# last time

storing page tables in memory

as array of *page table entries*

page table entries encoded as integer

*page table base register* — OS tells CPU base address of array

each program memory access = two real ones

preview: multi-level page table lookup

# lab this week

no lab

a subset of TAs will be in Rice 130 to give office-hour-type help

# anonymous feedback (1)

got this last Tuesday, missed addressing last Thursday (sorry!)

"Can you make the quizzes have more number of questions and/or each MCQ be worth less points? The percentage scores decreases drastically by making just one mistake. The quiz is also worth a huge portion of the final grade and being able to do better on it would help get a good grade."

> probably should be some more questions, but I think I'd get more complaints if I really embraced more questions
>
> would like to use comments to give more naunce in quiz grades

## quiz Q4

0x300010: movq %rax, (%rcx)

to execute:

    access 0x300010 to read machine code — VPN 0x300
    read %rax (no memory access) = 0x999000
    read %rcx (no memory access) = 0x123450
    write to 0x123450 to write %rax value — VPN 0x123

## quiz Q5

accessing 0x30110 = 0x30000 + 0x110 = 0x30000 + 0x44 * 4

VPN = 0x44

page offset is 0x9433 (same in physical + virtual)

# quiz Q6

usually I would say something that causes exception != can access

re: virtual/physical address bits
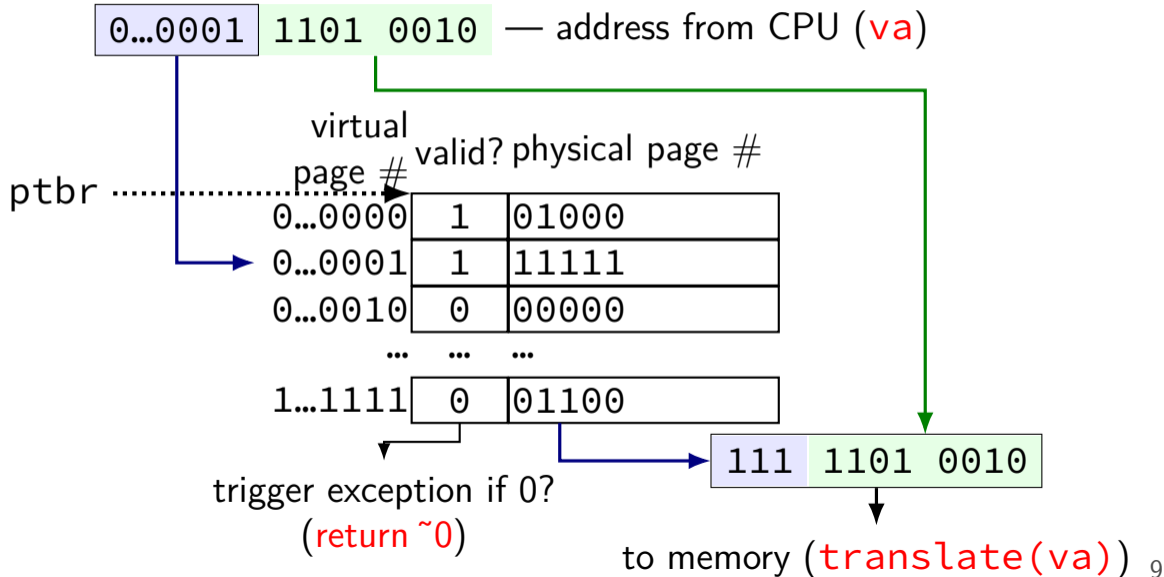   could have most page table entries be invalid (in both)

# anonymous feedback (2)

"I was hoping to have a little lecture time dedicated to explaining the nuances to the homework assignment, especially a clarification for what ptbr points to, and the difference between VPN and Physical page numbers in the context of this homework assignment"

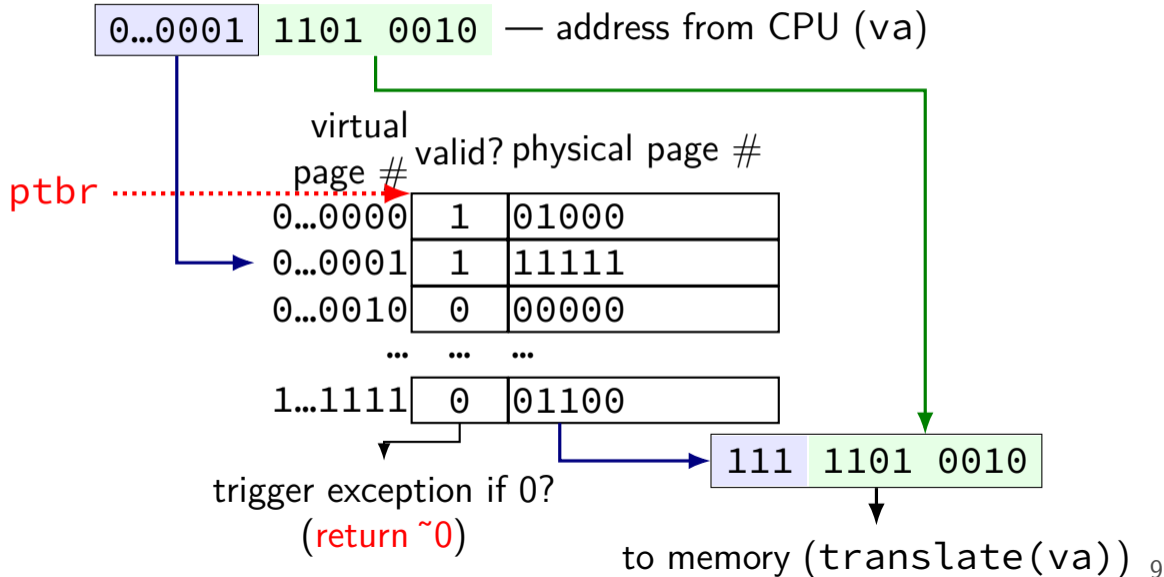"Do you mind re-visiting and explaining page_allocate()'s behavior?"

# page table lookup (and translate())

`0…0001` `1101 0010` — address from CPU (`va`)

virtual
page #     valid? physical page #

`ptbr` ·······→
`0…0000` | 1 | 01000
`0…0001` | 1 | 11111
`0…0010` | 0 | 00000
...   ...   ...
`1…1111` | 0 | 01100

trigger exception if 0?
(`return ~0`)

`111` `1101 0010`

to memory (`translate(va)`)

# page table lookup (and translate())

# page table lookup (and allocate)

`0…0001` `1101 0010` — address from CPU (`va`)

`ptbr` ·····················▸

page_allocate(va) — set ptbr if unset

trigger exception if 0?
(return ˜0)

`111` `1101 0010`

to memory (`translate(va)`) 10

# page table lookup (and allocate)



0…0001 1101 0010 — address from CPU (va)

ptbr

virtual page #
valid? physical page #

| virtual page # | valid? | physical page # |
|---|---|---|
| 0…0000 | 1 | 01000 |
| 0…0001 | 0 | 00000 |
| 0…0010 | 0 | 00000 |
| … | … | … |
| 1…1111 | 0 | 01100 |

page_allocate(va) — set page table entry if unset

trigger exception if 0?
(return ~0)

111 1101 0010

to memory (translate(va))

10

# ptbr in assignment

```
size_t ptbr;
```

points to beginning of (primary) page table

initially $0$ = doesn't point to anything

filled in by first call to page_allocate
or by testing code

# typical timings

| task | typical time (order of magnitude) |
| --- | --- |
| empty function | nanosecond |
| getppid (syscall) | microsecond |
| system(/bin/true) (run program) | milliseconds |
| start signal handler | microseconds |
| signal ping/pong (context switch?) | tens of microseconds |

# exercise: 64-bit system

my desktop: 39-bit physical addresses; 48-bit virtual addresses

4096 byte pages

# exercise: 64-bit system

my desktop: 39-bit physical addresses; 48-bit virtual addresses

4096 byte pages

top 16 bits of 64-bit addresses not used for translation

# exercise: 64-bit system

my desktop: 39-bit physical addresses; 48-bit virtual addresses

4096 byte pages

exercise: how many page table entries? (assuming page table like shown before)

exercise: how large are physical page numbers?

# exercise: 64-bit system

my desktop: 39-bit physical addresses; 48-bit virtual addresses

4096 byte pages

exercise: how many page table entries? (assuming page table like shown before) $2^{48}/2^{12} = 2^{36}$ entries

exercise: how large are physical page numbers? $39 - 12 = 27$ bits

# exercise: **64-bit system**

my desktop: 39-bit physical addresses; 48-bit virtual addresses

4096 byte pages

exercise: how many page table entries? (assuming page table like shown before) $2^{48}/2^{12} = 2^{36}$ entries

exercise: how large are physical page numbers? $39 - 12 = 27$ bits

page table entries are 8 bytes (room for expansion, metadata)
    trick: power of two size makes table lookup faster
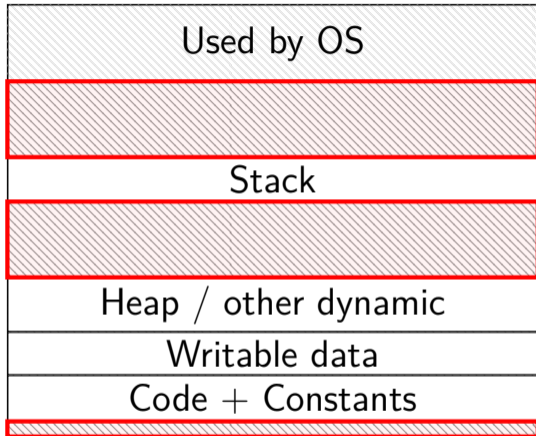
would take up $2^{39}$ bytes?? (512GB??)

# huge page tables

huge virtual address spaces!

impossible to store PTE for every page

how can we save space?

## holes

| |
|---|
| Used by OS |
| |
| Stack |
| |
| Heap / other dynamic |
| Writable data |
| Code + Constants |

most pages are invalid

## saving space

basic idea: don't store (most) invalid page table entries

use a data structure other than a flat array
    want a map — lookup key (virtual page number), get value (PTE)

options?

# saving space

basic idea: don't store (most) invalid page table entries

use a data structure other than a flat array
  want a map — lookup key (virtual page number), get value (PTE)

options?

hashtable
  actually used by some historical processors
  but never common

# saving space

basic idea: don't store (most) invalid page table entries

use a data structure other than a flat array
> want a map — lookup key (virtual page number), get value (PTE)

options?

hashtable
> actually used by some historical processors
> but never common

tree data structure
> but not quite a search tree

# search tree tradeoffs

lookup usually implemented in hardware

 lookup should be simple

 solution: lookup splits up address bits (no complex calculations)

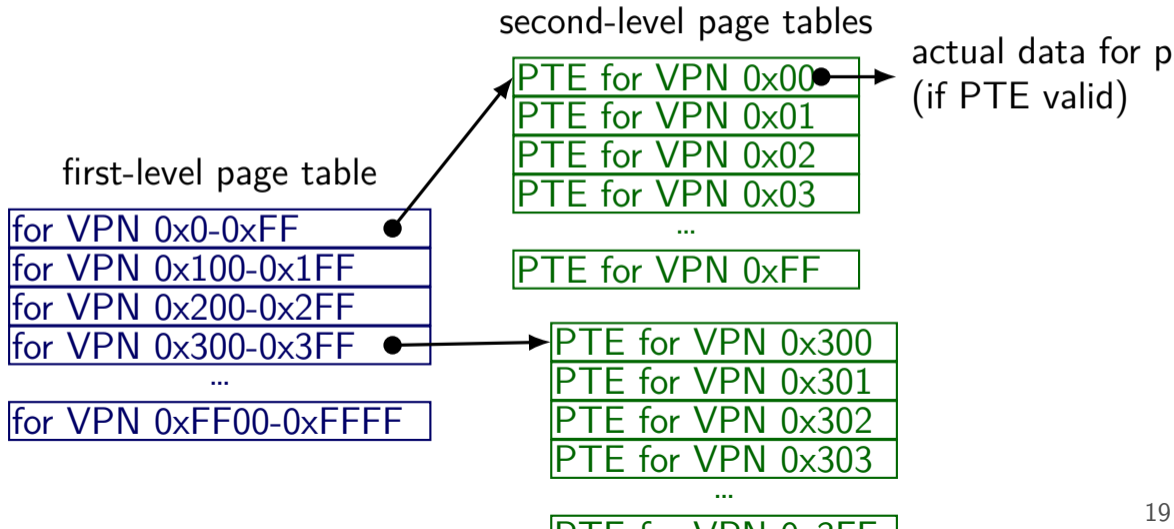lookup should not involve many memory accesses

 doing two memory accesses is already very slow

 solution: tree with many children from each node

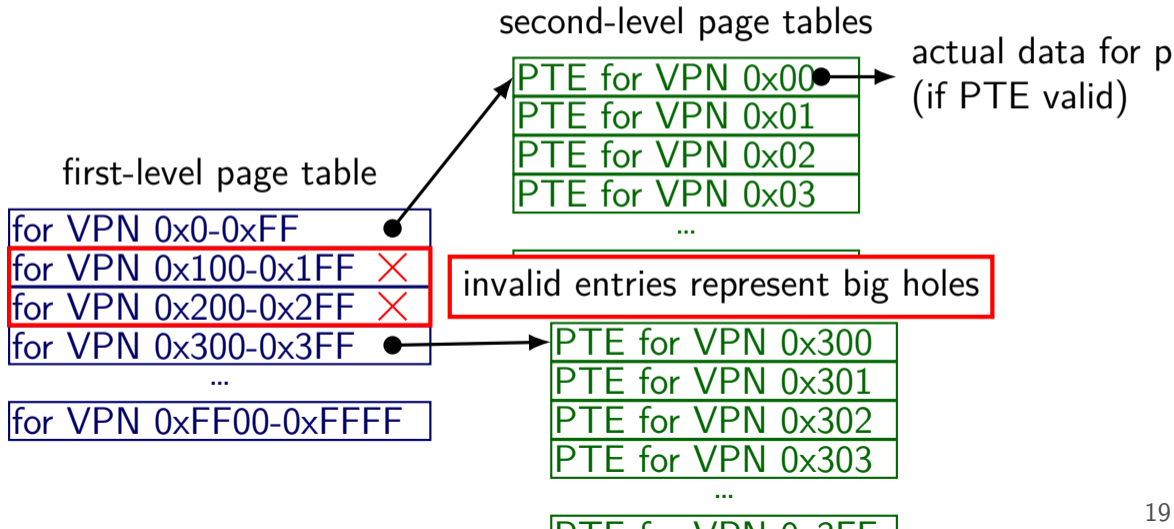  (far from binary tree's left/right child)

# two-level page tables

two-level page tables for 65536 pages (16-bit VPN; 256 entries/table)

second-level page tables

PTE for VPN 0x00 → actual data for p...
(if PTE valid)

PTE for VPN 0x01
PTE for VPN 0x02
PTE for VPN 0x03
...
PTE for VPN 0xFF

first-level page table

for VPN 0x0-0xFF
for VPN 0x100-0x1FF
for VPN 0x200-0x2FF
for VPN 0x300-0x3FF
...
for VPN 0xFF00-0xFFFF

PTE for VPN 0x300
PTE for VPN 0x301
PTE for VPN 0x302
PTE for VPN 0x303
...

19

# two-level page tables

two-level page tables for 65536 pages (16-bit VPN; 256 entries/table)

second-level page tables

PTE for VPN 0x00 → actual data for p...
(if PTE valid)

PTE for VPN 0x01
PTE for VPN 0x02
PTE for VPN 0x03
...

first-level page table

for VPN 0x0-0xFF ●
for VPN 0x100-0x1FF ✕
for VPN 0x200-0x2FF ✕

invalid entries represent big holes

for VPN 0x300-0x3FF ●→
PTE for VPN 0x300
PTE for VPN 0x301
PTE for VPN 0x302
PTE for VPN 0x303
...

...
for VPN 0xFF00-0xFFFF

# two-level page tables

two-level page tables for 65536 pages (16-bit VPN; 256 entries/table)

**first-level page table**

| VPN range | valid | … | physical page # (of next page table) |
|-----------|-------|---|--------------------------------------|
| 0x0000-0x00FF | 1 | … | 0x22343 |
| 0x0100-0x01FF | 0 | … | 0x00000 |
| 0x0200-0x02FF | 0 | … | 0x00000 |
| 0x0300-0x03FF | 1 | … | 0x33454 |
| 0x0400-0x04FF | 1 | … | 0xFF043 |
| … | … | … | … |
| 0xFF00-0xFFFF | 1 | … | 0xFF045 |

for p...
...d)

first-level pag...

for VPN 0x0-0xF...
for VPN 0x100-0...
for VPN 0x200-0...
for VPN 0x300-0...
…
for VPN 0xFF00...

PTE for VPN 0x303
…
PTE for VPN 0x3FF

# two-level page tables

two-level page tables for 65536 pages (16-bit VPN; 256 entries/table)

first-level pag[...]

for VPN 0x0-0xF[...]
for VPN 0x100-0[...]
for VPN 0x200-0[...]
for VPN 0x300-0[...]
...
for VPN 0xFF00[...]

**first-level page table**

| VPN range | valid | ... | physical page # (of next page table) |
|---|---|---|---|
| 0x0000–0x00FF | 1 | ... | 0x22343 |
| 0x0100–0x01FF | 0 | ... | 0x00000 |
| 0x0200–0x02FF | 0 | ... | 0x00000 |
| 0x0300–0x03FF | 1 | ... | 0x33454 |
| 0x0400–0x04FF | 1 | ... | 0xFF043 |
| ... | ... | ... | ... |
| 0xFF00–0xFFFF | 1 | ... | 0xFF045 |

for p[...]
d)

PTE for VPN 0x303
...
PTE f[...] VPN 0x3FF[...]

# two-level page tables

two-level page tables for 65536 pages (16-bit VPN; 256 entries/table)

first-level page ~~~~~~ for p~
d)

first-level pag~

| for VPN 0x0-0xF |
| for VPN 0x100-0 |
| for VPN 0x200-0 |
| for VPN 0x300-0 |
| ... |
| for VPN 0xFF00 |

## first-level page table

| VPN range | valid | ... | physical page #<br>(of next page table) |
|-----------|-------|-----|------------------------------------------|
| 0x0000-0x00FF | 1 | ... | 0x22343 |
| 0x0100-0x01FF | 0 | ... | 0x00000 |
| 0x0200-0x02FF | 0 | ... | 0x00000 |
| 0x0300-0x03FF | 1 | ... | 0x33454 |
| 0x0400-0x04FF | 1 | ... | 0xFF043 |
| ... | ... | ... | ... |
| 0xFF00-0xFFFF | 1 | ... | 0xFF045 |

PTE for VPN 0x303
...
PTE for VPN 0x3FF

# two-level page tables

two-level page tables for 65536 pages (16-bit VPN; 256 entries/table)
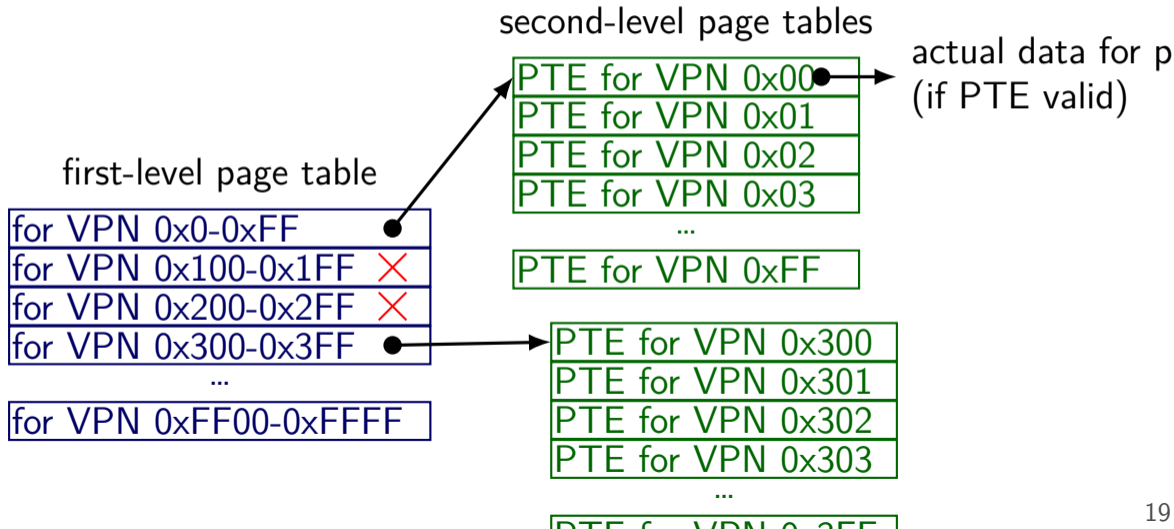
**a second-level page table**

| VPN | valid | ... | physical page # (of data) |
|-----|-------|-----|---------------------------|
| 0x300 | 0 | 1 | 0x42443 |
| 0x301 | 0 | 1 | 0x4A9DE |
| 0x302 | 0 | 1 | 0x5C001 |
| 0x303 | 0 | 1 | 0x00000 |
| 0x304 | 0 | 1 | 0x6C223 |
| ... | ... | ... | ... |
| 0x3FF | ... | 1 | 0x00000 |

first-level page table

| |
|---|
| for VPN 0x0-0xFF ● |
| for VPN 0x100-0x1FF ✗ |
| for VPN 0x200-0x2FF ✗ |
| for VPN 0x300-0x3FF ● |
| ... |
| for VPN 0xFF00-0xFFFF |

PTE for VPN 0x303
...

PTE for VPN 0x303

# two-level page tables

two-level page tables for 65536 pages (16-bit VPN; 256 entries/table)

first-level page table

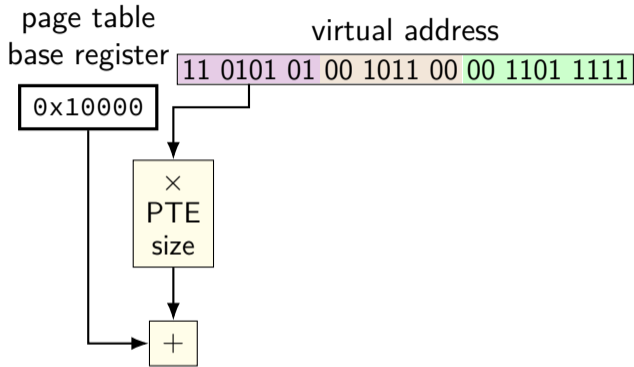| for VPN 0x0-0xFF | ● |
|---|---|
| for VPN 0x100-0x1FF | ✗ |
| for VPN 0x200-0x2FF | ✗ |
| for VPN 0x300-0x3FF | ● |
| ... | |
| for VPN 0xFF00-0xFFFF | |

**a second-level page table**

| VPN | valid | ... | physical page # (of data) |
|---|---|---|---|
| 0x300 | 0 | 1 | 0x42443 |
| 0x301 | 0 | 1 | 0x4A9DE |
| 0x302 | 0 | 1 | 0x5C001 |
| 0x303 | 0 | 1 | 0x00000 |
| 0x304 | 0 | 1 | 0x6C223 |
| ... | ... | ... | ... |
| 0x3FF | ... | 1 | 0x00000 |

PTE for VPN 0x303
...

# two-level page tables

two-level page tables for 65536 pages (16-bit VPN; 256 entries/table)

second-level page tables

PTE for VPN 0x00 ●——→ actual data for p
(if PTE valid)

PTE for VPN 0x01
PTE for VPN 0x02
PTE for VPN 0x03
…
PTE for VPN 0xFF

first-level page table

| for VPN 0x0-0xFF | ● |
| for VPN 0x100-0x1FF | ✗ |
| for VPN 0x200-0x2FF | ✗ |
| for VPN 0x300-0x3FF | ● |
| … | |
| for VPN 0xFF00-0xFFFF | |

PTE for VPN 0x300
PTE for VPN 0x301
PTE for VPN 0x302
PTE for VPN 0x303
…

# two-level page table lookup

virtual address

`11 0101 01 00 1011 00` `00 1101 1111`

VPN — split into two parts (one per level)

this example: parts equal sized — common, but not required

# two-level page table lookup

page table
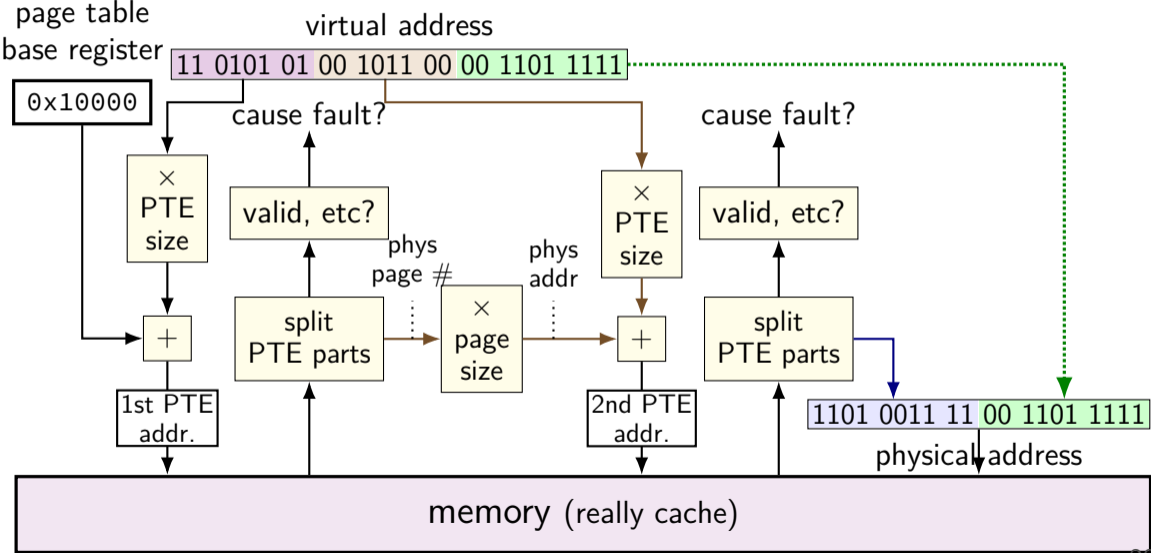base register    virtual address

11 0101 01 00 1011 00 00 1101 1111

0x10000

×
PTE
size

+

# two-level page table lookup

page table
base register

virtual address

0x10000

11 0101 01 00 1011 00 00 1101 1111

cause fault?

× PTE size

valid, etc?

+

split PTE parts
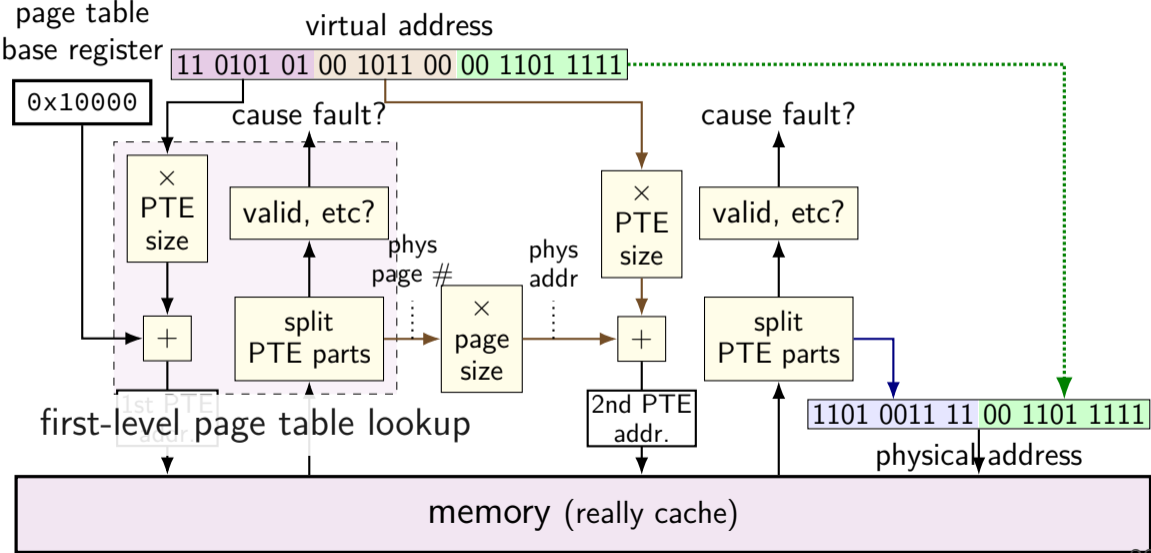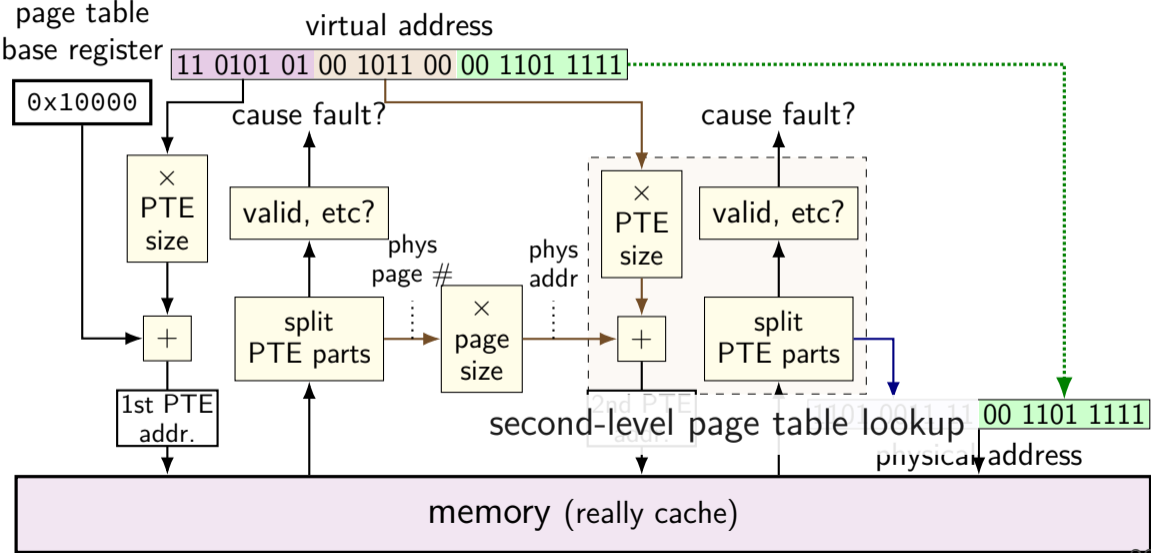
1st PTE addr.

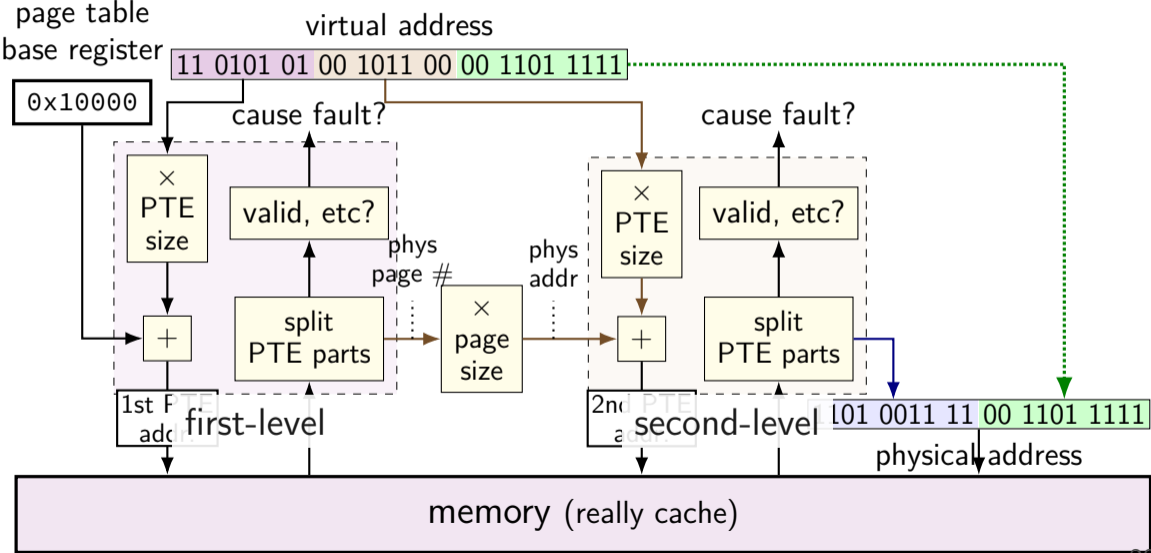physical address

memory (really cache)

# two-level page table lookup

# two-level page table lookup

# two-level page table lookup



page table base register

0x10000

virtual address

11 0101 01 00 1011 00 00 1101 1111

cause fault?

× PTE size

valid, etc?

split PTE parts

phys page #

× page size

phys addr

cause fault?

× PTE size

valid, etc?

split PTE parts

1st PTE addr.

+

2nd PTE addr.

+

1101 0011 11 00 1101 1111

physical address

memory (really cache)

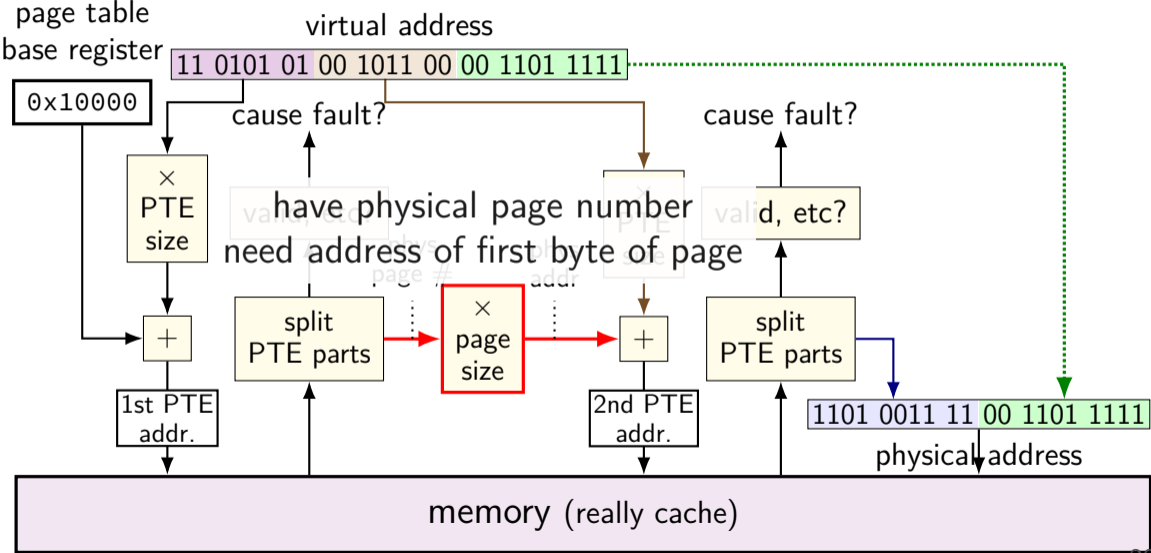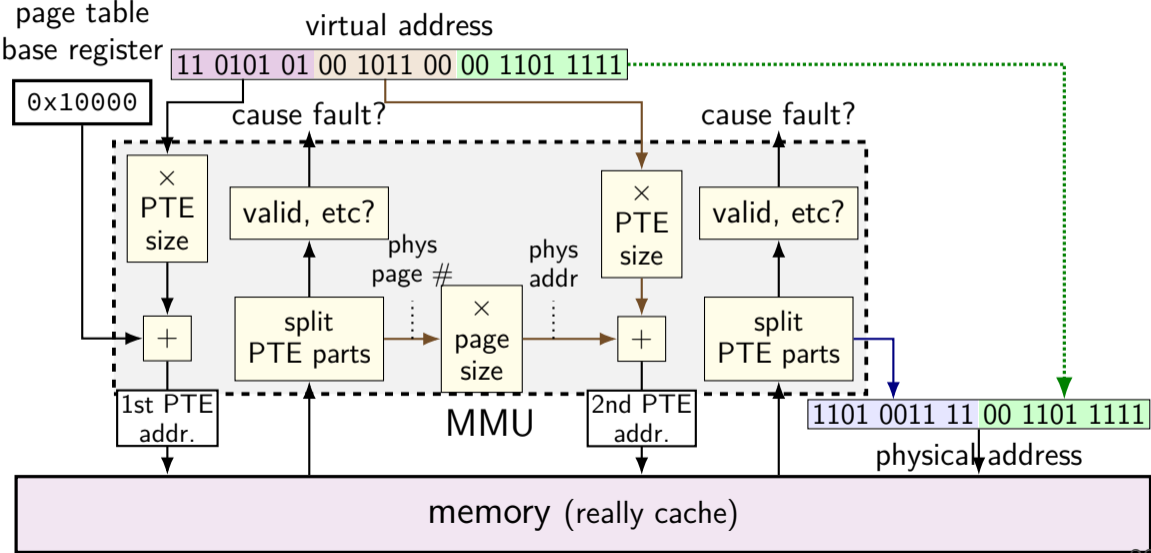# two-level page table lookup

# two-level page table lookup



page table base register

0x10000

virtual address

11 0101 01 00 1011 00 00 1101 1111

cause fault?

× PTE size

valid, etc?

phys page #

split PTE parts

+

1st PTE addr

first-level page table lookup

× page size

phys addr

cause fault?

× PTE size

valid, etc?

+

2nd PTE addr.

split PTE parts

1101 0011 11 00 1101 1111

physical address

memory (really cache)

# two-level page table lookup

# two-level page table lookup



page table base register: 0x10000

virtual address: 11 0101 01 00 1011 00 00 1101 1111

cause fault?

first-level: × PTE size, valid, etc?, split PTE parts, +, 1st PTE addr

phys page #

× page size

phys addr

cause fault?

second-level: × PTE size, valid, etc?, split PTE parts, +, 2nd PTE addr

physical address: 101 0011 11 00 1101 1111

memory (really cache)

20

# two-level page table lookup



page table base register

0x10000

virtual address

11 0101 01 00 1011 00 00 1101 1111

cause fault?

cause fault?

× PTE size

have physical page number
need address of first byte of page

valid, etc?

split PTE parts

× page size

+

split PTE parts

1st PTE addr.

2nd PTE addr.

1101 0011 11 00 1101 1111

physical address

memory (really cache)

# two-level page table lookup

## another view



| VPN part 1 | VPN part 2 | page offset |
|---|---|---|

first-level page table

second-level page table

page table entry

page table entry

physical page

page table base register

# 2-level splitting

9-bit virtual address

6-bit physical address

# 2-level splitting

9-bit virtual address

6-bit physical address

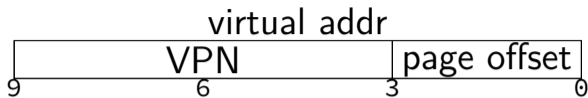8-byte pages → 3-bit page offset (bottom)

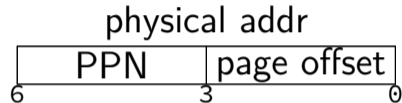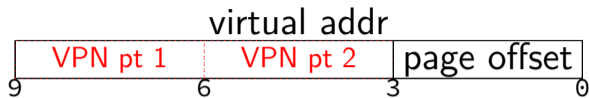9-bit VA: 6 bit VPN + 3 bit PO

6-bit PA: 3 bit PPN + 3 bit PO

virtual addr

| VPN | page offset |
|-----|-------------|

9        6        3        0

physical addr
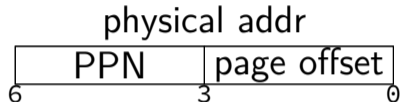
| PPN | page offset |
|-----|-------------|

6        3        0

# 2-level splitting

9-bit virtual address

6-bit physical address

8-byte pages → 3-bit page offset (bottom)

9-bit VA: 6 bit VPN + 3 bit PO

6-bit PA: 3 bit PPN + 3 bit PO

1 page page tables w/ 1 byte entry → 8 entry PTs

virtual addr

| VPN | | page offset |
|---|---|---|
9 | 6 | 3 | 0

physical addr

| PPN | page offset |
|---|---|
6 | 3 | 0

page table (either level)

| | valid? | PPN |
|---|---|---|
| 0 | | |
| 1 | | |
| 2 | | |
| ... | ... | ... |
| 7 | | |

# 2-level splitting

9-bit virtual address

6-bit physical address

8-byte pages → 3-bit page offset (bottom)

9-bit VA: 6 bit VPN + 3 bit PO

6-bit PA: 3 bit PPN + 3 bit PO

1 page page tables w/ 1 byte entry → 8 entry PTs

8 entry page tables → 3-bit VPN parts

9-bit VA: 3 bit VPN part 1; 3 bit VPN part 2

virtual addr

| VPN pt 1 | VPN pt 2 | page offset |
|---|---|---|

9        6        3       0

physical addr

| PPN | page offset |
|---|---|

6       3       0

page table (either level)

| | valid? | PPN |
|---|---|---|
| 0 | | |
| 1 | | |
| 2 | | |
| ... | ... | ... |
| 7 | | |

# 2-level example

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x20; translate virtual address 0x129

| physical addresses | bytes | | physical addresses | bytes |
|---|---|---|---|---|
| 0x00-3 | 00 11 22 33 | | 0x20-3 | 00 91 72 13 |
| 0x04-7 | 44 55 66 77 | | 0x24-7 | F4 A5 36 07 |
| 0x08-B | 88 99 AA BB | | 0x28-B | 89 9A AB BC |
| 0x0C-F | CC DD EE FF | | 0x2C-F | CD DE EF F0 |
| 0x10-3 | 1A 2A 3A 4A | | 0x30-3 | BA 0A BA 0A |
| 0x14-7 | 1B 2B 3B 4B | | 0x34-7 | DB 0B DB 0B |
| 0x18-B | 1C 2C 3C 4C | | 0x38-B | EC 0C EC 0C |
| 0x1C-F | 1C 2C 3C 4C | | 0x3C-F | AC DC DC 0C |

# 2-level example

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x20; translate virtual address 0x129

| physical addresses | bytes |
|---|---|
| 0x00–3 | 00 11 22 33 |
| 0x04–7 | 44 55 66 77 |
| 0x08–B | 88 99 AA BB |
| 0x0C–F | CC DD EE FF |
| 0x10–3 | 1A 2A 3A 4A |
| 0x14–7 | 1B 2B 3B 4B |
| 0x18–B | 1C 2C 3C 4C |
| 0x1C–F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20–3 | 00 91 72 13 |
| 0x24–7 | F4 A5 36 07 |
| 0x28–B | 89 9A AB BC |
| 0x2C–F | CD DE EF F0 |
| 0x30–3 | BA 0A BA 0A |
| 0x34–7 | DB 0B DB 0B |
| 0x38–B | EC 0C EC 0C |
| 0x3C–F | AC DC DC 0C |

```
0x129 = 1 0010 1001
0x20 + 0x4 ×1 = 0x24
```
*PTE 1 value:*
```
0xF4 = 1111 0100
```
PPN 111, valid 1

# 2-level example

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register `0x20`; translate virtual address `0x129`

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | 00 91 72 13 |
| 0x24-7 | F4 A5 36 07 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | AC DC DC 0C |

$0x129 = 1\ 0010\ 1001$

$0x20 + 0x4 \times 1 = 0x24$

*PTE 1 value:*

$0xF4 = 1111\ 0100$

PPN 111, valid 1

*PTE 2 addr:*

$111\ 000 + 101 \times 1 = 0x3D$

*PTE 2 value:* 0xDC

# 2-level example

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register `0x20`; translate virtual address `0x129`

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | 00 91 72 13 |
| 0x24-7 | F4 A5 36 07 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | AC DC DC 0C |

```
0x129 = 1 0010 1001
0x20 + 0x4 ×1 = 0x24
```
*PTE 1 value:*
```
0xF4 = 1111 0100
```
PPN 111, valid 1
*PTE 2 addr:*
```
111 000 + 101 × 1 = 0x3D
```
*PTE 2 value:* `0xDC`
PPN 110; valid 1
`M[110 001 (0x31)] = 0x0A`

23

# 2-level example

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register `0x20`; translate virtual address `0x129`

| physical addresses | bytes |
|---|---|
| 0x00–3 | 00 11 22 33 |
| 0x04–7 | 44 55 66 77 |
| 0x08–B | 88 99 AA BB |
| 0x0C–F | CC DD EE FF |
| 0x10–3 | 1A 2A 3A 4A |
| 0x14–7 | 1B 2B 3B 4B |
| 0x18–B | 1C 2C 3C 4C |
| 0x1C–F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20–3 | 00 91 72 13 |
| 0x24–7 | F4 A5 36 07 |
| 0x28–B | 89 9A AB BC |
| 0x2C–F | CD DE EF F0 |
| 0x30–3 | BA 0A BA 0A |
| 0x34–7 | DB 0B DB 0B |
| 0x38–B | EC 0C EC 0C |
| 0x3C–F | AC DC DC 0C |

`0x129` = 1 0010 1001

`0x20` + `0x4` ×1 = `0x24`

*PTE 1 value:*

`0xF4` = 1111 0100

PPN 111, valid 1

*PTE 2 addr:*

111 000 + 101 × 1 = `0x3D`

*PTE 2 value:* `0xDC`

PPN 110; valid 1

M[110 001 (`0x31`)] = `0x0A`

# 2-level example

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x20; translate virtual address 0x129

| physical addresses | bytes |
|---|---|
| 0x00–3 | 00 11 22 33 |
| 0x04–7 | 44 55 66 77 |
| 0x08–B | 88 99 AA BB |
| 0x0C–F | CC DD EE FF |
| 0x10–3 | 1A 2A 3A 4A |
| 0x14–7 | 1B 2B 3B 4B |
| 0x18–B | 1C 2C 3C 4C |
| 0x1C–F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20–3 | 00 91 72 13 |
| 0x24–7 | F4 A5 36 07 |
| 0x28–B | 89 9A AB BC |
| 0x2C–F | CD DE EF F0 |
| 0x30–3 | BA 0A BA 0A |
| 0x34–7 | DB 0B DB 0B |
| 0x38–B | EC 0C EC 0C |
| 0x3C–F | AC DC DC 0C |

0x129 = 1 0010 1001
0x20 + 0x4 ×1 = 0x24
*PTE 1 value:*
0xF4 = 1111 0100
PPN 111, valid 1
*PTE 2 addr:*
111 000 + 101 × 1 = 0x3D
*PTE 2 value:* 0xDC
PPN 110; valid 1
M[110 001 (0x31)] = 0x0A

# 2-level exercise (1)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register 0x08; translate virtual address 0x0FB

| physical addresses | bytes | | | | physical addresses | bytes | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0x00-3 | 00 | 11 | 22 | 33 | 0x20-3 | D0 | D1 | D2 | D3 |
| 0x04-7 | 44 | 55 | 66 | 77 | 0x24-7 | D4 | D5 | D6 | D7 |
| 0x08-B | 88 | 99 | AA | BB | 0x28-B | 89 | 9A | AB | BC |
| 0x0C-F | CC | DD | EE | FF | 0x2C-F | CD | DE | EF | F0 |
| 0x10-3 | 1A | 2A | 3A | 4A | 0x30-3 | BA | 0A | BA | 0A |
| 0x14-7 | 1B | 2B | 3B | 4B | 0x34-7 | DB | 0B | DB | 0B |
| 0x18-B | 1C | 2C | 3C | 4C | 0x38-B | EC | 0C | EC | 0C |
| 0x1C-F | 1C | 2C | 3C | 4C | 0x3C-F | FC | 0C | FC | 0C |

# 2-level exercise (1)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register 0x08; translate virtual address 0x0FB

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x0F3 = 011 111 011
(PTE 1 addr: 0x08 +
PTE size times 011 (3))
*PTE 1:* 0xBB at 0x0B
*PTE 1:* PPN 101 (5) valid 1
*PTE 2:* 0xF0 at 0x2F
*PTE 2:* PPN 111 (7) valid 1
111 011 = 0x3B → 0x0C

24

# 2-level exercise (1)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register 0x08; translate virtual address 0x0FB

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x0F3 = 011 111 011
(PTE 1 addr: 0x08 +
PTE size times 011 (3))
*PTE 1:* 0xBB at 0x0B
*PTE 1:* PPN 101 (5) valid 1
*PTE 2:* 0xF0 at 0x2F
*PTE 2:* PPN 111 (7) valid 1
111 011 = 0x3B → 0x0C

# 2-level exercise (1)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register 0x08; translate virtual address 0x0FB

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x0F3 = 011 111 011
(PTE 1 addr: 0x08 +
PTE size times 011 (3))
*PTE 1:* 0xBB at 0x0B
*PTE 1:* PPN 101 (5) valid 1
*PTE 2:* 0xF0 at 0x2F
*PTE 2:* PPN 111 (7) valid 1
111 011 = 0x3B → 0x0C

# 2-level exercise (1)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register `0x08`; translate virtual address `0x0FB`

| physical addresses | bytes |
| --- | --- |
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
| --- | --- |
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

`0x0F3 = 011 111 011`
(PTE 1 addr: `0x08` +
PTE size times 011 (3))
*PTE 1:* `0xBB` at `0x0B`
*PTE 1:* PPN 101 (5) valid 1
*PTE 2:* `0xF0` at `0x2F`
*PTE 2:* PPN 111 (7) valid 1
`111 011 = 0x3B → 0x0C`

# multi-level page tables

VPN split into pieces for each level of page table

top levels: page table entries point to next page table
     usually using physical page number of next page table

bottom level: page table entry points to destination page

validity checks at each level

# note on VPN splitting

indexes used for lookup parts of the virtual page number
   (there are not multiple VPNs)

## splitting addresses

if:

256-byte ($2^8$ byte) pages
4-byte page table entries
3 levels of page tables
page tables take up 1 page

Q1: page offset size (bits)
A. $<=4$  B. 5–7  C. 8–11  D. 12–15  E. $>15$

Q2: virtual page number size (bits)
A. $<=4$  B. 5–7  C. 8–11  D. 12–15  E. $>15$

Q3: split address 0x1234

# backup slides

# x86-64 page table splitting

48-bit virtual address

12-bit page offset (4KB pages)

36-bit virtual page number, split into four 9-bit parts

page tables at each level: $2^9$ entries, 8 bytes/entry
    page tables take up 4KB (1 page)

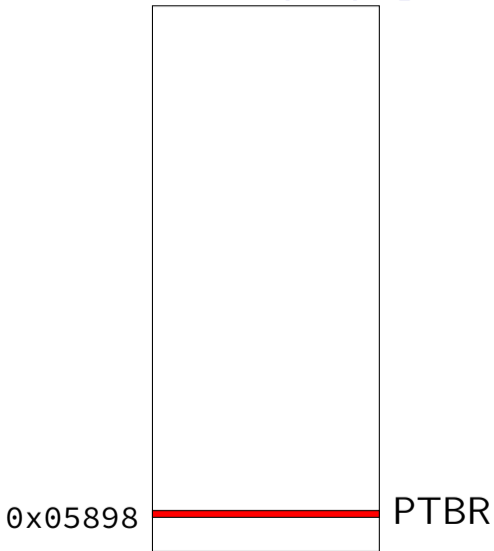# assignment part 2/3

supporting arbitrary numbers of LEVELS, POBITS

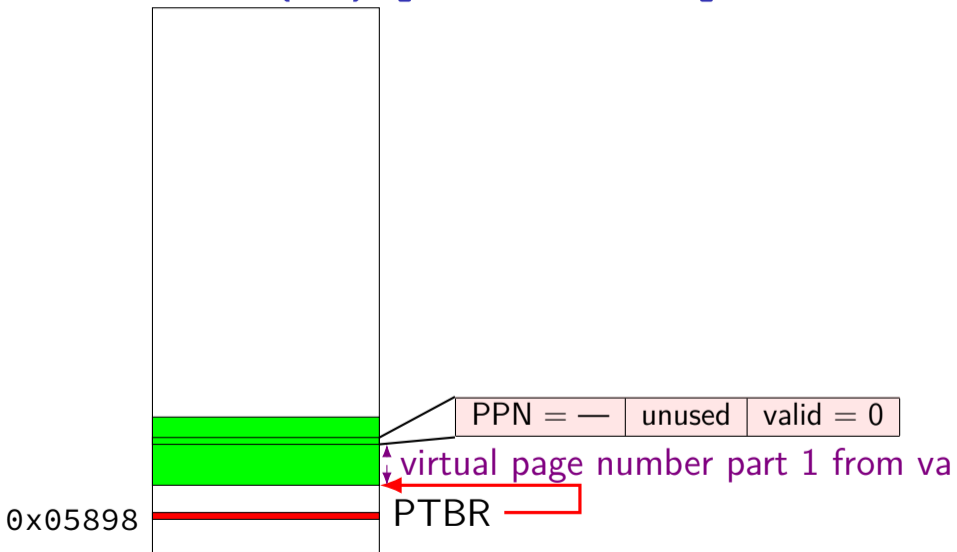code review in lab after reading days
    limited allowed collaboration
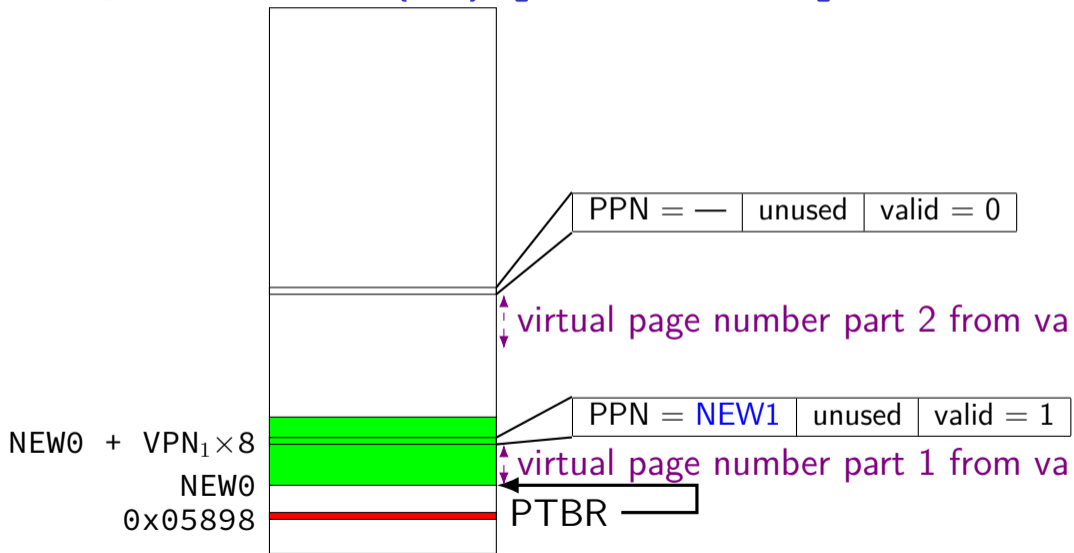
# pa = translate(va) [LEVELS=2]



translate(va)

0x23×page size

physical page 0x23
page offset from va

| PPN = 0x23 | unused | valid = 1 |

0x20× page size + VPN$_2$×8

0x20×page size

physical page 0x20
virtual page number part 2 from va

| PPN = 0x20 | unused | valid = 1 |

0x10000 + VPN$_1$×8

virtual page number part 1 from va

0x10000

0x05898

PTBR

# first page_allocate(va) [LEVELS=2]



0x05898    PTBR

# first page_allocate(va) [LEVELS=2]



| PPN = — | unused | valid = 0 |

virtual page number part 1 from va

0x05898

PTBR

# first page_allocate(va) [LEVELS=2]



PPN = — | unused | valid = 0

virtual page number part 2 from va

PPN = NEW1 | unused | valid = 1

NEW0 + VPN$_1$×8

virtual page number part 1 from va

NEW0

0x05898

PTBR

# first page_allocate(va) [LEVELS=2]



PPN = NEW2 | unused | valid = 1

physical page NEW1

NEW1× page size + VPN$_2$×8

virtual page number part 2 from va

NEW1×page size

PPN = NEW1 | unused | valid = 1

NEW0 + VPN$_1$×8

virtual page number part 1 from va

NEW0

0x05898

PTBR

# first page_allocate(va) [LEVELS=2]



translate(va)

NEW2×page size

physical page NEW2
page offset from va

| PPN = NEW2 | unused | valid = 1 |

NEW1× page size + VPN$_2$×8

NEW1×page size

physical page NEW1
virtual page number part 2 from va

| PPN = NEW1 | unused | valid = 1 |

NEW0 + VPN$_1$×8

NEW0

0x05898

virtual page number part 1 from va

PTBR

# later page allocates?

some of those allocations done earlier
  e.g. ptbr already set

should reuse existing allocation then

# 2-level exercise (2)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register `0x10`; translate virtual address `0x109`

| physical addresses | bytes | | physical addresses | bytes | |
|---|---|---|---|---|---|
| 0x00-3 | 00 11 22 33 | | 0x20-3 | D0 D1 D2 D3 | |
| 0x04-7 | 44 55 66 77 | | 0x24-7 | D4 D5 D6 D7 | |
| 0x08-B | 88 99 AA BB | | 0x28-B | 89 9A AB BC | |
| 0x0C-F | CC DD EE FF | | 0x2C-F | CD DE EF F0 | |
| 0x10-3 | 1A 2A 5A 4A | | 0x30-3 | BA 0A BA 0A | |
| 0x14-7 | 1B 2B 3B 4B | | 0x34-7 | DB 0B DB 0B | |
| 0x18-B | 1C 2C 3C 4C | | 0x38-B | EC 0C EC 0C | |
| 0x1C-F | 1C 2C 3C 4C | | 0x3C-F | FC 0C FC 0C | |

34

# 2-level exercise (2)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register 0x10; translate virtual address 0x109

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 5A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x109 = 100 011 001
(PTE 1 at:
0x10 + PTE size times 4 (100))
*PTE 1:* 0x1B at 0x14
*PTE 1:* PPN 000 (0) valid 1
(second table at:
0 (000) times page size = 0x00)
*PTE 2:* 0x33 at 0x03
*PTE 2:* PPN 001 (1) valid 1
001 001 = 0x09 → 0x99

34

# 2-level exercise (2)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register 0x10; translate virtual address 0x109

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 5A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x109 = 100 011 001
(PTE 1 at:
0x10 + PTE size times 4 (100))
*PTE 1:* 0x1B at 0x14
*PTE 1:* PPN 000 (0) valid 1
(second table at:
0 (000) times page size = 0x00)
*PTE 2:* 0x33 at 0x03
*PTE 2:* PPN 001 (1) valid 1
001 001 = 0x09 → 0x99

34

# 2-level exercise (2)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register `0x10`; translate virtual address `0x109`

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 5A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

`0x109 = 100 011 001`
(PTE 1 at:
`0x10` + PTE size times 4 (100))
*PTE 1:* `0x1B` at `0x14`
*PTE 1:* PPN `000` (0) valid `1`
(second table at:
0 (000) times page size = `0x00`)
*PTE 2:* `0x33` at `0x03`
*PTE 2:* PPN `001` (1) valid `1`
`001 001 = 0x09 → 0x99`

# 2-level exercise (2)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused;

page table base register `0x10`; translate virtual address `0x109`

| physical addresses | bytes |
|---|---|
| 0x00–3 | 00 11 22 33 |
| 0x04–7 | 44 55 66 77 |
| 0x08–B | 88 99 AA BB |
| 0x0C–F | CC DD EE FF |
| 0x10–3 | 1A 2A 5A 4A |
| 0x14–7 | 1B 2B 3B 4B |
| 0x18–B | 1C 2C 3C 4C |
| 0x1C–F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20–3 | D0 D1 D2 D3 |
| 0x24–7 | D4 D5 D6 D7 |
| 0x28–B | 89 9A AB BC |
| 0x2C–F | CD DE EF F0 |
| 0x30–3 | BA 0A BA 0A |
| 0x34–7 | DB 0B DB 0B |
| 0x38–B | EC 0C EC 0C |
| 0x3C–F | FC 0C FC 0C |

`0x109 = 100 011 001`
(PTE 1 at:
`0x10` + PTE size times 4 (100))
*PTE 1:* `0x1B` at `0x14`
*PTE 1:* PPN `000` (0) valid `1`
(second table at:
0 (000) times page size = `0x00`)
*PTE 2:* `0x33` at `0x03`
*PTE 2:* PPN `001` (1) valid `1`
`001 001 = 0x09 → 0x99`

34

# 2-level exercise (3)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x08; translate virtual address 0x00B

| physical addresses | bytes | | physical addresses | bytes |
|---|---|---|---|---|
| 0x00-3 | 00 11 22 33 | | 0x20-3 | D0 D1 D2 D3 |
| 0x04-7 | 44 55 66 77 | | 0x24-7 | D4 D5 D6 D7 |
| 0x08-B | 88 99 AA BB | | 0x28-B | 89 9A AB BC |
| 0x0C-F | CC DD EE FF | | 0x2C-F | CD DE EF F0 |
| 0x10-3 | 1A 2A 3A 4A | | 0x30-3 | BA 0A BA 0A |
| 0x14-7 | 1B 2B 3B 4B | | 0x34-7 | DB 0B DB 0B |
| 0x18-B | 1C 2C 3C 4C | | 0x38-B | EC 0C EC 0C |
| 0x1C-F | 1C 2C 3C 4C | | 0x3C-F | FC 0C FC 0C |

# 2-level exercise (3)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x08; translate virtual address 0x00B

| physical addresses | bytes | | physical addresses | bytes | |
|---|---|---|---|---|---|
| 0x00-3 | 00 11 22 33 | | 0x20-3 | D0 D1 D2 D3 | |
| 0x04-7 | 44 55 66 77 | | 0x24-7 | D4 D5 D6 D7 | 0x0F3 = 000 001 011 |
| 0x08-B | 88 99 AA BB | | 0x28-B | 89 9A AB BC | PTE 1: 0x88 at 0x08 |
| 0x0C-F | CC DD EE FF | | 0x2C-F | CD DE EF F0 | PTE 1: PPN 100 (5) valid 0 |
| 0x10-3 | 1A 2A 3A 4A | | 0x30-3 | BA 0A BA 0A | page fault! |
| 0x14-7 | 1B 2B 3B 4B | | 0x34-7 | DB 0B DB 0B | |
| 0x18-B | 1C 2C 3C 4C | | 0x38-B | EC 0C EC 0C | |
| 0x1C-F | 1C 2C 3C 4C | | 0x3C-F | FC 0C FC 0C | |

# 2-level exercise (3)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register `0x08`; translate virtual address `0x00B`

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x0F3 = 000 001 011
PTE 1: 0x88 at 0x08
PTE 1: PPN 100 (5) valid 0
page fault!

# 2-level exercise (4)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x08; translate virtual address 0x1CB

| physical addresses | bytes | | | | physical addresses | bytes | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0x00-3 | 00 | 11 | 22 | 33 | 0x20-3 | D0 | D1 | D2 | D3 |
| 0x04-7 | 44 | 55 | 66 | 77 | 0x24-7 | D4 | D5 | D6 | D7 |
| 0x08-B | 88 | 99 | AA | BB | 0x28-B | 89 | 9A | AB | BC |
| 0x0C-F | CC | DD | EE | FF | 0x2C-F | CD | DE | EF | F0 |
| 0x10-3 | 1A | 2A | 3A | 4A | 0x30-3 | BA | 0A | BA | 0A |
| 0x14-7 | 1B | 2B | 3B | 4B | 0x34-7 | DB | 0B | DB | 0B |
| 0x18-B | 1C | 2C | 3C | 4C | 0x38-B | EC | 0C | EC | 0C |
| 0x1C-F | 1C | 2C | 3C | 4C | 0x3C-F | FC | 0C | FC | 0C |

# 2-level exercise (4)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x08; translate virtual address 0x1CB

| physical addresses | bytes | | physical addresses | bytes | |
|---|---|---|---|---|---|
| 0x00-3 | 00 11 22 33 | | 0x20-3 | D0 D1 D2 D3 | 0x1CB = 111 001 011 |
| 0x04-7 | 44 55 66 77 | | 0x24-7 | D4 D5 D6 D7 | PTE 1: 0xFF at 0x0F |
| 0x08-B | 88 99 AA BB | | 0x28-B | 89 9A AB BC | PTE 1: PPN 111 (7) valid 1 |
| 0x0C-F | CC DD EE FF | | 0x2C-F | CD DE EF F0 | PTE 2: 0x0C at 0x39 |
| 0x10-3 | 1A 2A 3A 4A | | 0x30-3 | BA 0A BA 0A | PTE 2: PPN 000 (0) valid 0 |
| 0x14-7 | 1B 2B 3B 4B | | 0x34-7 | DB 0B DB 0B | page fault! |
| 0x18-B | 1C 2C 3C 4C | | 0x38-B | EC 0C EC 0C | |
| 0x1C-F | 1C 2C 3C 4C | | 0x3C-F | FC 0C FC 0C | |

# 2-level exercise (4)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x08; translate virtual address 0x1CB

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | 1C 2C 3C 4C |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 D1 D2 D3 |
| 0x24-7 | D4 D5 D6 D7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x1CB = 111 001 011
PTE 1: 0xFF at 0x0F
PTE 1: PPN 111 (7) valid 1
PTE 2: 0x0C at 0x39
PTE 2: PPN 000 (0) valid 0
page fault!

36

# 2-level exercise (4)

9-bit virtual addresses, 6-bit physical; 8 byte pages, 1 byte PTE

page tables 1 page; PTE: 3 bit PPN (MSB), 1 valid bit, 4 unused

page table base register 0x08; translate virtual address 0x1CB

| physical addresses | bytes | physical addresses | bytes |
|---|---|---|---|
| 0x00-3 | 00 11 22 33 | 0x20-3 | D0 D1 D2 D3 |
| 0x04-7 | 44 55 66 77 | 0x24-7 | D4 D5 D6 D7 |
| 0x08-B | 88 99 AA BB | 0x28-B | 89 9A AB BC |
| 0x0C-F | CC DD EE FF | 0x2C-F | CD DE EF F0 |
| 0x10-3 | 1A 2A 3A 4A | 0x30-3 | BA 0A BA 0A |
| 0x14-7 | 1B 2B 3B 4B | 0x34-7 | DB 0B DB 0B |
| 0x18-B | 1C 2C 3C 4C | 0x38-B | EC 0C EC 0C |
| 0x1C-F | 1C 2C 3C 4C | 0x3C-F | FC 0C FC 0C |

0x1CB = 111 001 011
PTE 1: 0xFF at 0x0F
PTE 1: PPN 111 (7) valid 1
PTE 2: 0x0C at 0x39
PTE 2: PPN 000 (0) valid 0
page fault!

# 2-level exercise (5)

10-bit virtual addresses, 6-bit physical; 16 byte pages, 2 byte PTE

page tables 1 page; PTE 1st byte: (MSB) 2-bit PPN, valid bit; rest unused

page table base register `0x10`; translate virtual address `0x376`

| physical addresses | bytes | physical addresses | bytes |
|---|---|---|---|
| 0x00-3 | 00 11 22 33 | 0x20-3 | D0 E1 D2 D3 |
| 0x04-7 | 44 55 66 77 | 0x24-7 | D4 E5 D6 E7 |
| 0x08-B | 88 99 AA BB | 0x28-B | 89 9A AB BC |
| 0x0C-F | CC DD EE FF | 0x2C-F | CD DE EF F0 |
| 0x10-3 | 1A 2A 3A 4A | 0x30-3 | BA 0A BA 0A |
| 0x14-7 | 1B 2B 3B 4B | 0x34-7 | DB 0B DB 0B |
| 0x18-B | 1C 2C 3C 4C | 0x38-B | EC 0C EC 0C |
| 0x1C-F | AC BC DC EC | 0x3C-F | FC 0C FC 0C |

# 2-level exercise (5)

10-bit virtual addresses, 6-bit physical; 16 byte pages, 2 byte PTE

page tables 1 page; PTE 1st byte: (MSB) 2-bit PPN, valid bit; rest unused

page table base register 0x10; translate virtual address 0x376

| physical addresses | bytes |
| --- | --- |
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | AC BC DC EC |

| physical addresses | bytes |
| --- | --- |
| 0x20-3 | D0 E1 D2 D3 |
| 0x24-7 | D4 E5 D6 E7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x376 = 110 111 0110
PTE 1: 0x10 + 6 × 2 = 0x1C:
AC BC
PTE 1: PPN 10 valid 1
PTE 2: 0x20 + 7 × 2 = 0x2E:
EF F0
PTE 2: PPN 11 valid 1
11 0110 = 0x36 → DB

37

# 2-level exercise (5)

10-bit virtual addresses, 6-bit physical; 16 byte pages, 2 byte PTE

page tables 1 page; PTE 1st byte: (MSB) 2-bit PPN, valid bit; rest unused

page table base register `0x10`; translate virtual address `0x376`

| physical addresses | bytes |
|---|---|
| 0x00–3 | 00 11 22 33 |
| 0x04–7 | 44 55 66 77 |
| 0x08–B | 88 99 AA BB |
| 0x0C–F | CC DD EE FF |
| 0x10–3 | 1A 2A 3A 4A |
| 0x14–7 | 1B 2B 3B 4B |
| 0x18–B | 1C 2C 3C 4C |
| 0x1C–F | AC BC DC EC |

| physical addresses | bytes |
|---|---|
| 0x20–3 | D0 E1 D2 D3 |
| 0x24–7 | D4 E5 D6 E7 |
| 0x28–B | 89 9A AB BC |
| 0x2C–F | CD DE EF F0 |
| 0x30–3 | BA 0A BA 0A |
| 0x34–7 | DB 0B DB 0B |
| 0x38–B | EC 0C EC 0C |
| 0x3C–F | FC 0C FC 0C |

$0x376 = 110\ 111\ 0110$
PTE 1: $0x10 + 6 \times 2 = 0x1C$:
AC BC
PTE 1: PPN 10 valid 1
PTE 2: $0x20 + 7 \times 2 = 0x2E$:
EF F0
PTE 2: PPN 11 valid 1
$11\ 0110 = 0x36 \rightarrow DB$

# 2-level exercise (5)

10-bit virtual addresses, 6-bit physical; 16 byte pages, 2 byte PTE

page tables 1 page; PTE 1st byte: (MSB) 2-bit PPN, valid bit; rest unused

page table base register `0x10`; translate virtual address `0x376`

| physical addresses | bytes |
|---|---|
| 0x00–3 | 00 11 22 33 |
| 0x04–7 | 44 55 66 77 |
| 0x08–B | 88 99 AA BB |
| 0x0C–F | CC DD EE FF |
| 0x10–3 | 1A 2A 3A 4A |
| 0x14–7 | 1B 2B 3B 4B |
| 0x18–B | 1C 2C 3C 4C |
| 0x1C–F | AC BC DC EC |

| physical addresses | bytes |
|---|---|
| 0x20–3 | D0 E1 D2 D3 |
| 0x24–7 | D4 E5 D6 E7 |
| 0x28–B | 89 9A AB BC |
| 0x2C–F | CD DE EF F0 |
| 0x30–3 | BA 0A BA 0A |
| 0x34–7 | DB 0B DB 0B |
| 0x38–B | EC 0C EC 0C |
| 0x3C–F | FC 0C FC 0C |

0x376 = 110 111 0110
PTE 1: 0x10 + 6 × 2 = 0x1C:
AC BC
PTE 1: PPN 10 valid 1
PTE 2: 0x20 + 7 × 2 = 0x2E:
EF F0
PTE 2: PPN 11 valid 1
11 0110 = 0x36 → DB

37

# 2-level exercise (5)

10-bit virtual addresses, 6-bit physical; 16 byte pages, 2 byte PTE

page tables 1 page; PTE 1st byte: (MSB) 2-bit PPN, valid bit; rest unused

page table base register 0x10; translate virtual address 0x376

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | AC BC DC EC |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 E1 D2 D3 |
| 0x24-7 | D4 E5 D6 E7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x376 = 110 111 0110
PTE 1: 0x10 + 6 × 2 = 0x1C:
AC BC
PTE 1: PPN 10 valid 1
PTE 2: 0x20 + 7 × 2 = 0x2E:
EF F0
PTE 2: PPN 11 valid 1
11 0110 = 0x36 → DB

37

# 2-level exercise (5)

10-bit virtual addresses, 6-bit physical; 16 byte pages, 2 byte PTE

page tables 1 page; PTE 1st byte: (MSB) 2-bit PPN, valid bit; rest unused

page table base register `0x10`; translate virtual address `0x376`

| physical addresses | bytes |
|---|---|
| 0x00–3 | 00 11 22 33 |
| 0x04–7 | 44 55 66 77 |
| 0x08–B | 88 99 AA BB |
| 0x0C–F | CC DD EE FF |
| 0x10–3 | 1A 2A 3A 4A |
| 0x14–7 | 1B 2B 3B 4B |
| 0x18–B | 1C 2C 3C 4C |
| 0x1C–F | AC BC DC EC |

| physical addresses | bytes |
|---|---|
| 0x20–3 | D0 E1 D2 D3 |
| 0x24–7 | D4 E5 D6 E7 |
| 0x28–B | 89 9A AB BC |
| 0x2C–F | CD DE EF F0 |
| 0x30–3 | BA 0A BA 0A |
| 0x34–7 | DB 0B DB 0B |
| 0x38–B | EC 0C EC 0C |
| 0x3C–F | FC 0C FC 0C |

$0x376 = 110\ 111\ 0110$
PTE 1: $0x10 + 6 \times 2 = 0x1C$:
AC BC
PTE 1: PPN 10 valid 1
PTE 2: $0x20 + 7 \times 2 = 0x2E$:
EF F0
PTE 2: PPN 11 valid 1
11 0110 = 0x36 → DB

# 2-level exercise (5)
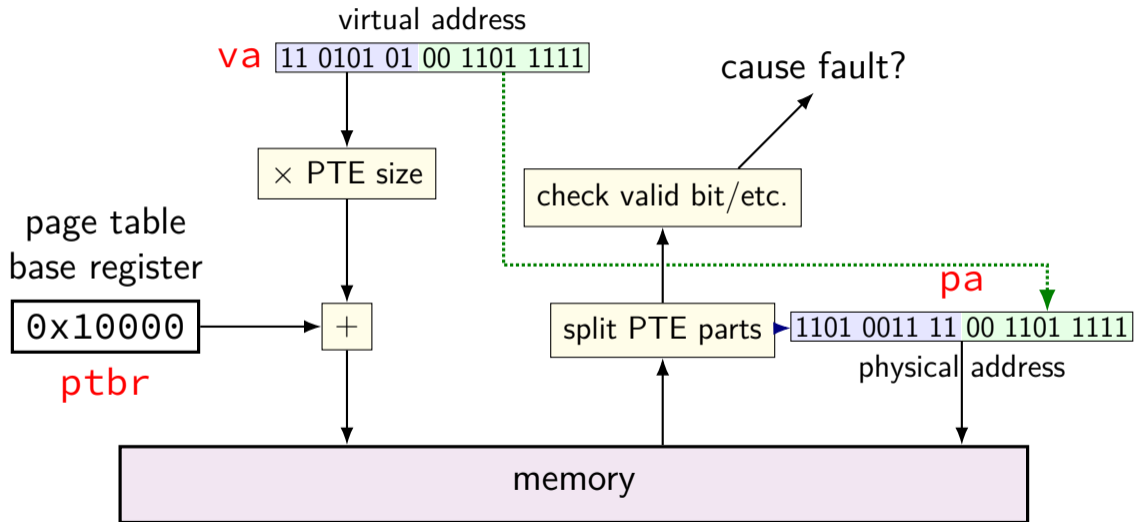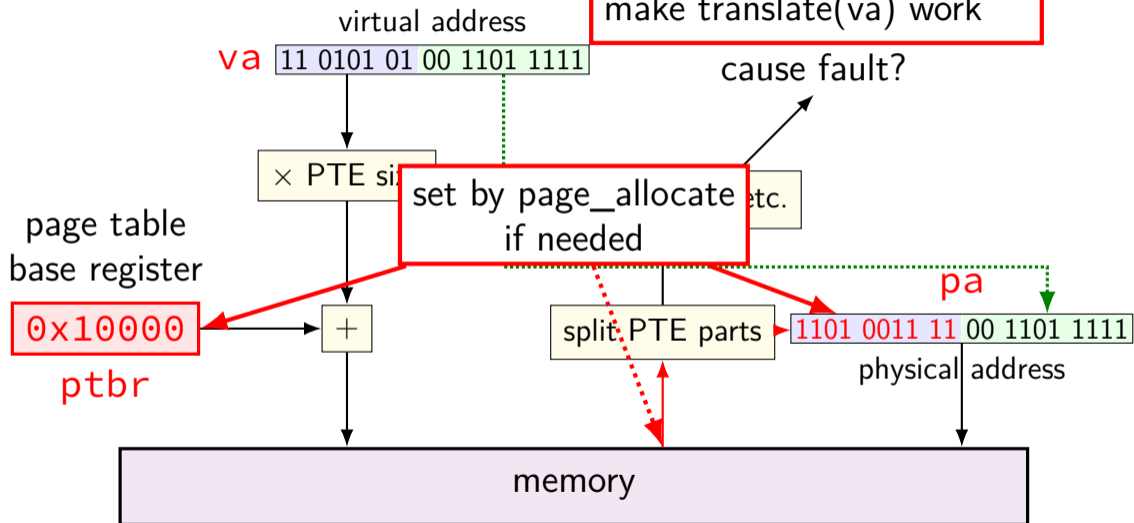
10-bit virtual addresses, 6-bit physical; 16 byte pages, 2 byte PTE

page tables 1 page; PTE 1st byte: (MSB) 2-bit PPN, valid bit; rest unused

page table base register `0x10`; translate virtual address `0x376`

| physical addresses | bytes |
|---|---|
| 0x00–3 | 00 11 22 33 |
| 0x04–7 | 44 55 66 77 |
| 0x08–B | 88 99 AA BB |
| 0x0C–F | CC DD EE FF |
| 0x10–3 | 1A 2A 3A 4A |
| 0x14–7 | 1B 2B 3B 4B |
| 0x18–B | 1C 2C 3C 4C |
| 0x1C–F | AC BC DC EC |

| physical addresses | bytes |
|---|---|
| 0x20–3 | D0 E1 D2 D3 |
| 0x24–7 | D4 E5 D6 E7 |
| 0x28–B | 89 9A AB BC |
| 0x2C–F | CD DE EF F0 |
| 0x30–3 | BA 0A BA 0A |
| 0x34–7 | DB 0B DB 0B |
| 0x38–B | EC 0C EC 0C |
| 0x3C–F | FC 0C FC 0C |

0x376 = 110 111 0110
PTE 1: 0x10 + 6 × 2 = 0x1C:
AC BC
PTE 1: PPN 10 valid 1
PTE 2: 0x20 + 7 × 2 = 0x2E:
EF F0
PTE 2: PPN 11 valid 1
11 0110 = 0x36 → DB

37

# 2-level exercise (5)

10-bit virtual addresses, 6-bit physical; 16 byte pages, 2 byte PTE

page tables 1 page; PTE 1st byte: (MSB) 2-bit PPN, valid bit; rest unused

page table base register 0x10; translate virtual address 0x376

| physical addresses | bytes |
|---|---|
| 0x00-3 | 00 11 22 33 |
| 0x04-7 | 44 55 66 77 |
| 0x08-B | 88 99 AA BB |
| 0x0C-F | CC DD EE FF |
| 0x10-3 | 1A 2A 3A 4A |
| 0x14-7 | 1B 2B 3B 4B |
| 0x18-B | 1C 2C 3C 4C |
| 0x1C-F | AC BC DC EC |

| physical addresses | bytes |
|---|---|
| 0x20-3 | D0 E1 D2 D3 |
| 0x24-7 | D4 E5 D6 E7 |
| 0x28-B | 89 9A AB BC |
| 0x2C-F | CD DE EF F0 |
| 0x30-3 | BA 0A BA 0A |
| 0x34-7 | DB 0B DB 0B |
| 0x38-B | EC 0C EC 0C |
| 0x3C-F | FC 0C FC 0C |

0x376 = 110 111 0110
PTE 1: 0x10 + 6 × 2 = 0x1C:
AC BC
PTE 1: PPN 10 valid 1
PTE 2: 0x20 + 7 × 2 = 0x2E:
EF F0
PTE 2: PPN 11 valid 1
11 0110 = 0x36 → DB

# x86-64 page table entries (1)

| 6 6 6 6 5 5 5 5 5 5 5 5 | M¹ | M-1 | 3 3 3 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 2 1 0 9 8 7 6 5 4 3 2 1 | | | 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 | | | | | | | | | | | | |

| X D | Prot. Key[4] | Ignored | Rsvd. | Address of 4KB page frame | Ign. | G | P A T | D | A | P C D | P W T | U /S | R /W | 1 | PTE: 4KB page |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ignored | | | | | | | | | | | | | | 0 | PTE: not present |

present = valid

R/W = writes allowed?

U/S = user-mode allowed? ("user/supervisor")

XD = execute-disable?

A = accessed? (MMU sets to 1 on page read/write)

D = dirty? (MMU sets to 1 on page write)

# x86-64 page table entries (1)

| 6 6 6 6 5 5 5 5 5 5 5 5 | M¹ | M-1 | 3 3 3 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 2 1 0 9 8 7 6 5 4 3 2 | | | 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 | 9 8 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 | |
| X D | Prot. Key⁴ | Ignored | Rsvd. | Address of 4KB page frame | Ign. | G | P A T | D | A | P C D | P W T | U / S | R / W | 1 | PTE: 4KB page |
| | | | Ignored | | | | | | | | | | | 0 | PTE: not present |

present = valid

R/W = writes allowed?

U/S = user-mode allowed? ("user/supervisor")

XD = execute-disable?

A = accessed? (MMU sets to 1 on page read/write)

D = helps support replacement policies for swapping

# x86-64 page table entries (1)

| 6 6 6 6 5 5 5 5 5 5 5 5 | | | 3 3 3 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 | | | | |
|---|---|---|---|---|---|---|---|
| 3 2 1 0 9 8 7 6 5 4 3 2 1 | M¹ | M-1 | 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 | | | | |
| X D | Prot. Key⁴ | Ignored | Rsvd. | Address of 4KB page frame | Ign. | G A D A C W / S W 1 | PTE: 4KB page |
| | | | Ignored | | | | 0 | PTE: not present |

present = valid

R/W = writes allowed?

U/S = user-mode allowed? ("user/supervisor")

XD = execute-disable?

A = accessed? (MMU sets to 1 on page read/write)

D = dirty? (MMU sets to 1 on page write)

helps support writeback policy for swapping

# x86-64 page table entries (2)



| 63 62 61 60 59 58 57 56 55 54 53 52 | M[1] | M-1 ... | 32 31 30 ... 12 | 11 10 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| X D | Prot. Key[4] | Ignored | Rsvd. | Address of 4KB page frame | Ign. | G | PAT | D | A | PCD | PWT | U/S | R/W | 1 | PTE: 4KB page |
| | | | Ignored | | | | | | | | | | | 0 | PTE: not present |

G = global? (shared between all page tables)

PWT, PCD, PAT = control how caches work when accessing physical page:

    can disable using the cache entirely
    can disable write-back (use write-through instead)
    multicore-related cache settings
    (and some other settings)

# x86-64 page table entries (2)

| 6 6 6 6 5 5 5 5 5 5 5 5 | M[1] | M-1 | 3 3 3 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 2 1 0 9 8 7 6 5 4 3 2 1 | | | 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 | | | | | | | | | | | | | | | |

| X D | Prot. Key[4] | Ignored | Rsvd. | Address of 4KB page frame | Ign. | G | P A T | D | A | P C D | P W T | U / S | R / W | 1 | PTE: 4KB page |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ignored | | | | | | | | | | | | | | 0 | PTE: not present |

G = global? (shared between all page tables)

P | CPU won't evict TLB entries on most page table base registers changes

    can disable using the cache entirely
    can disable write-back (use write-through instead)
    multicore-related cache settings
    (and some other settings)

# pa=translate(va)



virtual address
va `11 0101 01 00 1101 1111`

cause fault?

× PTE size

check valid bit/etc.

page table
base register

pa

`0x10000`

split PTE parts    `1101 0011 11 00 1101 1111`

ptbr

physical address

+

memory

# pa=translate(va)



page_allocate(va) needs to make translate(va) work

virtual address

va `11 0101 01 00 1101 1111`

cause fault?

× PTE si...    etc.

set by page_allocate if needed

page table base register

`0x10000`

ptbr

+

split PTE parts

pa

`1101 0011 11 00 1101 1111`

physical address

memory

# toy program memory

```
11 1111 1111 = 0x3FF ──→  ┌──────────────────┐
                          │      stack        │
11 0000 0000 = 0x300 ──→  ├──────────────────┤
                          │ empty/more heap?  │
10 0000 0000 = 0x200 ──→  ├──────────────────┤
                          │    data/heap      │
01 0000 0000 = 0x100 ──→  ├──────────────────┤
                          │      code         │
00 0000 0000 = 0x000 ──→  └──────────────────┘
```

# toy program memory

```
11 1111 1111 = 0x3FF ──→  ┌─────────────────┐
                          │      stack      │  virtual page# 3
11 0000 0000 = 0x300 ──→  ├─────────────────┤
                          │ empty/more heap?│  virtual page# 2
10 0000 0000 = 0x200 ──→  ├─────────────────┤
                          │    data/heap    │  virtual page# 1
01 0000 0000 = 0x100 ──→  ├─────────────────┤
                          │      code       │  virtual page# 0
00 0000 0000 = 0x000 ──→  └─────────────────┘
```

# toy program memory



```
11 1111 1111 = 0x3FF →
11 0000 0000 = 0x300 →
10 0000 0000 = 0x200 →
01 0000 0000 = 0x100 →
00 0000 0000 = 0x000 →
```

| stack | virtual page# 3 |
| empty/more heap? | virtual page# 2 |
| data/heap | virtual page# 1 |
| code | virtual page# 0 |

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

# toy program memory

| | | |
|---|---|---|
| 11 1111 1111 = 0x3FF → | stack | virtual page# 3 |
| 11 0000 0000 = 0x300 → | empty/more heap? | virtual page# 2 |
| 10 0000 0000 = 0x200 → | data/heap | virtual page# 1 |
| 01 0000 0000 = 0x100 → | code | virtual page# 0 |
| 00 0000 0000 = 0x000 → | | |

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

# toy program memory



```
11 1111 1111 = 0x3FF →
11 0000 0000 = 0x300 →
10 0000 0000 = 0x200 →
01 0000 0000 = 0x100 →
00 0000 0000 = 0x000 →
```

| stack | virtual page# 3 |
| empty/more heap? | virtual page# 2 |
| data/heap | virtual page# 1 |
| code | virtual page# 0 |

rest of address is called page offset

# toy physical memory

real memory
physical addresses

| |
|---|
| 111 0000 0000 to<br>111 1111 1111 |
| |
| |
| |
| |
| |
| 001 0000 0000 to<br>001 1111 1111 |
| 000 0000 0000 to<br>000 1111 1111 |

program memory
virtual addresses

| |
|---|
| 11 0000 0000 to<br>11 1111 1111 |
| 10 0000 0000 to<br>10 1111 1111 |
| 01 0000 0000 to<br>01 1111 1111 |
| 00 0000 0000 to<br>00 1111 1111 |

# toy physical memory

real memory
physical addresses

| | |
|---|---|
| `111 0000 0000 to`<br>`111 1111 1111` | physical page 7 |
| | |
| | |
| | |
| | |
| | |
| `001 0000 0000 to`<br>`001 1111 1111` | physical page 1 |
| `000 0000 0000 to`<br>`000 1111 1111` | physical page 0 |

program memory
virtual addresses

| |
|---|
| `11 0000 0000 to`<br>`11 1111 1111` |
| `10 0000 0000 to`<br>`10 1111 1111` |
| `01 0000 0000 to`<br>`01 1111 1111` |
| `00 0000 0000 to`<br>`00 1111 1111` |

# toy physical memory

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 |

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 |

42

# toy physical memory

real memory
physical addresses

| virtual page # | physical page # |
|---|---|
| 00 | 010 (2) |
| 01 | 111 (7) |
| 10 | *none* |
| 11 | 000 (0) |

program memory
virtual addresses

| |
|---|
| 111 0000 0000 to<br>111 1111 1111 |
| |
| |
| |
| |
| |
| 001 0000 0000 to<br>001 1111 1111 |
| 000 0000 0000 to<br>000 1111 1111 |

| |
|---|
| 11 0000 0000 to<br>11 1111 1111 |
| 10 0000 0000 to<br>10 1111 1111 |
| 01 0000 0000 to<br>01 1111 1111 |
| 00 0000 0000 to<br>00 1111 1111 |

# toy physical memory

**page table!** real memory

physical addresses

| virtual page # | physical page # |
|---|---|
| 00 | 010 (2) |
| 01 | 111 (7) |
| 10 | *none* |
| 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 |

| |
|---|
| 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 memory

# 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) |

"page table entry"

`111` `1101 0010`

trigger exception if 0?

to memory

# t "virtual page number" 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 memory

# 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 memory

# toy page "page offset" lookup

`01` `1101 0010` — address from CPU

virtual
page #   valid?  physical 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 memory

# exit statuses

```
int main() {
    return 0;  /* or exit(0);  */
}
```

# the status

```
#include <sys/wait.h>
...
  waitpid(child_pid, &status, 0);
  if (WIFEXITED(status)) {
    printf("main returned or exit called with %d\n",
           WEXITSTATUS(status));
  } else if (WIFSIGNALED(status)) {
    printf("killed by signal %d\n", WTERMSIG(status));
  } else {
      ...
  }
```

"status code" encodes both return value and if exit was abnormal

W* macros to decode it

# the status

```
#include <sys/wait.h>
...
  waitpid(child_pid, &status, 0);
  if (WIFEXITED(status)) {
    printf("main returned or exit called with %d\n",
           WEXITSTATUS(status));
  } else if (WIFSIGNALED(status)) {
    printf("killed by signal %d\n", WTERMSIG(status));
  } else {
      ...
  }
```

"status code" encodes both return value and if exit was abnormal

W* macros to decode it

# kernel buffering (reads)

program

operating system

keyboard

disk

# kernel buffering (reads)



46

# kernel buffering (reads)



46

# kernel buffering (reads)

# kernel buffering (reads)

# kernel buffering (reads)

# kernel buffering (writes)

| program |
|---|

| operating system |
|---|

| network |   | disk |
|---|---|---|

# kernel buffering (writes)



program

print char
to remote machine

operating system

network

disk

47

# kernel buffering (writes)

# kernel buffering (writes)



47

# kernel buffering (writes)



program

print char
to remote machine

write char
to file

operating system

buffer: output
waiting for network

buffer: data waiting
to be written on disk

(when ready)
send data

(when ready)
write *block* of data from disk

network

disk

# read/write operations

read()/write(): move data into/out of buffer

possibly wait if buffer is empty (read)/full (write)

actual I/O operations — wait for device to be ready
    trigger process to stop waiting if needed

# layering



application

standard library ——— cout/printf — and their own buffers

system calls ——— read/write

kernel's file interface ——— kernel's buffers

device drivers

hardware interfaces

# why the extra layer

better (but more complex to implement) interface:
    read line
    formatted input (scanf, cin into integer, etc.)
    formatted output

less system calls (bigger reads/writes) sometimes faster
    buffering can combine multiple in/out library calls into one system call

more portable interface
    cin, printf, etc. defined by C and C++ standards

# exercise

```
pid_t p = fork();
int pipe_fds[2];
pipe(pipe_fds);
if (p == 0) { /* child */
  close(pipe_fds[0]);
  char c = 'A';
  write(pipe_fds[1], &c, 1);
  exit(0);
} else { /* parent */
  close(pipe_fds[1]);
  char c;
  int count = read(pipe_fds[0], &c, 1);
  printf("read %d bytes\n", count);
}
```

The child is trying to send the character A to the parent, but the above code outputs read 0 bytes instead of read 1 bytes. What happened?

# exercise solution

pipe() is after fork — two pipes, one in child, one in parent

# pipe example (1)

```
int pipe_fd[2];
if (pipe(pipe_fd) < 0)
    handle_error(); /* e.g. out of file descriptors */
int read_fd = pipe_fd[0];
int write_fd = pipe_fd[1];
child_pid = fork();
if (child_pid  == 0) {
    /* in child process, write to pipe */
    close(read_fd);
    write_to_pipe(write_fd); /* function not shown */
    exit(EXIT_SUCCESS);
} else if (child_pid > 0) {
    /* in parent process, read from pipe */
    close(write_fd);
    read_from_pipe(read_fd); /* function not shown */
    waitpid(child_pid, NULL, 0);
    close(read_fd);
} else { /* fork error */ }
```

# pipe example (1)

```
int pipe_fd[2];
if (pipe(pipe_fd) < 0)
    handle_error(); /* e.g. out of file descriptors */
int read_fd = pipe_fd[0];
int write_fd = pipe_fd[1];
child_pid = fork();
if (child_pid  == 0) {
    /* in child process, write to pipe */
    close(read_fd);
    write_to_pipe(write_fd); /* function not shown */
    exit(EXIT_SUCCESS);
} else if (child_pid > 0) {
    /* in parent process, read from pipe */
    close(write_fd);
    read_from_pipe(read_fd); /* function not shown */
    waitpid(child_pid, NULL, 0);
    close(read_fd);
} else { /* fork error */ }
```

# pipe example (1)

```
int pipe_fd[2];
if (pipe(pipe_fd) < 0)
    handle_error(); /* e.g. out of file
int read_fd = pipe_fd[0];
int write_fd = pipe_fd[1];
child_pid = fork();
if (child_pid  == 0) {
    /* in child process, write to pipe */
    close(read_fd);
    write_to_pipe(write_fd); /* function not shown */
    exit(EXIT_SUCCESS);
} else if (child_pid > 0) {
    /* in parent process, read from pipe */
    close(write_fd);
    read_from_pipe(read_fd); /* function not shown */
    waitpid(child_pid, NULL, 0);
    close(read_fd);
} else { /* fork error */ }
```

read() will not indicate
end-of-file if write fd is open
(any copy of it)

# pipe example (1)

```
int pipe_fd[2];
if (pipe(pipe_fd) < 0)
    handle_error(); /* e.g. out of fi│you can run out
int read_fd = pipe_fd[0];
int write_fd = pipe_fd[1];
child_pid = fork();
if (child_pid  == 0) {
    /* in child process, write to pipe */
    close(read_fd);
    write_to_pipe(write_fd); /* function not shown */
    exit(EXIT_SUCCESS);
} else if (child_pid > 0) {
    /* in parent process, read from pipe */
    close(write_fd);
    read_from_pipe(read_fd); /* function not shown */
    waitpid(child_pid, NULL, 0);
    close(read_fd);
} else { /* fork error */ }
```

have habit of closing
to avoid 'leaking' file descriptors
you can run out

# pipe() and blocking

BROKEN example:

```
int pipe_fd[2];
if (pipe(pipe_fd) < 0)
    handle_error();
int read_fd = pipe_fd[0];
int write_fd = pipe_fd[1];
write(write_fd, some_buffer, some_big_size);
read(read_fd, some_buffer, some_big_size);
```

This is likely to not terminate. What's the problem?

# pattern with multiple?

# this class: focus on Unix

Unix-like OSes will be our focus

we have source code

used to from 2150, etc.?

have been around for a while

xv6 imitates Unix

# Unix history

# POSIX: standardized Unix

Portable Operating System Interface (POSIX)
>  "standard for Unix"

current version online:
https://pubs.opengroup.org/onlinepubs/9699919799/

(almost) followed by most current Unix-like OSes

...but OSes add extra features

...and POSIX doesn't specify everything

# what POSIX defines

POSIX specifies the library and shell interface
>    source code compatibility

doesn't care what is/is not a system call…

doesn't specify binary formats…

idea: write applications for POSIX, recompile and run on all implementations
>    this was a very important goal in the 80s/90s
>    at the time, no dominant Unix-like OS (Linux was very immature)

# getpid

```
pid_t my_pid = getpid();
printf("my pid is %ld\n", (long) my_pid);
```

## process ids in ps

```
cr4bd@machine:~$ ps
  PID TTY          TIME CMD
14777 pts/3    00:00:00 bash
14798 pts/3    00:00:00 ps
```

# read/write

```
ssize_t read(int fd, void *buffer, size_t count);
ssize_t write(int fd, void *buffer, size_t count);
```

read/write up to *count* bytes to/from *buffer*

returns number of bytes read/written or -1 on error
    ssize_t is a signed integer type
    error code in errno

read returning 0 means end-of-file (*not an error*)
    can read/write less than requested (end of file, broken I/O device, …)

# read'ing one byte at a time

```
string s;
ssize_t amount_read;
char c;
/* cast to void * not needed in C */
while ((amount_read = read(STDIN_FILENO, (void*) &c, 1)) > 0)
    /* amount_read must be exactly 1 */
    s += c;
}
if (amount_read == -1) {
    /* some error happened */
    perror("read"); /* print out a message about it */
} else if (amount_read == 0) {
    /* reached end of file */
}
```

# write example

```
/* cast to void * optional in C */
write(STDOUT_FILENO, (void *) "Hello, World!\n", 14);
```

# read/write

```
ssize_t read(int fd, void *buffer, size_t count);
ssize_t write(int fd, void *buffer, size_t count);
```

read/write up to *count* bytes to/from *buffer*

returns number of bytes read/written or -1 on error
>    ssize_t is a signed integer type
>    error code in errno

read returning 0 means end-of-file (*not an error*)
>    can read/write less than requested (end of file, broken I/O device, …)

# read'ing a fixed amount

```
ssize_t offset = 0;
const ssize_t amount_to_read = 1024;
char result[amount_to_read];
do {
    /* cast to void * optional in C */
    ssize_t amount_read =
        read(STDIN_FILENO,
             (void *) (result + offset),
             amount_to_read − offset);
    if (amount_read < 0) {
        perror("read"); /* print error message */
        ... /* abort??? */
    } else {
        offset += amount_read;
    }
```

# partial reads

on regular file: read reads what you request

but otherwise: usually gives you what's known to be available
    after waiting for something to be available

# partial reads

on regular file: read reads what you request

but otherwise: usually gives you what's known to be available
    after waiting for something to be available

reading from network — what's been received

reading from keyboard — what's been typed

# write example (with error checking)

```
const char *ptr = "Hello, World!\n";
ssize_t remaining = 14;
while (remaining > 0) {
    /* cast to void * optional in C */
    ssize_t amount_written = write(STDOUT_FILENO,
                                   ptr,
                                   remaining);
    if (amount_written < 0) {
        perror("write"); /* print error message */
        ... /* abort??? */
    } else {
        remaining -= amount_written;
        ptr += amount_written;
    }
}
```

# partial writes

usually only happen on error or interruption
>
> but can request "non-blocking"
>
> (interruption: via *signal*)

*usually*: write waits until it completes

> = until remaining part fits in buffer in kernel
>
> does not mean data was sent on network, shown to user yet, etc.

# aside: environment variables (1)

key=value pairs associated with every process:

```
$ printenv
MODULE_VERSION_STACK=3.2.10
MANPATH=:/opt/puppetlabs/puppet/share/man
XDG_SESSION_ID=754
HOSTNAME=labsrv01
SELINUX_ROLE_REQUESTED=
TERM=screen
SHELL=/bin/bash
HISTSIZE=1000
SSH_CLIENT=128.143.67.91 58432 22
SELINUX_USE_CURRENT_RANGE=
QTDIR=/usr/lib64/qt-3.3
OLDPWD=/zf14/cr4bd
QTINC=/usr/lib64/qt-3.3/include
SSH_TTY=/dev/pts/0
QT_GRAPHICSSYSTEM_CHECKED=1
USER=cr4bd
LS_COLORS=rs=0:di=01;34:ln=01;36:mh=00:pi=40;33:so=01;35:do=01;35:bd=40;33;01:cd=40;33;01:or
MODULE_VERSION=3.2.10
MAIL=/var/spool/mail/cr4bd
PATH=/zf14/cr4bd/.cargo/bin:/zf14/cr4bd/bin:/usr/lib64/qt-3.3/bin:/usr/local/bin:/usr/bin:/u
PWD=/zf14/cr4bd
```

## aside: environment variables (2)

environment variable library functions:
   getenv("KEY") → *value*
   putenv("KEY=*value*") (sets KEY to *value*)
   setenv("KEY", "value") (sets KEY to *value*)

```
int execve(char *path, char **argv, char **envp)
    char *envp[] = { "KEY1=value1", "KEY2=value2", NULL };
    char *argv[] = { "somecommand", "some arg", NULL };
    execve("/path/to/somecommand", argv, envp);
```

normal exec versions — keep same environment variables

# aside: environment variables (3)

interpretation up to programs, but common ones…

PATH=/bin:/usr/bin
    to run a program 'foo', look for an executable in /bin/foo, then
    /usr/bin/foo

HOME=/zf14/cr4bd
    current user's home directory is '/zf14/cr4bd'

TERM=screen-256color
    your output goes to a 'screen-256color'-style terminal

…

## multiple processes?

```
while (...) {
    pid = fork();
    if (pid == 0) {
        exec ...
    } else if (pid > 0) {
        pids.push_back(pid);
    }
}

/* retrieve exit statuses in order */
for (pid_t pid : pids) {
    waitpid(pid, ...);
    ...
}
```

# waiting for all children

```
#include <sys/wait.h>
...
  while (true) {
    pid_t child_pid = waitpid(-1, &status, 0);
    if (child_pid == (pid_t) -1) {
      if (errno == ECHILD) {
        /* no child process to wait for */
        break;
      } else {
        /* some other error */
      }
    }
    /* handle child_pid exiting */
  }
```

# multiple processes?

```
while (...) {
    pid = fork();
    if (pid == 0) {
        exec ...
    } else if (pid > 0) {
        pids.push_back(pid);
    }
}

/* retrieve exit statuses as processes finish */
while ((pid = waitpid(-1, ...)) != -1) {
    handleProcessFinishing(pid);
}
```

# 'waiting' without waiting

```
#include <sys/wait.h>
...
  pid_t return_value = waitpid(child_pid, &status, WNOHANG);
  if (return_value == (pid_t) 0) {
    /* child process not done yet */
  } else if (child_pid == (pid_t) -1) {
    /* error */
  } else {
    /* handle child_pid exiting */
  }
```

# parent and child processes

every process (but process id 1) has a *parent process* (getppid())

this is the process that can wait for it

creates tree of processes (Linux pstree command):

# parent and child questions...

what if parent process exits before child?
    child's parent process becomes process id 1 (typically called *init*)

what if parent process never `waitpid()`s (or equivalent) for child?
    child process stays around as a "zombie"
    can't reuse pid in case parent wants to use `waitpid()`

what if non-parent tries to `waitpid()` for child?
    waitpid fails

# kernel buffering (reads)

program

operating system

keyboard

disk

# kernel buffering (reads)



program

operating system

buffer: keyboard input
waiting for program

① keypress happens, read

keyboard          disk

# kernel buffering (reads)



program

② read char from terminal
③ ...via buffer

operating system

buffer: keyboard input waiting for program

① keypress happens, read

keyboard

disk

80

# kernel buffering (reads)

# kernel buffering (reads)

# kernel buffering (reads)

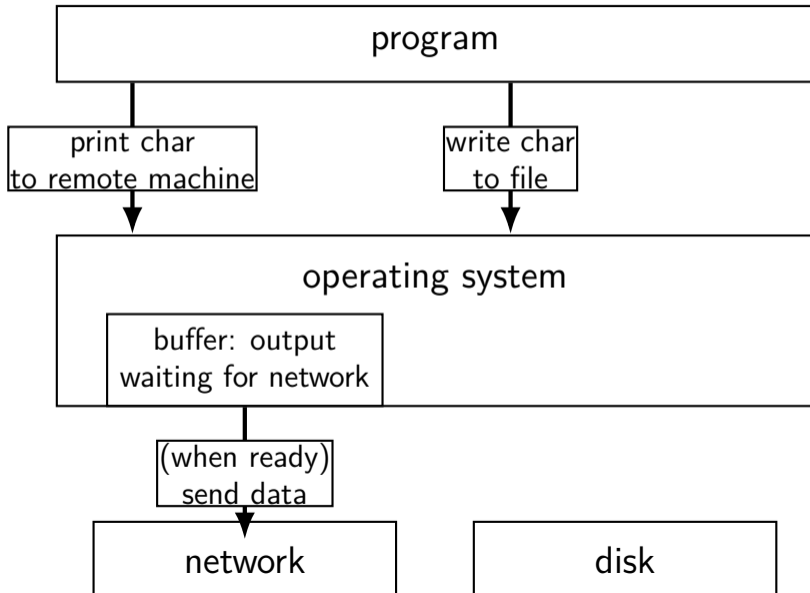# kernel buffering (writes)

program

operating system

network

disk

# kernel buffering (writes)

# kernel buffering (writes)



program

print char
to remote machine

operating system

buffer: output
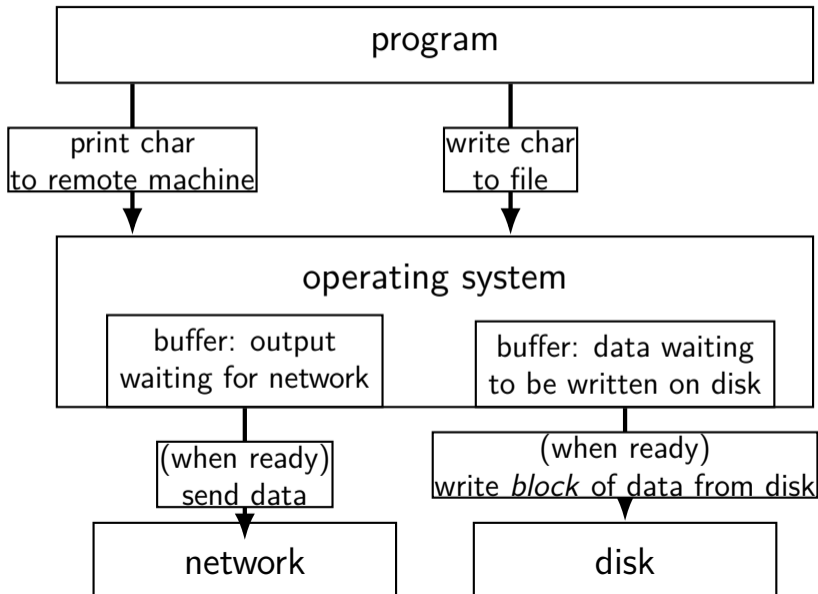waiting for network

(when ready)
send data

network          disk

# kernel buffering (writes)

# kernel buffering (writes)

# read/write operations

read()/write(): move data into/out of buffer

possibly wait if buffer is empty (read)/full (write)

actual I/O operations — wait for device to be ready
    trigger process to stop waiting if needed

# filesystem abstraction

regular files — named collection of bytes
    also: size, modification time, owner, access control info, …

directories — folders containing files and directories
    hierarchical naming: `/net/zf14/cr4bd/fall2018/cs4414`
    *mostly* contains regular files or directories

## open

```
int open(const char *path, int flags);
int open(const char *path, int flags, int mode);
...

int read_fd = open("dir/file1", O_RDONLY);
int write_fd = open("/other/file2",
        O_WRONLY | O_CREAT | O_TRUNC, 0666);
int rdwr_fd = open("file3", O_RDWR);
```

# open

```
int open(const char *path, int flags);
int open(const char *path, int flags, int mode);
```

path = filename

e.g. "/foo/bar/file.txt"
    file.txt in
    directory bar in
    directory foo in
    "the root directory"

e.g. "quux/other.txt
    other.txt in
    directory quux in
    "the current working directory" (set with chdir())

# open: file descriptors

```
int open(const char *path, int flags);
int open(const char *path, int flags, int mode);
```

return value = file descriptor (or -1 on error)

index into table of *open file descriptions* for each process

used by system calls that deal with open files

# POSIX: everything is a file

the file: one interface for
  devices (terminals, printers, …)
  regular files on disk
  networking (sockets)
  local interprocess communication (pipes, sockets)


basic operations: open(), read(), write(), close()

# exercise

```
int pipe_fds[2]; pipe(pipe_fds);
pid_t p = fork();
if (p == 0) {
  close(pipe_fds[0]);
  for (int i = 0; i < 10; ++i) {
    char c = '0' + i;
    write(pipe_fds[1], &c, 1);
  }
  exit(0);
}
close(pipe_fds[1]);
char buffer[10];
ssize_t count = read(pipe_fds[0], buffer, 10);
for (int i = 0; i < count; ++i) {
  printf("%c", buffer[i]);
}
```

Which of these are possible outputs (if pipe, read, write, fork don't fail)?

A. 0123456789   B. 0        C. (nothing)
D. A and B       E. A and C  F. A, B, and C

# exercise

```
int pipe_fds[2]; pipe(pipe_fds);
pid_t p = fork();
if (p == 0) {
  close(pipe_fds[0]);
  for (int i = 0; i < 10; ++i) {
    char c = '0' + i;
    write(pipe_fds[1], &c, 1);
  }
  exit(0);
}
close(pipe_fds[1]);
char buffer[10];
ssize_t count = read(pipe_fds[0], buffer, 10);
for (int i = 0; i < count; ++i) {
  printf("%c", buffer[i]);
}
```

Which of these are possible outputs (if pipe, read, write, fork don't fail)?
A. 0123456789   B. 0         C. (nothing)
D. A and B         E. A and C   F. A, B, and C

# empirical evidence

```
    8 0
  374 01
  210 012
   30 0123
   12 01234
    3 012345
    1 0123456
    2 01234567
    1 012345678
  359 0123456789
```

# partial reads

read returning 0 always means end-of-file
    by default, read always waits *if no input available yet*
    but can set read to return *error* instead of waiting

read can return less than requested if not available
    e.g. child hasn't gotten far enough

# pipe: closing?

if all write ends of pipe are closed
    can get end-of-file (read() returning 0) on read end
    exit()ing closes them

$\rightarrow$ close write end when not using

generally: limited number of file descriptors per process

$\rightarrow$ good habit to close file descriptors not being used

(but probably didn't matter for read end of pipes in example)

# swapping almost mmap

access mapped file for first time, read from disk
    (like swapping when memory was swapped out)

write "mapped" memory, write to disk eventually
    (like writeback policy in swapping)
    use "dirty" bit


extra detail: other processes should see changes
    all accesses to file use same physical memory