

acquire resources in **consistent order**

*no circular wait*

request **all resources at once**

*no hold and wait*

# resource allocation graphs

nodes: resources or threads

edge thread→resource: thread waiting for resource

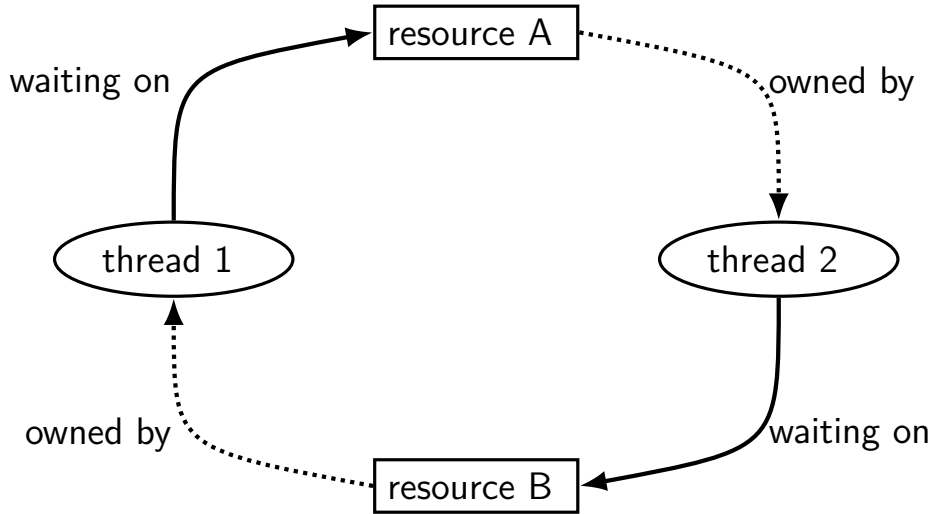
edge resource→thread: resource is “owned” by thread

holds lock on

will be deallocated by

...

# resource allocate graphs



# searching for cycles

cycle  $\rightarrow$  deadlock happened!

finding cycles: recall 2150 topological sort (maybe???)

## using deadlock detection for prevention

suppose you know the *maximum resources* a process could request  
make decision **when starting process** (“*admission control*”)



# using deadlock detection for prevention

suppose you know the *maximum resources* a process could request  
make decision **when starting process** (“*admission control*”)

ask “what if every process was waiting for maximum resources”  
including the one we’re starting

would it cause deadlock? then **don’t let it start**

called Baker’s algorithm

# beyond threads: event based programming

writing server that servers multiple clients?

e.g. multiple web browsers at a time

maybe don't really need multiple processors/cores

one network, not that fast

idea: one thread handles multiple connections

# beyond threads: event based programming

writing server that servers multiple clients?

e.g. multiple web browsers at a time

maybe don't really need multiple processors/cores

one network, not that fast

idea: one thread handles multiple connections

issue: read from/write to multiple streams at once?

# event loops

```
while (true) {
    event = WaitForNextEvent();
    switch (event.type) {
    case NEW_CONNECTION:
        handleNewConnection(event); break;
    case CAN_READ_DATA_WITHOUT_WAITING:
        connection = LookupConnection(event.fd);
        handleRead(connection);
        break;
    case CAN_WRITE_DATA_WITHOUT_WAITING:
        connection = LookupConnection(event.fd);
        handleWrite(connection);
        break;
        ...
    }
}
```

## some single-threaded processing code

```
void ProcessRequest(int fd) {
    while (true) {
        char command[1024] = {};
        size_t command_length = 0;
        do {
            ssize_t read_result =
                read(fd, command + command_length,
                    sizeof(command) - command_length);
            if (read_result <= 0) handle_error();
            command_length += read_result;
        } while (command[command_length - 1] != '\n');
        if (IsExitCommand(command)) { return; }
        char response[1024];
        computeResponse(response, commmand);
        size_t total_written = 0;
        while (total_written < sizeof(response)) {
            ...
        }
    }
}
```

## some single-threaded processing code

original code: loop to handle one request

reads/writes multiple times; each read/write can block

```
void Pro
while (true) {
    char command[1024] = {};
    size_t command_length = 0;
    do {
        ssize_t read_result =
            read(fd, command + command_length,
                sizeof(command) - command_length);
        if (read_result <= 0) handle_error();
        command_length += read_result;
    } while (command[command_length - 1] != '\n');
    if (IsExitCommand(command)) { return; }
    char response[1024];
    computeResponse(response, commmand);
    size_t total_written = 0;
    while (total_written < sizeof(response)) {
        ...
    }
}
```

## some single-threaded processing code

```
void ProcessRequest(int fd) {
    while (true) {
        char command[1024] = {};
        size_t command_length = 0;
        do {
            ssize_t read_result =
                read(fd, command + command_length,
                    sizeof(command) - command_length);
            if (read_result <= 0) handle_error();
            command_length += read_result;
        } while (command[command_length - 1] != '\n');
        if (IsExitCommand(command)) { return; }
        char response[1024];
        computeResponse(response, command);
        size_t total_written = 0;
        while (total_written < sizeof(response)) {
            ...
        }
    }
}
```

```
struct Connection {
    int fd;
    char command[1024];
    size_t command_length;
    char response[1024];
    size_t total_written;
    ...
};
```

## as event code

```
handleRead(Connection *c) {
    ssize_t read_result =
        read(fd, c->command + command_length,
            sizeof(command) - c->command_length);
    if (read_result <= 0) handle_error();
    c->command_length += read_result;

    if (c->command[c->command_length - 1] == '\n') {
        computeResponse(c->response, c->command);
        StopWaitingToRead(c->fd);
        StartWaitingToWrite(c->fd);
    }
}
```

new code: one read step per handleRead call  
Connection struct: info between write calls



## as event code

```
handleRead(Connection *c) {
    ssize_t read_result =
        read(fd, c->command + command_length,
             sizeof(command) - c->command_length);
    if (read_result <= 0) handle_error();
    c->command_length += read_result;

    if (c->command[c->command_length - 1] == '\n') {
        computeResponse(c->response, c->command);
        StopWaitingToRead(c->fd);
        StartWaitingToWrite(c->fd);
    }
}
```

# POSIX support for event loops

## `select` and `poll` functions

take list(s) of file descriptors to read and to write  
wait for them to be read/writeable without waiting  
(or for new connections associated with them, etc.)

many OS-specific extensions/improvements/alternatives:

examples: Linux `epoll`, Windows IO completion ports

better ways of managing list of file descriptors

enqueue read/write instead of learning when read/write okay

# message passing

instead of having variables, locks between threads...

send messages between threads/processes

what you need anyways between machines

big 'supercomputers' = really many machines together

arguably an easier model to program

can't have locking issues

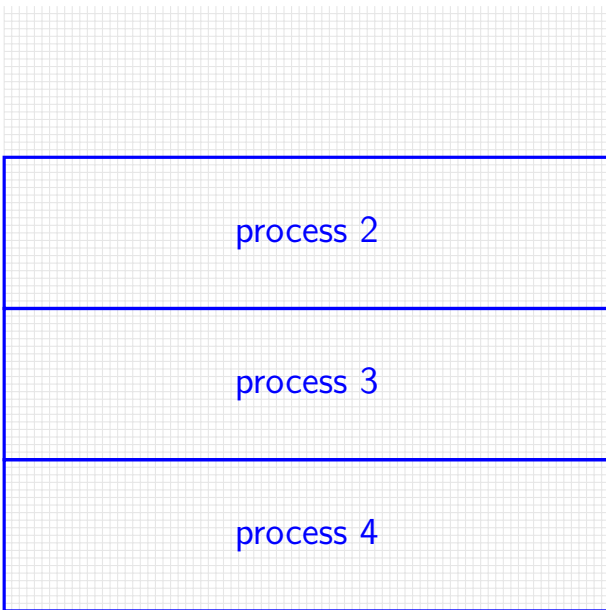
# message passing API

core functions: Send(toId, data)/Recv(fromId, data)

simplest(?) version: functions wait for other processes/threads

```
if (thread_id == 0) {
    for (int i = 1; i < MAX_THREAD; ++i) {
        Send(i, getWorkForThread(i));
    }
    for (int i = 1; i < MAX_THREAD; ++i) {
        WorkResult result;
        Recv(i, &result);
        handleResultForThread(i, result);
    }
} else {
    WorkInfo work;
    Recv(0, &work);
    Send(0, ComputeResultFor(work));
}
```

# message passing game of life

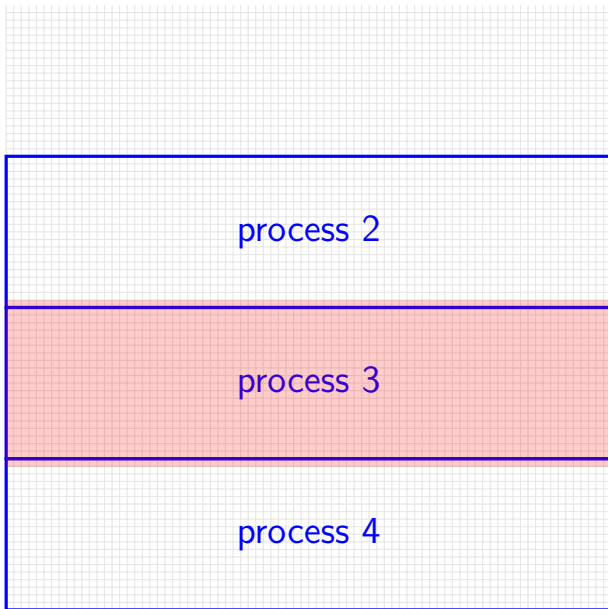


divide grid  
like you would for normal threads

each process **stores cells**  
in that part of grid

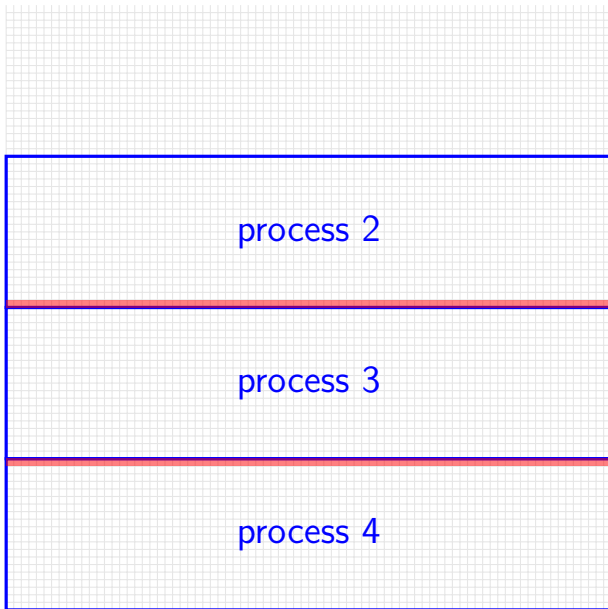
(no shared memory!)

# message passing game of life



process 3 only needs values of cells around its area (values of cells adjacent to the ones it computes)

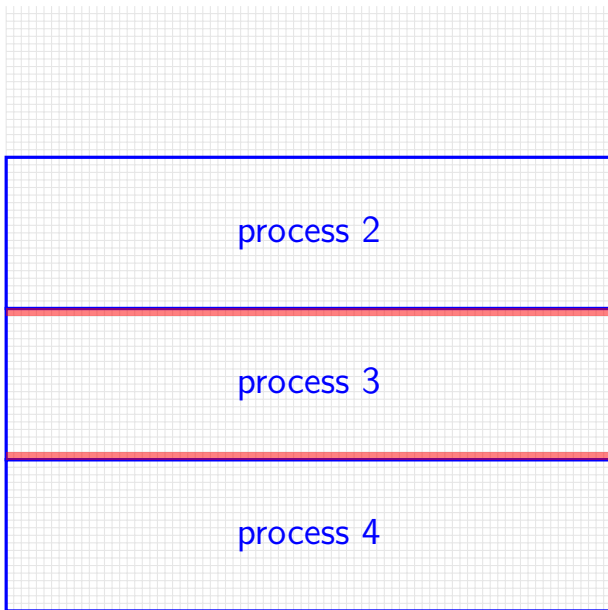
# message passing game of life



small slivers of  
other process's cells needed

solution: process 2, 4  
send messages with cells every iterat

# message passing game of life

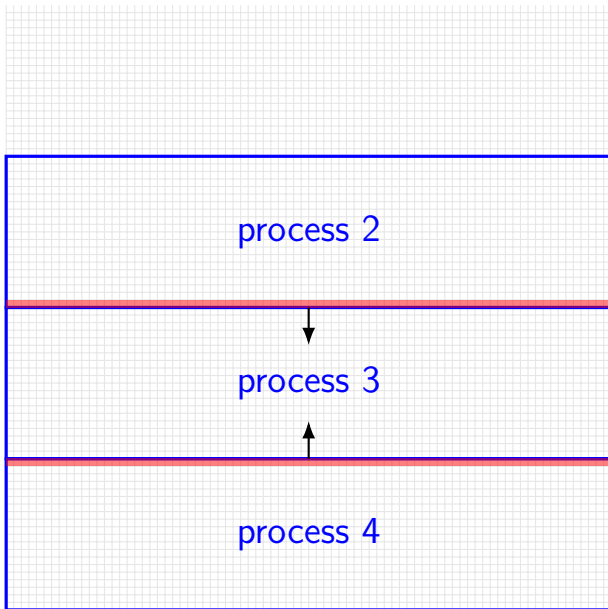


some of process 3's cells  
also needed by process 2/4

so process 3 also sends messages

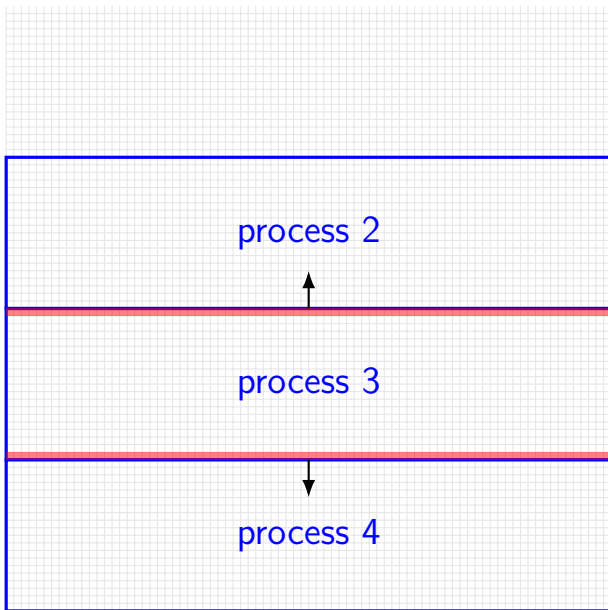


# message passing game of life



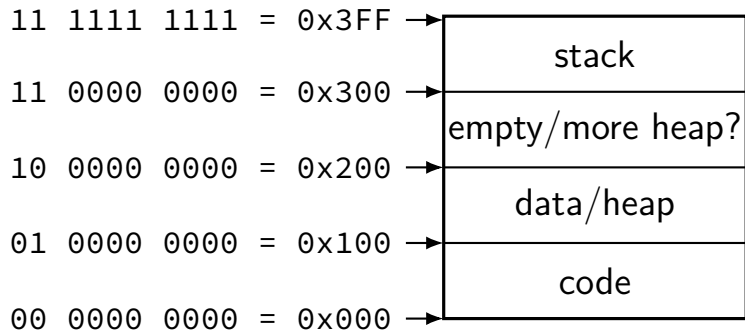
one possible pseudocode:  
all **even processes send messages**  
(while odd receives), then  
all odd processes send messages  
(while even receives)

# message passing game of life

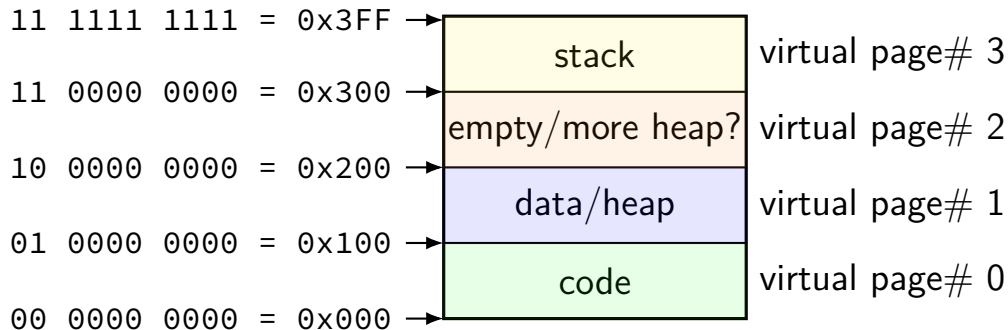


one possible pseudocode:  
all even processes send messages  
(while odd receives), then  
all **odd processes send messages**  
(while even receives)

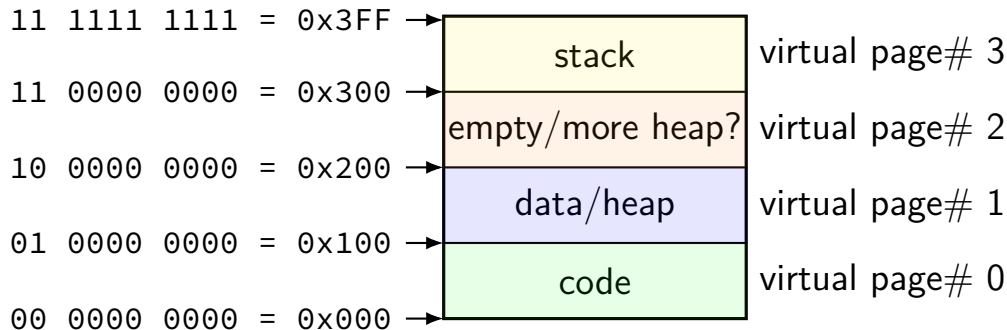
# toy program memory



# toy program memory

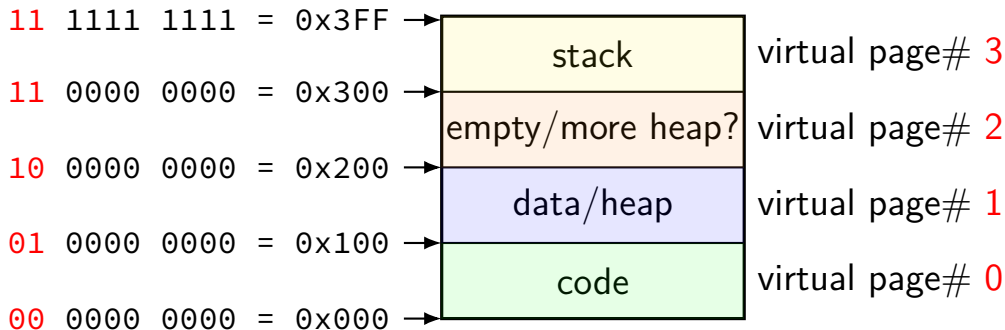


# toy program memory



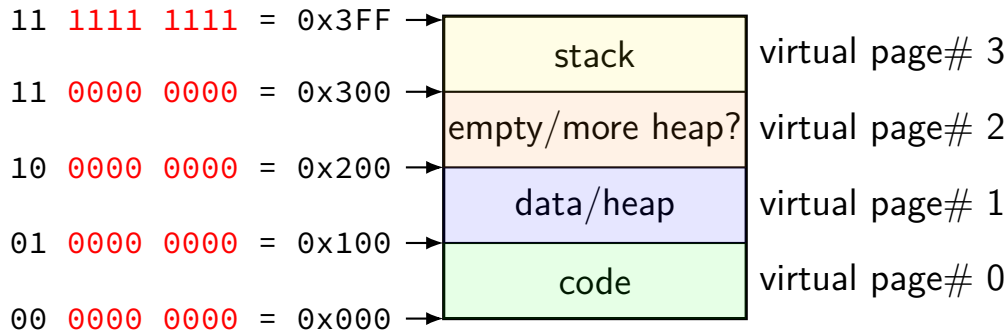
divide memory into **pages** ( $2^8$  bytes in this case)  
“virtual” = addresses the program sees

# toy program memory



page number is upper bits of address  
(because page size is power of two)

# toy program memory



rest of address is called **page offset**

# toy physical memory

program memory  
virtual addresses

11 0000 0000 to 11 1111 1111
10 0000 0000 to 10 1111 1111
01 0000 0000 to 01 1111 1111
00 0000 0000 to 00 1111 1111

real memory  
physical addresses

111 0000 0000 to 111 1111 1111
001 0000 0000 to 001 1111 1111
000 0000 0000 to 000 1111 1111



# toy physical memory

program memory  
virtual addresses

11 0000 0000 to
11 1111 1111
10 0000 0000 to
10 1111 1111
01 0000 0000 to
01 1111 1111
00 0000 0000 to
00 1111 1111

real memory  
physical addresses

111 0000 0000 to
111 1111 1111
001 0000 0000 to
001 1111 1111
000 0000 0000 to
000 1111 1111

physical page 7

physical page 1

physical page 0

# toy physical memory

program memory  
virtual addresses

11 0000 0000 to 11 1111 1111
10 0000 0000 to 10 1111 1111
01 0000 0000 to 01 1111 1111
00 0000 0000 to 00 1111 1111

real memory  
physical addresses

111 0000 0000 to 111 1111 1111
001 0000 0000 to 001 1111 1111
000 0000 0000 to 000 1111 1111

# toy physical memory

virtual page #	physical page #
00	010 (2)
01	111 (7)
10	<i>none</i>
11	000 (0)

program memory  
virtual addresses

11 0000 0000 to 11 1111 1111
10 0000 0000 to 10 1111 1111
01 0000 0000 to 01 1111 1111
00 0000 0000 to 00 1111 1111

real memory  
physical addresses

111 0000 0000 to 111 1111 1111
001 0000 0000 to 001 1111 1111
000 0000 0000 to 000 1111 1111

# toy physical memory

page table!

virtual page #	physical page #
00	010 (2)
01	111 (7)
10	<i>none</i>
11	000 (0)

program memory  
virtual addresses

11 0000 0000 to 11 1111 1111
10 0000 0000 to 10 1111 1111
01 0000 0000 to 01 1111 1111
00 0000 0000 to 00 1111 1111

real memory  
physical addresses

111 0000 0000 to 111 1111 1111
001 0000 0000 to 001 1111 1111
000 0000 0000 to 000 1111 1111

# toy page table lookup

virtual page #	valid?	physical page #
00	1	010 (2, code)
01	1	111 (7, data)
10	0	??? (ignored)
11	1	000 (0, stack)

# toy page table lookup

01 1101 0010 — address from CPU

virtual  
page # valid? physical page #

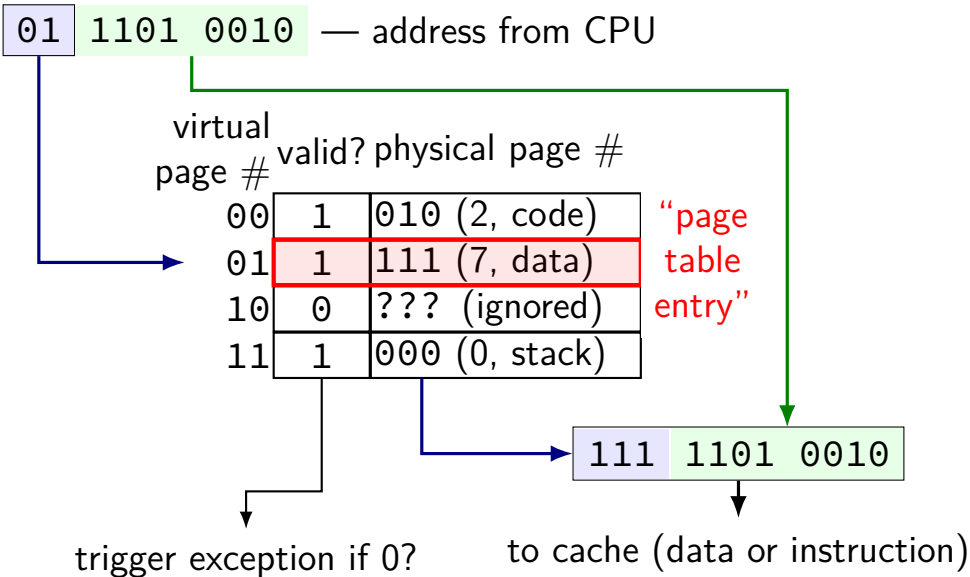
00	1	010 (2, code)
01	1	111 (7, data)
10	0	??? (ignored)
11	1	000 (0, stack)

111 1101 0010

trigger exception if 0?

to cache (data or instruction)

# toy page table lookup



# tov page table lookup

“virtual page number”

01 1101 0010 — address from CPU

virtual  
page # valid? physical page #

00	1	010 (2, code)
01	1	111 (7, data)
10	0	??? (ignored)
11	1	000 (0, stack)

trigger exception if 0?

to cache (data or instruction)



# toy page table lookup

01 1101 0010 — address from CPU

virtual  
page # valid? physical page #

00	1	010 (2, code)
01	1	111 (7, data)
10	0	??? (ignored)
11	1	000 (0, stack)

“physical page number”

111 1101 0010

trigger exception if 0?

to cache (data or instruction)

# toy page table lookup

“page offset”

01 1101 0010 — address from CPU

virtual  
page #    valid?    physical page #

00	1	010 (2, code)
01	1	111 (7, data)
10	0	??? (ignored)
11	1	000 (0, stack)

“page offset”

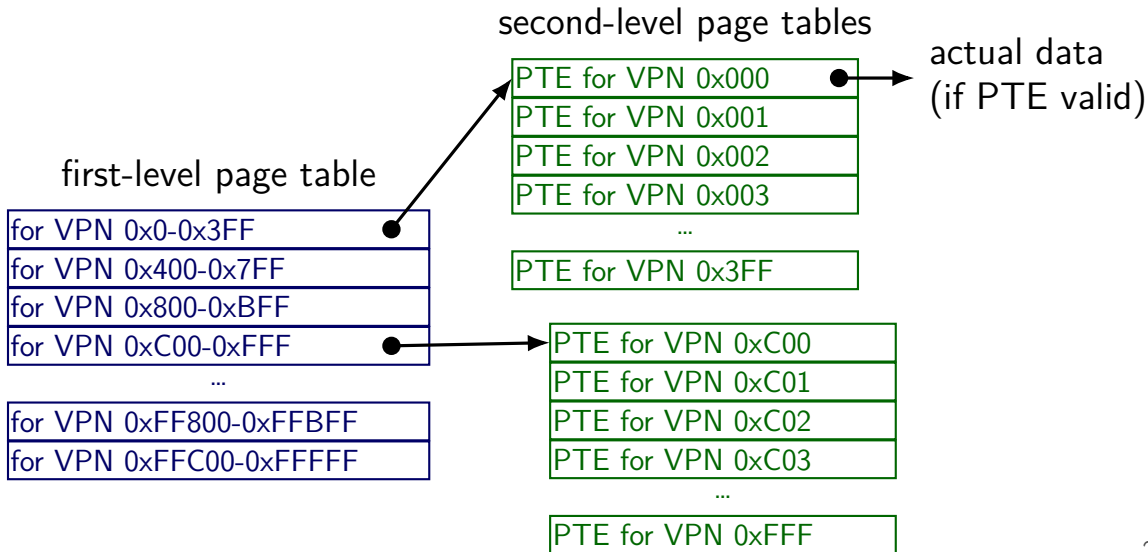
111 1101 0010

trigger exception if 0?

to cache (data or instruction)

# two-level page tables

two-level page table;  $2^{20}$  pages total;  $2^{10}$  entries per table



# two-level page tables

two-level page table;  $2^{20}$  pages total;  $2^{10}$  entries per table

x86-32: arrays of  $2^{10}$  32-bit  
*page table entries*

first-level page table

for VPN 0x0-0x3FF
for VPN 0x400-0x7FF
for VPN 0x800-0xBFF
for VPN 0xC00-0xFFF
...
for VPN 0xFF800-0xFFBFF
for VPN 0xFFC00-0xFFFFF

second-level page tables

PTE for VPN 0x000
PTE for VPN 0x001
PTE for VPN 0x002
PTE for VPN 0x003
...
PTE for VPN 0x3FF

PTE for VPN 0xC00
PTE for VPN 0xC01
PTE for VPN 0xC02
PTE for VPN 0xC03
...
PTE for VPN 0xFFF

actual data  
(if PTE valid)

# two-level page tables

two-level page table;  $2^{20}$  pages total;  $2^{10}$  entries per table

second-level page tables

actual data  
(if PTE valid)

first-level page table

for VPN 0x0-0x3FF	●
for VPN 0x400-0x7FF	✗
for VPN 0x800-0xBFF	✗
for VPN 0xC00-0xFFF	●
...	
for VPN 0xFF800-0xFFBFF	
for VPN 0xFFC00-0xFFFFF	

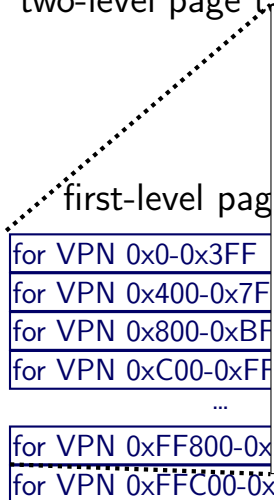
PTE for VPN 0x000
PTE for VPN 0x001
PTE for VPN 0x002
PTE for VPN 0x003
...

invalid entries represent big holes

PTE for VPN 0xC00
PTE for VPN 0xC01
PTE for VPN 0xC02
PTE for VPN 0xC03
...
PTE for VPN 0xFFF

# two-level page tables

two-level page table:  $2^{20}$  pages total,  $2^{10}$  entries per table



## first-level page table

VPN range	valid	user?	write?	physical page # (of next page table)
0x0-0x3FF	1	1	1	0x22343
0x400-0x7FF	0	0	1	0x00000
0x800-0xBFF	0	0	0	0x00000
0xC00-0xFFF	1	1	0	0x33454
0x1000-0x13FF	1	1	0	0xFF043
...	...	...	...	...
0xFFC00-0xFFFF	1	1	0	0xFF045

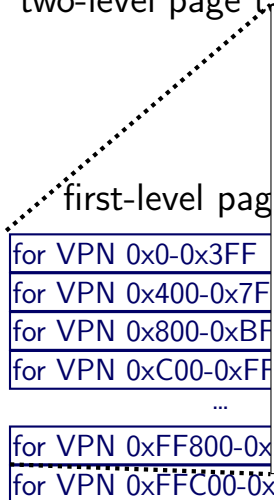
PTE for VPN 0xC03

...

PTE for VPN 0xFFF

# two-level page tables

two-level page table:  $2^{20}$  pages total,  $2^{10}$  entries per table



## first-level page table

VPN range	valid	user?	write?	physical page # (of next page table)
0x0-0x3FF	1	1	1	0x22343
0x400-0x7FF	0	0	1	0x00000
0x800-0xBFF	0	0	0	0x00000
0xC00-0xFFF	1	1	0	0x33454
0x1000-0x13FF	1	1	0	0xFF043
...	...	...	...	...
0xFFC00-0xFFFF	1	1	0	0xFF045

PTE for VPN 0xC03

...

PTE for VPN 0xFFF

# two-level page tables

two-level page table:  $2^{20}$  pages total:  $2^{10}$  entries per table

first-level page table  
 for VPN 0x0-0x3FF  
 for VPN 0x400-0x7FF  
 for VPN 0x800-0xBF  
 for VPN 0xC00-0xFF  
 ...  
 for VPN 0xFF800-0x  
 for VPN 0xFFC00-0xFFFF

## first-level page table

VPN range	valid	user?	write?	physical page # (of next page table)
0x0-0x3FF	1	1	1	0x22343
0x400-0x7FF	0	0	1	0x00000
0x800-0xBF	0	0	0	0x00000
0xC00-0xFFF	0	0	0	0x33454
0x1000-0x13FF	0	0	0	0xFF043
...	...	...	...	...
0xFFC00-0xFFFF	1	1	0	0xFF045

pointers to page tables  
 (arrays of PTEs)  
 but using page number  
 (not byte number)

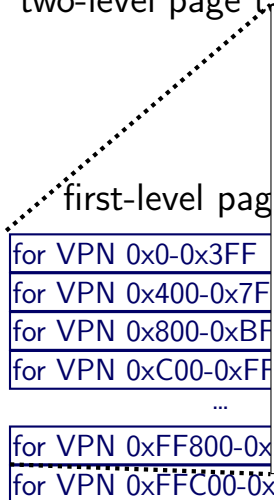
PTE for VPN 0xC03

PTE for VPN 0xFF



# two-level page tables

two-level page table:  $2^{20}$  pages total,  $2^{10}$  entries per table



## first-level page table

VPN range	valid	user?	write?	physical page # (of next page table)
0x0-0x3FF	1	1	1	0x22343
0x400-0x7FF	0	0	1	0x00000
0x800-0xBFF	0	0		
0xC00-0xFFF	1	1		
0x1000-0x13FF	1	1		
...	...	...	...	...
0xFFC00-0xFFFF	1	1	0	0xFF045

valid bits indicate "holes"  
note: physical page 0 is valid  
so can't use NULL ptrs

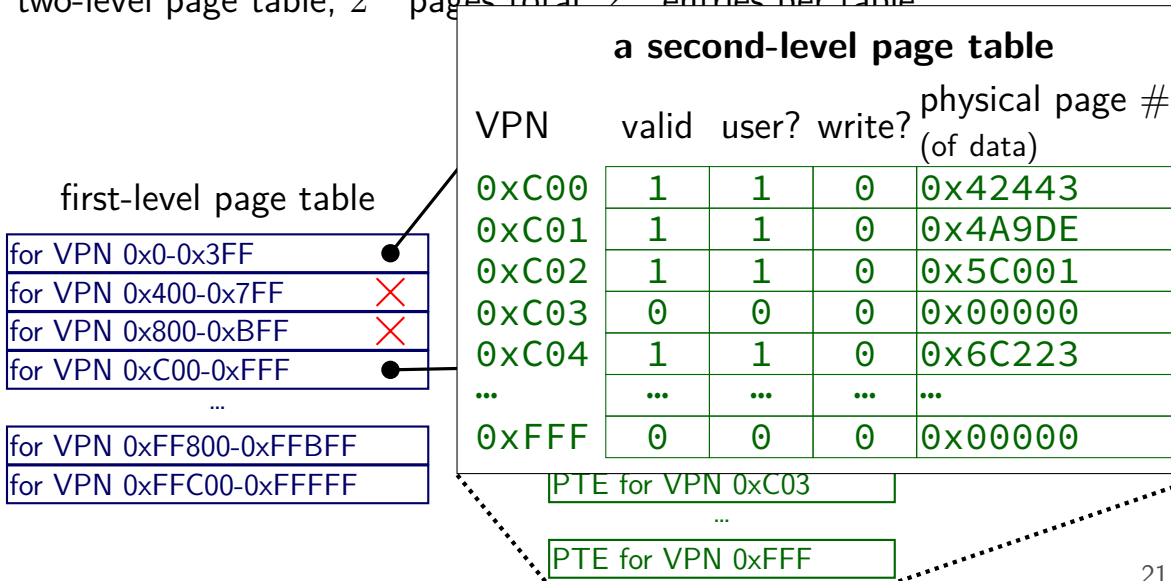
PTE for VPN 0xC03

...

PTE for VPN 0xFFF

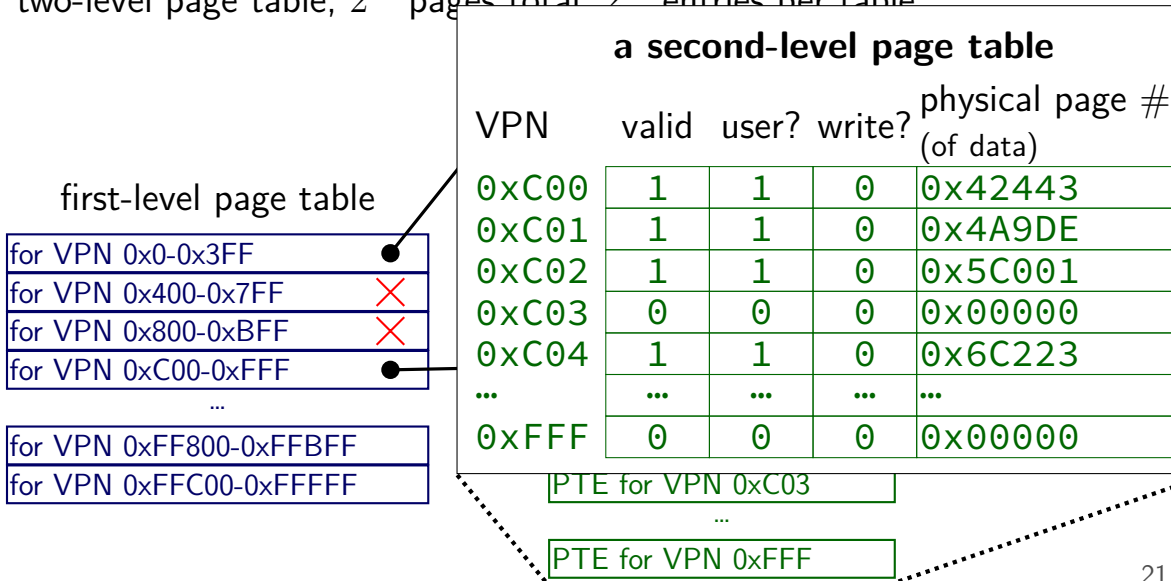
# two-level page tables

two-level page table;  $2^{20}$  pages total.  $2^{10}$  entries per table



# two-level page tables

two-level page table;  $2^{20}$  pages total;  $2^{10}$  entries per table



# x86-32 pagetables: overall structure

xv6 header: mmu.h

```
// A virtual address 'va' has a three-part structure as follows:
```

```
//  
// +-----10-----+-----10-----+-----12-----+  
// | Page Directory | Page Table | Offset within Page |  
// |      Index      |      Index      |                   |  
// +-----+-----+-----+  
// \--- PDX(va) ---/ \--- PTX(va) ---/
```

```
// page directory index
```

```
#define PDX(va) ((uint)(va) >> PDXSHIFT) & 0x3FF)
```

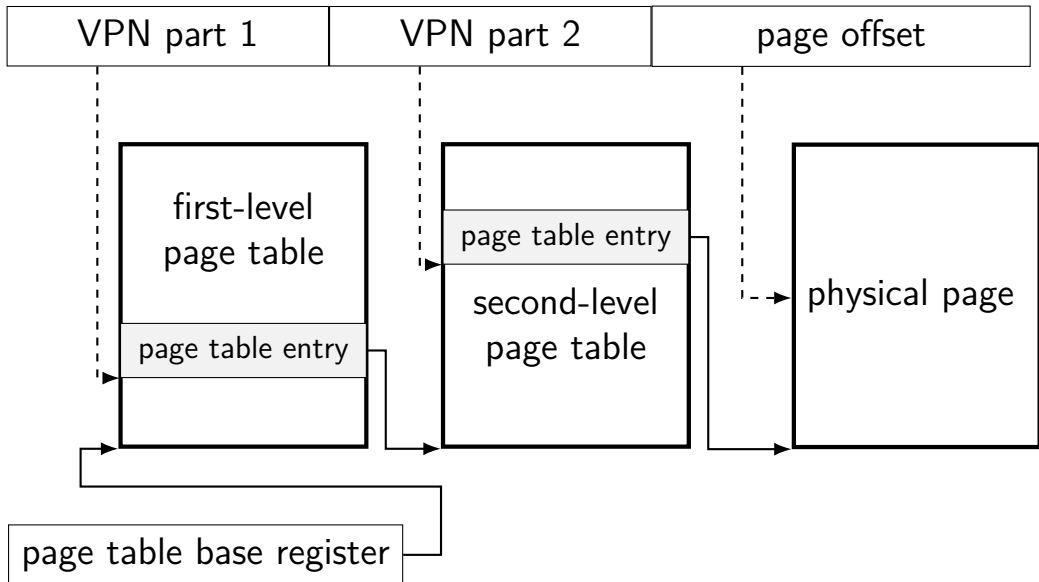
```
// page table index
```

```
#define PTX(va) ((uint)(va) >> PTXSHIFT) & 0x3FF)
```

```
// construct virtual address from indexes and offset
```

```
#define PGADDR(d, t, o) ((uint)((d) << PDXSHIFT | (t) << PTXSHIFT |
```

## another view



## 32-bit x86 paging

4096 ( $= 2^{12}$ ) byte pages

4-byte page table entries — stored in memory

two-level table:

- first 10 bits lookup in first level (“page directory”)

- second 10 bits lookup in second level

remaining 12 bits: which byte of 4096 in page?

## exercise

4096 ( $= 2^{12}$ ) byte pages

4-byte page table entries — stored in memory

two-level table:

- first 10 bits lookup in first level (“page directory”)

- second 10 bits lookup in second level

exercise: how big is...

- a process's x86-32 page tables with 1 valid 4K page?

- a process's x86-32 page table with all 4K pages populated?

# x86-32 page table entries

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									
Address of page directory <sup>1</sup>											Ignored						P C D	PW T	Ignored			CR3																		
Bits 31:22 of address of 4MB page frame					Reserved (must be 0)				Bits 39:32 of address <sup>2</sup>				P A T	Ignored	G	<u>1</u>	D	A	P C D	PW T	U / S	R / W	<u>1</u>	PDE: 4MB page																
Address of page table											Ignored						<u>0</u>	I g n	A	P C D	PW T	U / S	R / W	<u>1</u>	PDE: page table															
Ignored																	<u>0</u>																						PDE: not present	
Address of 4KB page frame											Ignored						G	P A T	D	A	P C D	PW T	U / S	R / W	<u>1</u>	PTE: 4KB page														
Ignored																	<u>0</u>																							PTE: not present

Figure 4-4. Formats of CR3 and Paging-Structure Entries with 32-Bit Paging



# x86-32 page table entries

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	
Address of page directory <sup>1</sup>											Ignored						P C D	PW T	Ignored			CR3										
Bits 31:22 of address of 4MB page frame											page table base register (CR3)											P C D	PW T	U / S	R / W	1	PDE: 4MB page					
Address of page table											Ignored						0	I g n	A	P C D	PW T	U / S	R / W	1	PDE: page table							
Ignored											Ignored						0												PDE: not present			
Address of 4KB page frame											Ignored						G	P A T	D	A	P C D	PW T	U / S	R / W	1	PTE: 4KB page						
Ignored											Ignored						0												PTE: not present			

Figure 4-4. Formats of CR3 and Paging-Structure Entries with 32-Bit Paging

# x86-32 page table entries

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			
Address of page																												P C D	PW T	Ignored			CR3	
Bits 31:22 of address of 4MB page frame										Reserved (must be 0)					Bits 39:32 of address <sup>2</sup>				P A T	Ignored	G	<u>1</u>	D	A	P C D	PW T	U / S	R / W	<u>1</u>	PDE: 4MB page				
Address of page table										Ignored					<u>0</u>	Ignored	G	<u>0</u>	I g n	A	P C D	PW T	U / S	R / W	<u>1</u>	PDE: page table								
Ignored																												<u>0</u>	PDE: not present					
Address of 4KB page frame										Ignored					G	P A T	D	A	P C D	PW T	U / S	R / W	<u>1</u>	PTE: 4KB page										
Ignored																												<u>0</u>	PTE: not present					

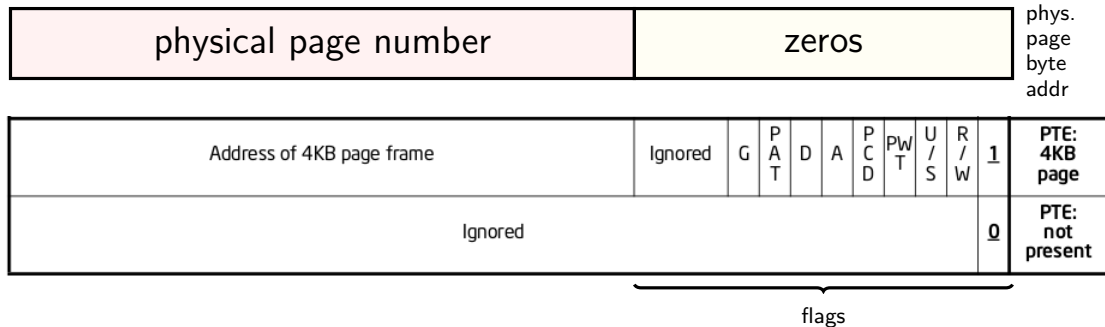
Figure 4-4. Formats of CR3 and Paging-Structure Entries with 32-Bit Paging

# x86-32 page table entries

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	
Address of page directory <sup>1</sup>											Ignored						P C D	PW T	Ignored			CR3										
Bits 31:22 of address of 4MB page frame					Reserved (must be 0)				Bits 39:32 of address <sup>2</sup>				P A T	Ignored	G	<u>1</u>	D	A	P C D	PW T	U / S	R / W	<u>1</u>	PDE: 4MB page								
Address of page table											Ignored						<u>0</u>	I g n	A	P C D	PW T	U / S	R / W	<u>1</u>	PDE: page table							
second-level page table entries																				<u>0</u>	PDE: not present											
Address of 4KB page frame											Ignored						G	P A T	D	A	P C D	PW T	U / S	R / W	<u>1</u>	PTE: 4KB page						
Ignored																				<u>0</u>	PTE: not present											

Figure 4-4. Formats of CR3 and Paging-Structure Entries with 32-Bit Paging

# x86-32 page table entry v addresses



trick: page table entry with lower bits zeroed =  
physical *byte* address of corresponding page  
page # is address of page ( $2^{12}$  byte units)

makes constructing page table entries simpler:  
physicalAddress | flagsBits

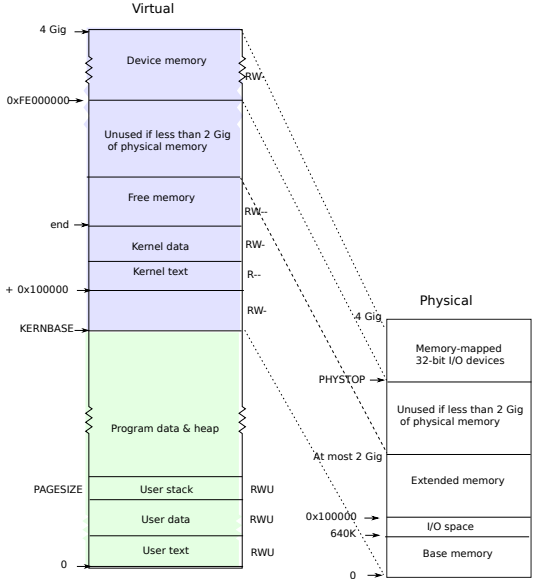
# x86-32 pagetables: page table entries

xv6 header: mmu.h

```
// Page table/directory entry flags.
#define PTE_P           0x001    // Present
#define PTE_W           0x002    // Writeable
#define PTE_U           0x004    // User
#define PTE_PWT         0x008    // Write-Through
#define PTE_PCD         0x010    // Cache-Disable
#define PTE_A           0x020    // Accessed
#define PTE_D           0x040    // Dirty
#define PTE_PS          0x080    // Page Size
#define PTE_MBZ         0x180    // Bits must be zero

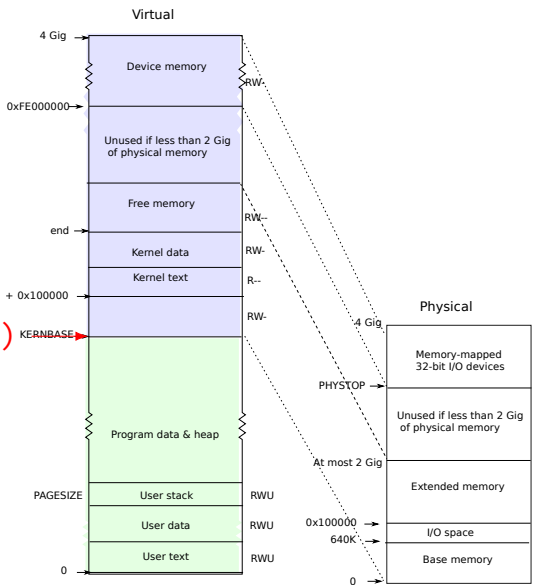
// Address in page table or page directory entry
#define PTE_ADDR(pte)   ((uint)(pte) & ~0xFFF)
#define PTE_FLAGS(pte) ((uint)(pte) &  0xFFF)
```

# xv6 memory layout

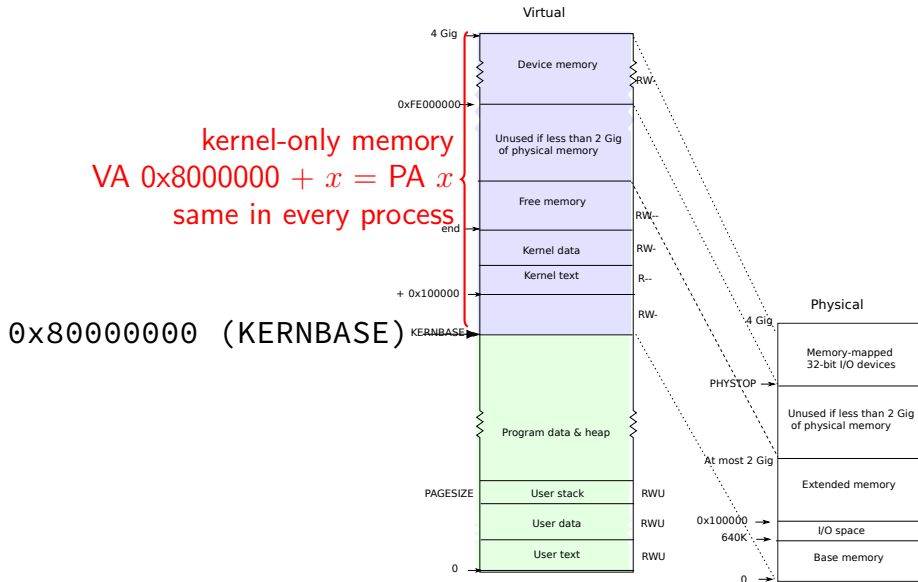


# xv6 memory layout

0x80000000 (KERNBASE)



# xv6 memory layout







# xv6 kernel memory

virtual memory  $>$  KERNBASE ( $0x8000\ 0000$ ) is for kernel

always mapped as kernel-mode only

protection fault for user-mode programs to access

physical memory address 0 is mapped to  $\text{KERNBASE}+0$

physical memory address  $N$  is mapped to  $\text{KERNBASE}+N$

not done by hardware — just page table entries OS sets up on boot  
very convenient for manipulating page tables with physical addresses

kernel code loaded into contiguous physical addresses

# P2V/V2P

V2P( $x$ ) (virtual to physical)

convert *kernel* address  $x$  to physical address

subtract KERNBASE (0x8000 0000)

have user address? need full page table lookup instead

P2V( $x$ ) (physical to virtual)

convert *physical* address  $x$  to kernel address

add KERNBASE (0x8000 0000)

xv6 convention: virtual addresses represented using pointers

xv6 convention: physical addresses represented using integers

# P2V/V2P

V2P( $x$ ) (virtual to physical)

convert *kernel* address  $x$  to physical address

subtract KERNBASE (0x8000 0000)

have user address? need full page table lookup instead

P2V( $x$ ) (physical to virtual)

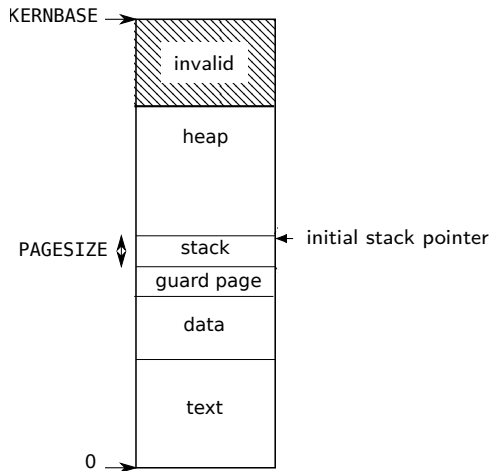
convert *physical* address  $x$  to kernel address

add KERNBASE (0x8000 0000)

xv6 convention: virtual addresses represented using pointers

xv6 convention: physical addresses represented using integers

# xv6 program memory



## xv6 page table-related functions

`walkpgdir` — get pointer to second-level page table entry  
...to set make it valid/invalid/point somewhere/etc.

`mappages` — set range of page table entries  
implementation: loop using `walkpgdir`

`allocvm` — create first-level page table, set kernel (high) part

`allocvm` — allocate new user memory

`kalloc` — allocate physical page, return kernel address

`deallocvm` — deallocate user memory

# xv6 page table-related functions

`walkpgdir` — get pointer to second-level page table entry  
...to set make it valid/invalid/point somewhere/etc.

`mappages` — set range of page table entries  
implementation: loop using `walkpgdir`

`allocvm` — create first-level page table, set kernel (high) part

`allocvm` — allocate new user memory

`kalloc` — allocate physical page, return kernel address

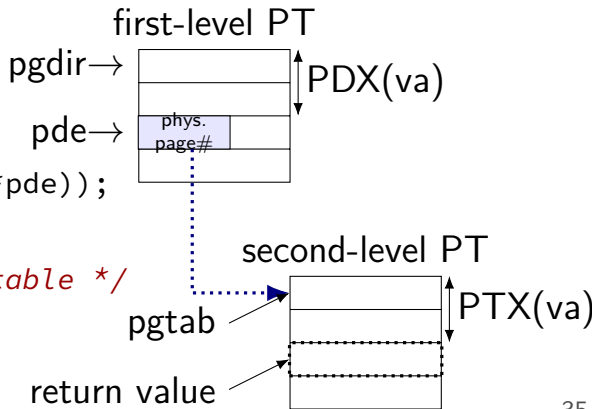
`deallocvm` — deallocate user memory

# xv6: finding page table entries

```
// Return the address of the PTE in page table pgdir  
// that corresponds to virtual address va. If alloc!=0,  
// create any required page table pages.
```

```
static pte_t *  
walkpgdir(pde_t *pgdir, const void *va, int alloc)
```

```
{  
    pde_t *pde;  
    pte_t *pgtab;  
  
    pde = &pgdir[PDX(va)];  
    if(*pde & PTE_P){  
        pgtab = (pte_t*)P2V(PTE_ADDR(*pde));  
    } else {  
        ... /* create new  
                second-level page table */  
    }  
    return &pgtab[PTX(va)];  
}
```



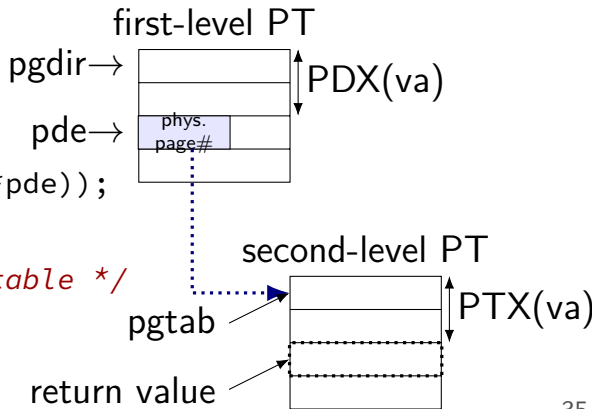


# xv6: finding page table entries

pde\_t — page directory entry  
pte\_t — page table entry  
both aliases for uint (32-bit unsigned int)

```
// Return the address of
// that corresponds to va
// create any required page table pages.
static pte_t *
walkpgdir(pde_t *pgdir, const void *va, int alloc)
{
    pde_t *pde;
    pte_t *pgtab;

    pde = &pgdir[PDX(va)];
    if(*pde & PTE_P){
        pgtab = (pte_t*)P2V(PTE_ADDR(*pde));
    } else {
        ... /* create new
              second-level page table */
    }
    return &pgtab[PTX(va)];
}
```

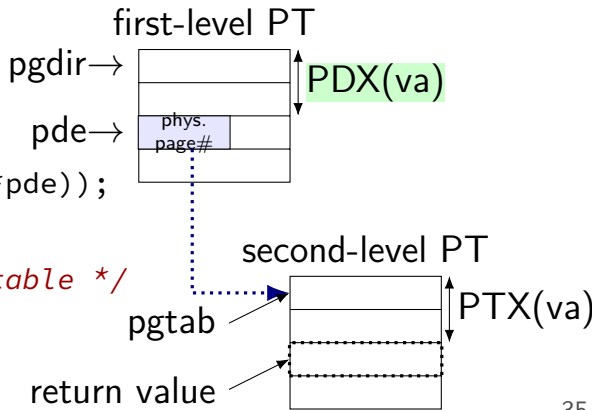


# xv6: finding page table entries

PDX(va) — extract top 10 bits of va  
used to index into first-level page table

```
// Return the address of the page table entry  
// that corresponds to virtual address va. If alloc!=0,  
// create any required page table pages.
```

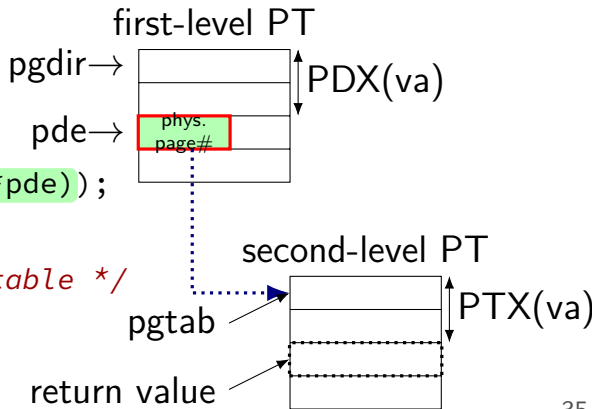
```
static pte_t *  
walkpgdir(pde_t *pgdir, const void *va, int alloc)  
{  
    pde_t *pde;  
    pte_t *pgtab;  
  
    pde = &pgdir[PDX(va)];  
    if(*pde & PTE_P){  
        pgtab = (pte_t*)P2V(PTE_ADDR(*pde));  
    } else {  
        ... /* create new  
            second-level page table */  
    }  
    return &pgtab[PTX(va)];  
}
```



## xv6: finding page table entries

`PTE_ADDR(*pde)` — return second-level page table address  
from first-level page table entry `*pde`  
returns *physical address*

```
// Return physical address  
// that returns physical address  
// create any required page table pages.  
static pte_t *  
walkpgdir(pde_t *pgdir, const void *va, int alloc)  
{  
    pde_t *pde;  
    pte_t *pgtab;  
  
    pde = &pgdir[PDX(va)];  
    if(*pde & PTE_P){  
        pgtab = (pte_t*)P2V(PTE_ADDR(*pde));  
    } else {  
        ... /* create new  
            second-level page table */  
    }  
    return &pgtab[PTX(va)];  
}
```



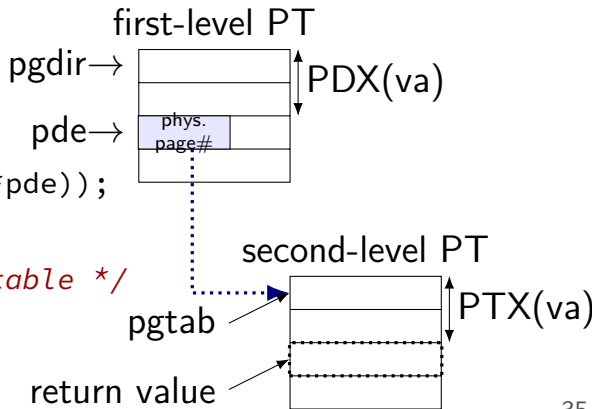
# xv6: finding page table entries

P2V — physical address to virtual address  
by convention, kernel maps physical memory at address  
KERNBASE (will show setup later)

```
// Return the physical address of the page table entry  
// that corresponds to the virtual address va  
// create an entry if it does not exist  
static pte_t *walkpgdir(pd_t *pd, va_t va)  
{
```

result is address that can access second-level page table

```
    pde_t *pde;  
    pte_t *pgtab;  
  
    pde = &pgdir[PDX(va)];  
    if(*pde & PTE_P){  
        pgtab = (pte_t*)P2V(PTE_ADDR(*pde));  
    } else {  
        ... /* create new  
            second-level page table */  
    }  
    return &pgtab[PTX(va)];  
}
```

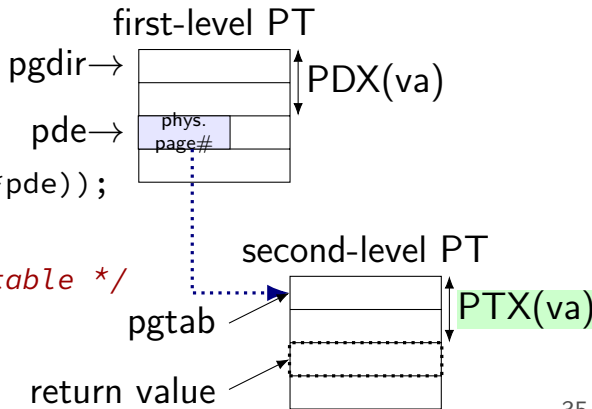


# xv6: finding page table entries

```
// Return the address of the page table entry  
// that corresponds to va  
// create any required page table pages.
```

```
static pte_t *  
walkpgdir(pde_t *pgdir, const void *va, int alloc)  
{  
    pde_t *pde;  
    pte_t *pgtab;  
  
    pde = &pgdir[PDX(va)];  
    if(*pde & PTE_P){  
        pgtab = (pte_t*)P2V(PTE_ADDR(*pde));  
    } else {  
        ... /* create new  
            second-level page table */  
    }  
    return &pgtab[PTX(va)];  
}
```

lookup in second-level page table  
PTX retrieves second-level page table index  
(= bits 10-20 of va)



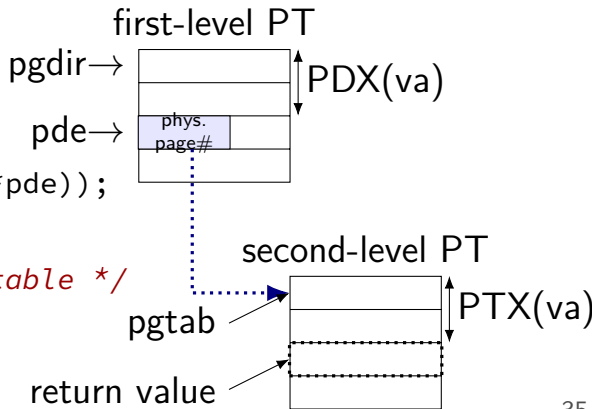
# xv6: finding page table entries

```
// Return the address of the PTE in  
// that corresponds to virtual addr  
// create any required page table p
```

```
static pte_t *  
walkpgdir(pde_t *pgdir, const void *va, int alloc)
```

```
{  
    pde_t *pde;  
    pte_t *pgtab;  
  
    pde = &pgdir[PDX(va)];  
    if(*pde & PTE_P){  
        pgtab = (pte_t*)P2V(PTE_ADDR(*pde));  
    } else {  
        ... /* create new  
            second-level page table */  
    }  
    return &pgtab[PTX(va)];  
}
```

if no second-level page table  
(present bit in first-level = 0)  
create one (if alloc=1)  
or return null (if alloc=0)



## xv6: creating second-level page tables

```
...
if(*pde & PTE_P){
    pgtab = (pte_t*)P2V(PTE_ADDR(*pde));
} else {
    if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
        return 0;
    // Make sure all those PTE_P bits are zero.
    memset(pgtab, 0, PGSIZE);
    // The permissions here are overly generous, but they can
    // be further restricted by the permissions in the page table
    // entries, if necessary.
    *pde = V2P(pgtab) | PTE_P | PTE_W | PTE_U;
}
```

## xv6: creating second-level page tables

return NULL if not trying to make new page table  
otherwise use kalloc to allocate it

```
...
if(*pde & PTE_P){
    pgtab = (pte_t*)P2V(PTE_ADDR(*pde));
} else {
    if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
        return 0;
    // Make sure all those PTE_P bits are zero.
    memset(pgtab, 0, PGSIZE);
    // The permissions here are overly generous, but they can
    // be further restricted by the permissions in the page table
    // entries, if necessary.
    *pde = V2P(pgtab) | PTE_P | PTE_W | PTE_U;
}
```



## xv6: creating second-level page tables

clear the page table  
PTE = 0  $\rightarrow$  present = 0

```
...
if(*pde & PTE_P){
    pgtab = (pte_t*)P2V(PTE_ADDR(*pde));
} else {
    if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
        return 0;
    // Make sure all those PTE_P bits are zero.
    memset(pgtab, 0, PGSIZE);
    // The permissions here are overly generous, but they can
    // be further restricted by the permissions in the page table
    // entries, if necessary.
    *pde = V2P(pgtab) | PTE_P | PTE_W | PTE_U;
}
```

## xv6: creating second-level page tables

```
...
if(*pde & PTE_P for "present" (valid)
    pgtab = (pt W for "writable" (pages access via may be writable)
} else {
    if(!alloc | U for "user-mode" (in addition to kernel)
        return 0;
    // Make sur second-level permission bits can restrict further
    memset(pgtab // The permissions here are overly generous, but they can
    // be further restricted by the permissions in the page table
    // entries, if necessary.
    *pde = V2P(pgtab) | PTE_P | PTE_W | PTE_U;
}
```

## xv6: creating second-level page tables

```
...
if(*pde & PTE_P){
    pgtab = (pte_t*)P2V(PTE_ADDR(*pde));
} else {
    if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
        return 0;
    // Make sure all those PTE_P bits are zero.
    memset(pgtab, 0, PGSIZE);
    // The permissions here are overly generous, but they can
    // be further restricted by the permissions in the page table
    // entries, if necessary.
    *pde = V2P(pgtab) | PTE_P | PTE_W | PTE_U;
}
```

# xv6 page table-related functions

`walkpgdir` — get pointer to second-level page table entry  
...to set make it valid/invalid/point somewhere/etc.

`mappages` — set range of page table entries  
implementation: loop using `walkpgdir`

`allocvm` — create first-level page table, set kernel (high) part

`allocvm` — allocate new user memory

`kalloc` — allocate physical page, return kernel address

`deallocvm` — deallocate user memory

## xv6: setting last-level page entries

```
static int
mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
{
    char *a, *last; pte_t *pte;

    a = (char*)PGROUNDDOWN((uint)va);
    last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
    for(;;){
        if((pte = walkpgdir(pgdir, a, 1)) == 0)
            return -1;
        if(*pte & PTE_P)
            panic("remap");
        *pte = pa | perm | PTE_P;
        if(a == last)
            break;
        a += PGSIZE;
        pa += PGSIZE;
    }
    return 0;
}
```

## xv6: setting last-level page entries

for each virtual page in range (va to va + size)  
get its page table entry  
(or fail if out of memory)

```
static int
mappages(pde_t *pgdir)
{
    char *a, *last; pte_t *pte;

    a = (char*)PGROUNDDOWN((uint)va);
    last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
    for(;;){
        if((pte = walkpgdir(pgdir, a, 1)) == 0)
            return -1;
        if(*pte & PTE_P)
            panic("remap");
        *pte = pa | perm | PTE_P;
        if(a == last)
            break;
        a += PGSIZE;
        pa += PGSIZE;
    }
    return 0;
}
```

## xv6: setting last-level page entries

make sure it's not already set

```
static int
mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
{
    char *a, *last; pte_t *pte;

    a = (char*)PGROUNDDOWN((uint)va);
    last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
    for(;;){
        if((pte = walkpgdir(pgdir, a, 1)) == 0)
            return -1;
        if(*pte & PTE_P)
            panic("remap");
        *pte = pa | perm | PTE_P;
        if(a == last)
            break;
        a += PGSIZE;
        pa += PGSIZE;
    }
    return 0;
}
```

## xv6: setting last-level page entries

```
static int  
mappages(pde  
{  
    char *a, *
```

create page table entry

pointing to physical page at pa

with specified permission bits (write and/or user-mode)

and P for present

```
    a = (char*)PGROUNDDOWN((uint)va);  
    last = (char*)PGROUNDDOWN(((uint)va) + size - 1);  
    for(;;){  
        if((pte = walkpgdir(pgdir, a, 1)) == 0)  
            return -1;  
        if(*pte & PTE_P)  
            panic("remap");  
        *pte = pa | perm | PTE_P;  
        if(a == last)  
            break;  
        a += PGSIZE;  
        pa += PGSIZE;  
    }  
    return 0;
```



## xv6: setting last-level page entries

advance to next physical page (pa)  
and next virtual page (va)

```
static int
mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
{
    char *a, *last; pte_t *pte;

    a = (char*)PGROUNDDOWN((uint)va);
    last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
    for(;;){
        if((pte = walkpgdir(pgdir, a, 1)) == 0)
            return -1;
        if(*pte & PTE_P)
            panic("remap");
        *pte = pa | perm | PTE_P;
        if(a == last)
            break;
        a += PGSIZE;
        pa += PGSIZE;
    }
    return 0;
}
```

# xv6 page table-related functions

`walkpgdir` — get pointer to second-level page table entry  
...to set make it valid/invalid/point somewhere/etc.

`mappages` — set range of page table entries  
implementation: loop using `walkpgdir`

`allocvm` — create first-level page table, set kernel (high) part

`allocvm` — allocate new user memory

`kaalloc` — allocate physical page, return kernel address

`deallocvm` — deallocate user memory

# xv6: setting process page tables (exec())

exec step 1: create new page table with kernel mappings

setupkvm()

(recall: kernel mappings — high addresses)

exec step 2a: allocate memory for executable pages

allocvm() in loop

new physical pages chosen by kalloc()

exec step 2b: load executable pages from executable file

loadvm() in a loop

copy from disk into newly allocated pages (in loadvm())

exec step 3: allocate pages for heap, stack (allocvm() calls)

# xv6: setting process page tables (exec())

exec step 1: create new page table with kernel mappings

setupkvm()

(recall: kernel mappings — high addresses)

exec step 2a: allocate memory for executable pages

allocvm() in loop

new physical pages chosen by kalloc()

exec step 2b: load executable pages from executable file

loadvm() in a loop

copy from disk into newly allocated pages (in loadvm())

exec step 3: allocate pages for heap, stack (allocvm() calls)

# create new page table (kernel mappings)

```
pde_t*
setupkvm(void)
{
    pde_t *pgdir;
    struct kmap *k;

    if((pgdir = (pde_t*)kalloc()) == 0)
        return 0;
    memset(pgdir, 0, PGSIZE);
    if (P2V(PHYSTOP) > (void*)DEVSPACE)
        panic("PHYSTOP_too_high");
    for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
        if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
                    (uint)k->phys_start, k->perm) < 0) {
            freevm(pgdir);
            return 0;
        }
    return pgdir;
}
```

```
}
```

# create new page table (kernel mappings)

allocate first-level page table  
("page directory")

```
pde_t*
setupkvm(void)
{
    pde_t *pgdir;
    struct kmap *k;

    if((pgdir = (pde_t*)kalloc()) == 0)
        return 0;
    memset(pgdir, 0, PGSIZE);
    if (P2V(PHYSTOP) > (void*)DEVSPACE)
        panic("PHYSTOP_too_high");
    for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
        if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
                    (uint)k->phys_start, k->perm) < 0) {
            freevm(pgdir);
            return 0;
        }
    return pgdir;
}
```

# create new page table (kernel mappings)

initialize to 0 — every page invalid

```
pde_t*
setupkvm(void)
{
    pde_t *pgdir;
    struct kmap *k;

    if((pgdir = (pde_t*)kalloc()) == 0)
        return 0;
    memset(pgdir, 0, PGSIZE);
    if (P2V(PHYSTOP) > (void*)DEVSPACE)
        panic("PHYSTOP_too_high");
    for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
        if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
                    (uint)k->phys_start, k->perm) < 0) {
            freevm(pgdir);
            return 0;
        }
    return pgdir;
}
```

# create new page table (kernel mappings)

iterate through list of kernel-space mappings  
everything above address 0x8000 0000

```
pde_t*
setupkvm(void)
{
    pde_t *pgdir;
    struct kmap *k;

    if((pgdir = (pde_t*)kalloc()) == 0)
        return 0;
    memset(pgdir, 0, PGSIZE);
    if (P2V(PHYSTOP) > (void*)DEVSPACE)
        panic("PHYSTOP_too_high");
    for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
        if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
                    (uint)k->phys_start, k->perm) < 0) {
            freevm(pgdir);
            return 0;
        }
    return pgdir;
}
```



# create new page table (kernel mappings)

on failure (no space for new second-level page tables)  
free everything

```
pde_t*
setupkvm(void)
{
    pde_t *pgdir;
    struct kmap *k;

    if((pgdir = (pde_t*)kalloc()) == 0)
        return 0;
    memset(pgdir, 0, PGSIZE);
    if (P2V(PHYSTOP) > (void*)DEVSPACE)
        panic("PHYSTOP_too_high");
    for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
        if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
                    (uint)k->phys_start, k->perm) < 0) {
            freevm(pgdir);
            return 0;
        }
    return pgdir;
}
```

# xv6: setting process page tables (exec())

exec step 1: create new page table with kernel mappings

setupkvm()

(recall: kernel mappings — high addresses)

exec step 2a: allocate memory for executable pages

allocvm() in loop

new physical pages chosen by kalloc()

exec step 2b: **load executable pages from executable file**

loadvm() in a loop

copy from disk into newly allocated pages (in loadvm())

exec step 3: allocate pages for heap, stack (allocvm() calls)

# reading executables (headers)

xv6 executables contain list of sections to load, represented by:

```
struct proghdr {
    uint type;           /* <-- debugging-only or not? */
    uint off;           /* <-- location in file */
    uint vaddr;         /* <-- location in memory */
    uint paddr;         /* <-- confusing ignored field */
    uint filesz;        /* <-- amount to load */
    uint memsz;         /* <-- amount to allocate */
    uint flags;         /* <-- readable/writable (ignored) */
    uint align;
};
```

# reading executables (headers)

xv6 executables contain list of sections to load, represented by:

```
struct proghdr {
    uint type;           /* <-- debugging-only or not? */
    uint off;           /* <-- location in file */
    uint vaddr;         /* <-- location in memory */
    uint paddr;         /* <-- confusing ignored field */
    uint filesz;        /* <-- amount to load */
    uint memsz;         /* <-- amount to allocate */
    uint flags;         /* <-- readable/writable (ignored) */
    uint align;
};

...
if((sz = allocuvm(pgdir, sz, ph.vaddr + ph.memsz)) == 0)
    goto bad;

...
if(loaduvm(pgdir, (char*)ph.vaddr, ip, ph.off, ph.filesz) < 0)
    goto bad;
```

# allocating user pages

```
allocuvm(pde_t *pgdir, uint oldsz, uint newsz)
{
    ...
    a = PGROUNDUP(oldsz);
    for(; a < newsz; a += PGSIZE){
        mem = kalloc();
        if(mem == 0){
            cprintf("allocuvm_out_of_memory\n");
            deallocuvm(pgdir, newsz, oldsz);
            return 0;
        }
        memset(mem, 0, PGSIZE);
        if(mappages(pgdir, (char*)a, PGSIZE, V2P(mem), PTE_W|PTE_U) < 0)
            cprintf("allocuvm_out_of_memory_(2)\n");
        deallocuvm(pgdir, newsz, oldsz);
        kfree(mem);
        return 0;
    }
}
```

# allocating user pages

allocate a new, zero page

```
allocuvm(pde_t *pgdir, uint oldsz, uint newsz)
{
    ...
    a = PGROUNDUP(oldsz);
    for(; a < newsz; a += PGSIZE){
        mem = kalloc();
        if(mem == 0){
            cprintf("allocuvm_out_of_memory\n");
            deallocuvm(pgdir, newsz, oldsz);
            return 0;
        }
        memset(mem, 0, PGSIZE);
        if(mappages(pgdir, (char*)a, PGSIZE, V2P(mem), PTE_W|PTE_U) < 0)
            cprintf("allocuvm_out_of_memory_(2)\n");
        deallocuvm(pgdir, newsz, oldsz);
        kfree(mem);
        return 0;
    }
}
```

# allocating user pages

add page to second-level page table

```
allocuvm(pde_t *pgdir, uint oldsz, uint newsz)
{
    ...
    a = PGROUNDUP(oldsz);
    for(; a < newsz; a += PGSIZE){
        mem = kalloc();
        if(mem == 0){
            cprintf("allocuvm_out_of_memory\n");
            deallocuvm(pgdir, newsz, oldsz);
            return 0;
        }
        memset(mem, 0, PGSIZE);
        if(mappages(pgdir, (char*)a, PGSIZE, V2P(mem), PTE_W|PTE_U) < 0)
            cprintf("allocuvm_out_of_memory_(2)\n");
        deallocuvm(pgdir, newsz, oldsz);
        kfree(mem);
        return 0;
    }
}
```

# allocating user pages

same function used to allocate memory for heap

```
allocuvm(pde_t *pgdir, uint oldsz, uint newsz)
{
    ...
    a = PGROUNDUP(oldsz);
    for(; a < newsz; a += PGSIZE){
        mem = kalloc();
        if(mem == 0){
            cprintf("allocuvm_out_of_memory\n");
            deallocuvm(pgdir, newsz, oldsz);
            return 0;
        }
        memset(mem, 0, PGSIZE);
        if(mappages(pgdir, (char*)a, PGSIZE, V2P(mem), PTE_W|PTE_U) < 0)
            cprintf("allocuvm_out_of_memory_(2)\n");
        deallocuvm(pgdir, newsz, oldsz);
        kfree(mem);
        return 0;
    }
}
```



# reading executables (headers)

xv6 executables contain list of sections to load, represented by:

```
struct proghdr {
    uint type;           /* <-- debugging-only or not? */
    uint off;           /* <-- location in file */
    uint vaddr;         /* <-- location in memory */
    uint paddr;         /* <-- confusing ignored field */
    uint filesz;        /* <-- amount to load */
    uint memsz;         /* <-- amount to allocate */
    uint flags;         /* <-- readable/writable (ignored) */
    uint align;
};
```

```
...
if((sz = allocuvm(pgdir, sz, ph.vaddr + ph.memsz)) == 0)
    goto bad;
```

```
...
if(loaduvm(pgdir, (char*)ph.vaddr, ip, ph.off, ph.filesz) < 0)
    goto bad;
```

# loading user pages from executable

```
loadvm(pde_t *pgdir, char *addr, struct inode *ip, uint offset, uint sz)
{
    ...
    for(i = 0; i < sz; i += PGSIZE){
        if((pte = walkpgdir(pgdir, addr+i, 0)) == 0)
            panic("loadvm:_address_should_exist");
        pa = PTE_ADDR(*pte);
        if(sz - i < PGSIZE)
            n = sz - i;
        else
            n = PGSIZE;
        if(readi(ip, P2V(pa), offset+i, n) != n)
            return -1;
    }
    return 0;
}
```

# loading user pages from executable

get page table entry being loaded  
already allocated earlier  
look up address to load into

```
loaduvm(pde_t *pgdir, char *addr
```

, uir

```
{  
    ...  
    for(i = 0; i < sz; i += PGSIZE){  
        if((pte = walkpgdir(pgdir, addr+i, 0)) == 0)  
            panic("loaduvm:_address_should_exist");  
        pa = PTE_ADDR(*pte);  
        if(sz - i < PGSIZE)  
            n = sz - i;  
        else  
            n = PGSIZE;  
        if(readi(ip, P2V(pa), offset+i, n) != n)  
            return -1;  
    }  
    return 0;  
}
```

# loading user pages from executable

```
loadvm(pde_t *pgdir, char *addr, struct inode *ip, uint offset, uint sz)
{
    ...
    for(i = 0; i < sz; i += PGSIZE){
        if((pte = walkpgdir(pgdir, addr+i, 0)) == 0)
            panic("loadvm: address should exist");
        pa = PTE_ADDR(*pte);
        if(sz - i < PGSIZE)
            n = sz - i;
        else
            n = PGSIZE;
        if(readi(ip, P2V(pa), offset+i, n) != n)
            return -1;
    }
    return 0;
}
```

**exercise:** why don't we just use `addr` directly?  
(instead of turning it into a physical address,  
then into a virtual address again)

# loading user pages from executable

copy from file (represented by struct inode) into memory  
P2V(pa) — mapping of physical addresss in kernel memory

```
loaduv  
{
```

```
    ...  
    for(i = 0; i < sz; i += PGSIZE){  
        if((pte = walkpgdir(pgdir, addr+i, 0)) == 0)  
            panic("loaduvm:_address_should_exist");  
        pa = PTE_ADDR(*pte);  
        if(sz - i < PGSIZE)  
            n = sz - i;  
        else  
            n = PGSIZE;  
        if(readi(ip, P2V(pa), offset+i, n) != n)  
            return -1;  
    }  
    return 0;
```

```
}
```

# kalloc/kfree

kalloc/kfree — xv6's physical memory allocator

allocates/deallocates **whole pages only**

keep linked list of free pages

- list nodes — stored in corresponding free page itself

- kalloc — return first page in list

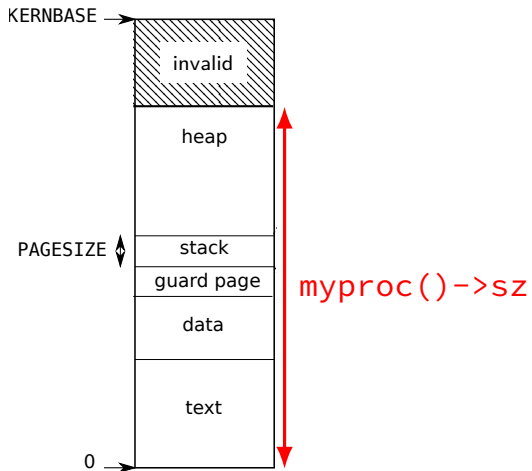
- kfree — add page to list

linked list created at boot

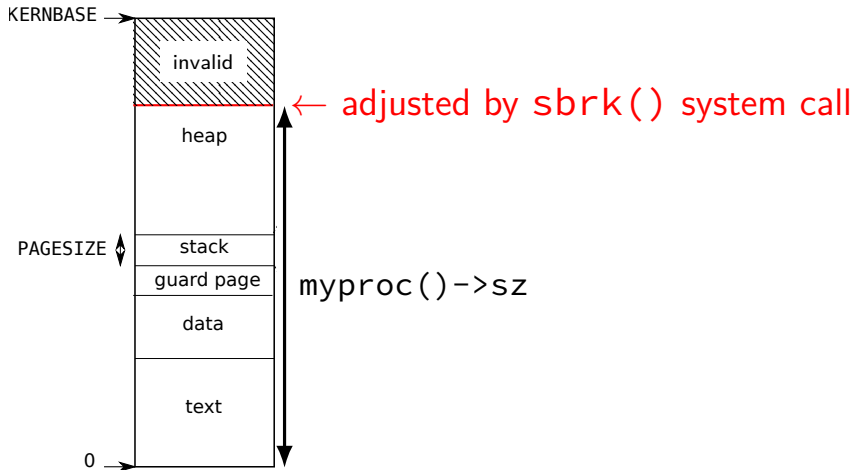
usable memory fixed size (224MB)

- determined by PHYSTOP in memlayout.h

# xv6 program memory



# xv6 program memory





# xv6 heap allocation

xv6: every process has a heap at the top of its address space  
yes, this is unlike Linux where heap is below stack

tracked in `struct proc` with `sz`  
= last valid address in process

position changed via `sbrk(amount)` system call  
sets `sz += amount`  
same call exists in Linux, etc. — but also others

# sbrk

```
sys_sbrk()
{
    if(argint(0, &n) < 0)
        return -1;
    addr = myproc()->sz;
    if(growproc(n) < 0)
        return -1;
    return addr;
}
```

# sbrk

SZ: current top of heap

```
sys_sbrk()
{
    if(argint(0, &n) < 0)
        return -1;
    addr = myproc()->sz;
    if(growproc(n) < 0)
        return -1;
    return addr;
}
```

# sbrk

sbrk(N): grow heap by  $N$  (shrink if negative)

```
sys_sbrk()
{
    if(argint(0, &n) < 0)
        return -1;
    addr = myproc()->sz;
    if(growproc(n) < 0)
        return -1;
    return addr;
}
```

# sbrk

returns old top of heap (or -1 on out-of-memory)

```
sys_sbrk()
{
    if(argint(0, &n) < 0)
        return -1;
    addr = myproc()->sz;
    if(growproc(n) < 0)
        return -1;
    return addr;
}
```

# growproc

```
growproc(int n)
{
    uint sz;
    struct proc *curproc = myproc();

    sz = curproc->sz;
    if(n > 0){
        if((sz = allocuvm(curproc->pgdir, sz, sz + n)) == 0)
            return -1;
    } else if(n < 0){
        if((sz = deallocuvm(curproc->pgdir, sz, sz + n)) == 0)
            return -1;
    }
    curproc->sz = sz;
    switchuvm(curproc);
    return 0;
}
```

# growproc

allocuvm — same function used to allocate initial space  
maps pages for addresses SZ to SZ + n  
calls kalloc to get each page

```
growproc(int n)
{
    uint sz;
    struct proc *curproc = myproc();

    sz = curproc->sz;
    if(n > 0){
        if((sz = allocuvm(curproc->pgdir, sz, sz + n)) == 0)
            return -1;
    } else if(n < 0){
        if((sz = deallocuvm(curproc->pgdir, sz, sz + n)) == 0)
            return -1;
    }
    curproc->sz = sz;
    switchuvm(curproc);
    return 0;
}
```

## xv6 page faults (now)

fault from accessing page table entry marked 'not-present'

xv6: prints an error and kills process:

```
*((int*) 0x800444) = 1;
```

```
...
```

```
/* in trap.c: */
```

```
    cprintf("pid_%d_%s: trap_%d_err_%d_on_cpu_%d_"  
            "eip_0x%x_addr_0x%x--kill_proc\n",  
            myproc()->pid, myproc()->name, tf->trapno,  
            tf->err, cpuid(), tf->eip, rcr2());  
    myproc()->killed = 1;
```

```
pid 4 processname: trap 14 err 6 on cpu 0 eip 0x1a addr 0x800444--k-
```

```
14 = T_PGFLT
```

special register CR2 contains faulting address



## xv6 page faults (now)

fault from accessing page table entry marked 'not-present'

xv6: prints an error and kills process:

```
*((int*) 0x800444) = 1;
```

```
...
```

```
/* in trap.c: */
```

```
    cprintf("pid_%d_%s:_trap_%d_err_%d_on_cpu_%d_"  
            "eip_0x%x_addr_0x%x--kill_proc\n",  
            myproc()->pid, myproc()->name, tf->trapno,  
            tf->err, cpuid(), tf->eip, rcr2());  
    myproc()->killed = 1;
```

```
pid 4 processname: trap 14 err 6 on cpu 0 eip 0x1a addr 0x800444--k-
```

```
14 = T_PGFLT
```

special register CR2 contains faulting address

## xv6 page faults (now)

fault from accessing page table entry marked 'not-present'

xv6: prints an error and kills process:

```
*((int*) 0x800444) = 1;
```

```
...
```

```
/* in trap.c: */
```

```
    cprintf("pid_%d_%s: trap_%d_err_%d_on_cpu_%d_"  
            "eip_0x%x_addr_0x%x--kill_proc\n",  
            myproc()->pid, myproc()->name, tf->trapno,  
            tf->err, cpuid(), tf->eip, rcr2());  
    myproc()->killed = 1;
```

```
pid 4 processname: trap 14 err 6 on cpu 0 eip 0x1a addr 0x800444--k-
```

```
14 = T_PGFLT
```

special register CR2 contains faulting address

## xv6: if one handled page faults

returning from page fault handler without killing process

...retries the failing instruction

can use to update the page table — “just in time”

```
if (tf->trapno == T_PGFLT) {
    void *address = (void *) rcr2();
    if (is_address_okay(myproc(), address)) {
        setup_page_table_entry_for(myproc(), address);
        // return from fault, retry access
    } else {
        // actual segfault, kill process
        cprintf("...");
        myproc()->killed = 1;
    }
}
```

## xv6: if one handled page faults

check *process control block* to see if access okay

returning from page fault handler without killing process

...retries the failing instruction

can use to update the page table — “just in time”

```
if (tf->trapno == T_PGFLT) {
    void *address = (void *) rcr2();
    if (is_address_okay(myproc(), address)) {
        setup_page_table_entry_for(myproc(), address);
        // return from fault, retry access
    } else {
        // actual segfault, kill process
        cprintf("...");
        myproc()->killed = 1;
    }
}
```

## xv6: if one handled page faults

returning from page

if so, setup the page table so it works next time  
i.e. immediately after returning from fault

...retries the failing instruction

can use to update the page table — “just in time”

```
if (tf->trapno == T_PGFLT) {
    void *address = (void *) rcr2();
    if (is_address_okay(myproc(), address)) {
        setup_page_table_entry_for(myproc(), address);
        // return from fault, retry access
    } else {
        // actual segfault, kill process
        cprintf("...");
        myproc()->killed = 1;
    }
}
```