# Operating Systems Practice

Virtual Memory – Copy-on-Write

Eunji Lee

(ejlee@ssu.ac.kr)

# Reference

- "xv6: a simple, Unix-like teaching operating system," Chapter 3. Page tables

# Paging Hardware

- RISC-V uses 39-bit virtual address

- A page table is stored in physical memory as a three-level tree

- Each pte has 44-bit PPN and 10-bit flag

**root**
- 4KB (1page)
- 512 entries

**register**: physicall address of the root page-table page

mmu.h

```
134 // Page table/directory entry flags.
135 #define PTE_P        0x001  // Present
136 #define PTE_W        0x002  // Writeable
137 #define PTE_U        0x004  // User
138 #define PTE_PWT      0x008  // Write-Through
139 #define PTE_PCD      0x010  // Cache-Disable
140 #define PTE_A        0x020  // Accessed
141 #define PTE_D        0x040  // Dirty
142 #define PTE_PS       0x080  // Page Size
143 #define PTE_MBZ      0x180  // Bits must be zero
144
```
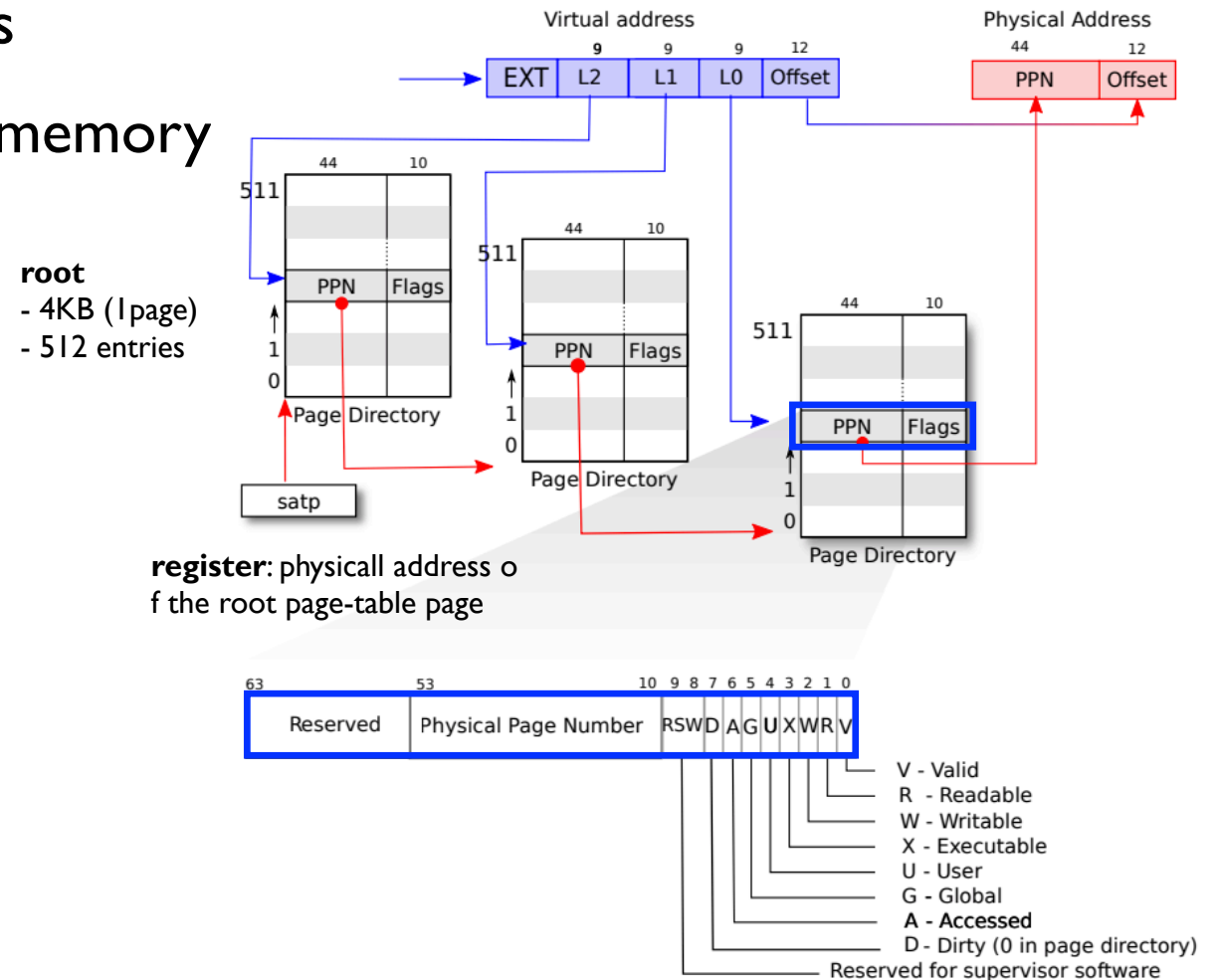
Figure 3.2: RISC-V page table hardware.

# Kernel Address Space

- Direct mapped
- Only kernel stack pages are not direct-mapped
- guard page
  - Prevent problem caused by kernel stack overflow
  - guard page's PTE is invalid



free pages

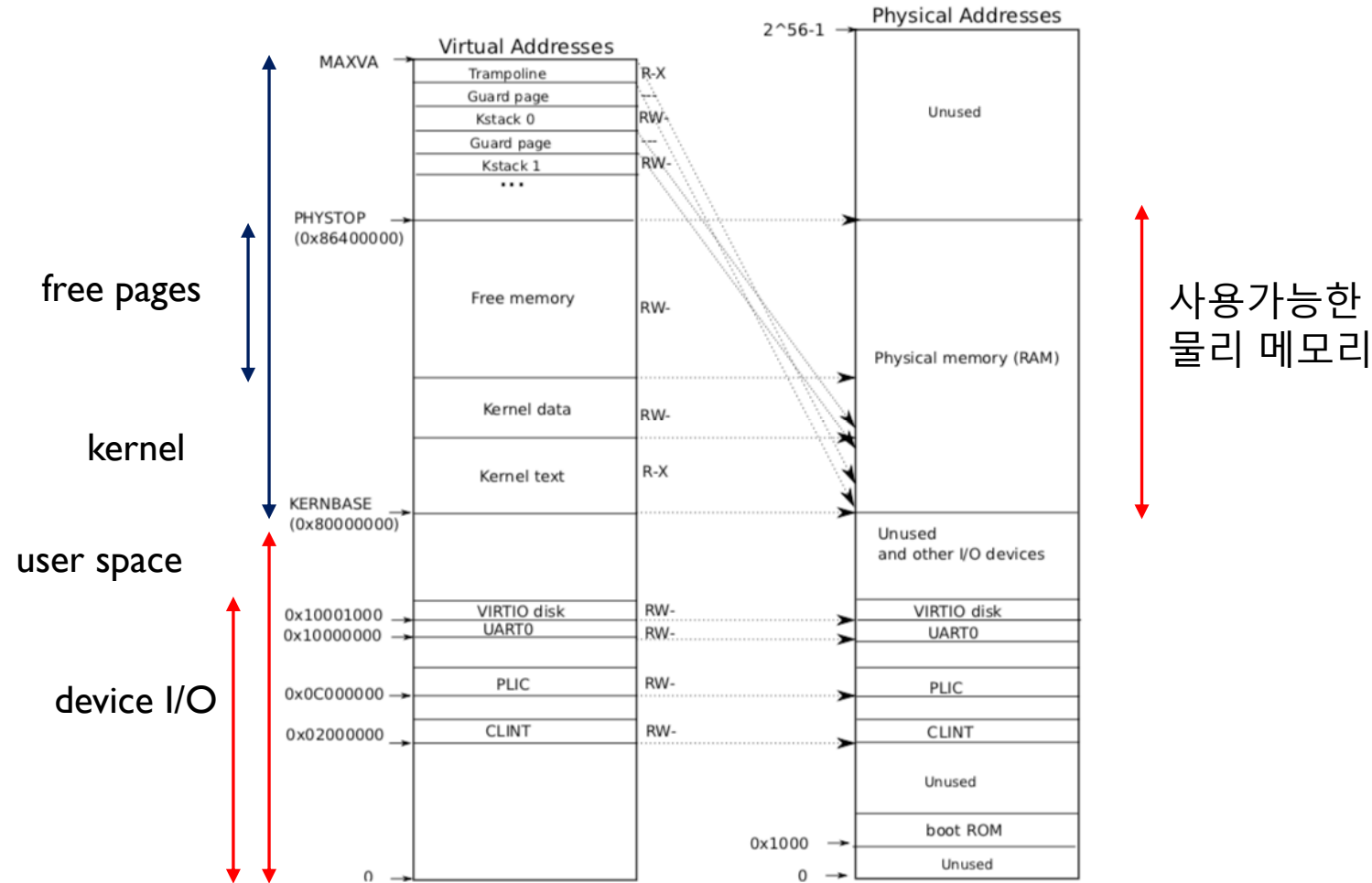kernel

user space

device I/O

사용가능한
물리 메모리

Figure 3.3: On the left, xv6's kernel address space. RWX refer to PTE read, write, and execute permissions. On the right, the RISC-V physical address space that xv6 expects to see.

# Creating an address space

- vm.c : code for manipulate address spaces and page table

proc.c

```
int
fork(void)
{
  int i, pid;
  struct proc *np;
  struct proc *curproc = myproc();

  // Allocate process.
  if((np = allocproc()) == 0){     struct proc 할당
    return -1;
  }

                                    parent 의 address space 를 복사
  // Copy process state from proc.
  if((np->pgdir = copyuvm(curproc->pgdir, curproc->sz)) == 0){
    kfree(np->kstack);
    np->kstack = 0;
    np->state = UNUSED;
    return -1;
  }
```

# Creating an address space

vm.c

```
pde_t*
copyuvm(pde_t *pgdir, uint sz)
{
  pde_t *d;
  pte_t *pte;
  uint pa, i, flags;
  char *mem;

  if((d = setupkvm()) == 0)          kernel stack 할당
    return 0;
  for(i = 0; i < sz; i += PGSIZE){   // parent address space size = sz
    if((pte = walkpgdir(pgdir, (void *) i, 0)) == 0)   // walkpgdir: virtual address 에 대한 pte 찾음
      panic("copyuvm: pte should exist");
    if(!(*pte & PTE_P))
      panic("copyuvm: page not present");
    pa = PTE_ADDR(*pte);            page table entry에서 user virtual address에 대한 physical address 와 flag 를
    flags = PTE_FLAGS(*pte);        읽어냄.
    if((mem = kalloc()) == 0)
      goto bad;
    memmove(mem, (char*)P2V(pa), PGSIZE);
    if(mappages(d, (void*)i, PGSIZE, V2P(mem), flags) < 0)
      goto bad;
  }
  return d;

bad:
  freevm(d);
  return 0;
}
```

pde_t *pgdir은 parent process의 page table
pde_t *d는 새로운 child process를 위한 page table

# Creating an address space

vm.c

```
pde_t*
copyuvm(pde_t *pgdir, uint sz)
{
  pde_t *d;
  pte_t *pte;
  uint pa, i, flags;
  char *mem;

  if((d = setupkvm()) == 0)
    return 0;
  for(i = 0; i < sz; i += PGSIZE){
    if((pte = walkpgdir(pgdir, (void *) i, 0)) == 0)
      panic("copyuvm: pte should exist");
    if(!(*pte & PTE_P))
      panic("copyuvm: page not present");
    pa = PTE_ADDR(*pte);
    flags = PTE_FLAGS(*pte);
    if((mem = kalloc()) == 0)          physical page 하나를 얻어옴.
      goto bad;
    memmove(mem, (char*)P2V(pa), PGSIZE);     parant 의 page 를 child 가 새로 할당받은 page 로 복사
    if(mappages(d, (void*)i, PGSIZE, V2P(mem), flags) < 0)     새로운 child 의 page 를 page table 에 저장
      goto bad;
  }
  return d;
                    d: child 의 page table
bad:                i: mapping 이 필요한 virtual address
  freevm(d);        V2P(mem) : physical address
  return 0;         Flags
}
```
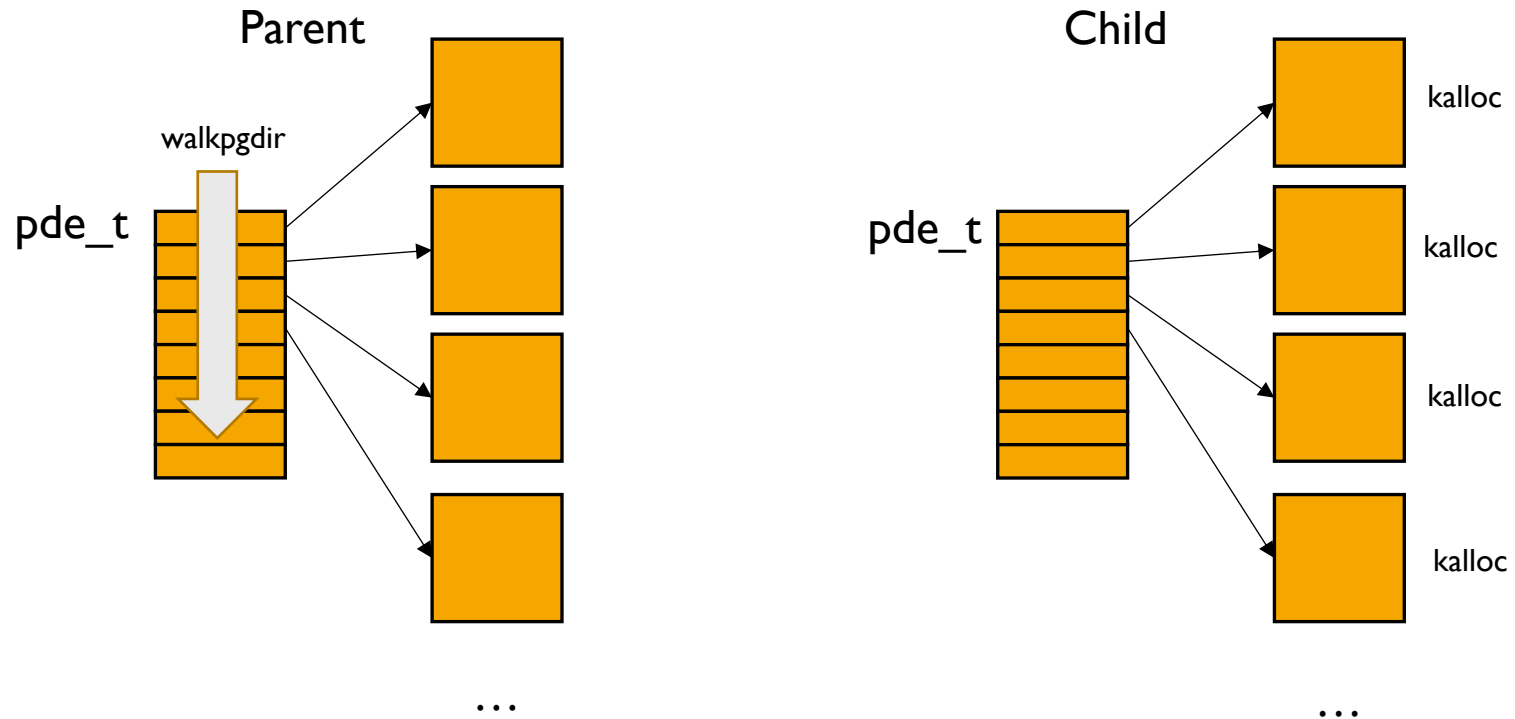
pde_t *pgdir은 parent process의 page table
pde_t *d는 새로운 child process를 위한 page table

# Creating an address space

- Duplication (copyuvm() in vm.c)

Parent

walkpgdir

pde_t

...

Child

pde_t

kalloc

kalloc

kalloc

kalloc

...

# Creating an address space

- Copy-on-write (copyuvm() in vm.c)



read-only

Parent          Child

walkpgdir

pde_t                    pde_t

...

read-only

Parent          Child

walkpgdir

pde_t                    pde_t

write

1) page fault
2) page handler invoked
3) allocate new page
4) copy original page into new page
5) update page table entry
6) TLB flush

# Modification – Reference Counter

- Add reference counter and associated functions
  - kalloc.c

```
15 #ifdef COW
16 struct {
17   struct spinlock lock;
18   int numfreepages;
19   uint ref[PHYSTOP >> PGSHIFT];
20 } pmem;
21 #endif
```

```
38 void
39 kinit1(void *vstart, void *vend)
40 {
41 #ifdef COW
42   initlock(&pmem.lock, "pmemlock");
43 #endif
44   initlock(&kmem.lock, "kmem");
45   kmem.use_lock = 0;
46   freerange(vstart, vend);
47 }
48
49 void
50 kinit2(void *vstart, void *vend)
51 {
52 #ifdef COW
53   memset(&pmem.ref, 1, sizeof(uint) * (PHYSTOP >> PGSHIFT));
54   acquire(&pmem.lock);
55   pmem.numfreepages = 0;
56   release(&pmem.lock);
57 #endif
58   freerange(vstart, vend);
59   kmem.use_lock = 1;
60 }
61
```

```
98 //PAGEBREAK: 21
99 // Free the page of physical memory pointed at by v,
100 // which normally should have been returned by a
101 // call to kalloc().  (The exception is when
102 // initializing the allocator; see kinit above.)
103 void
104 kfree(char *v)
105 {
106   struct run *r;
107
108   if((uint)v % PGSIZE || v < end || V2P(v) >= PHYSTOP)
109     panic("kfree");
110
111   // Fill with junk to catch dangling refs.
112   memset(v, 1, PGSIZE);
113
114   if(kmem.use_lock)
115     acquire(&kmem.lock);
116   r = (struct run*)v;
117   r->next = kmem.freelist;
118   kmem.freelist = r;
119   if(kmem.use_lock)
120     release(&kmem.lock);
121 #ifdef COW
122   if(pmem.use_lock)
123     acquire(&pmem.lock);
124
125   pmem.numfreepages++;
126
127   if(pmem.use_lock)
128     release(&pmem.lock);
129 #endif
130 }
```

```
135 char*
136 kalloc(void)
137 {
138   struct run *r;
139
140   if(kmem.use_lock)
141     acquire(&kmem.lock);
142   r = kmem.freelist;
143   if(r){
144     kmem.freelist = r->next;
145
146 #ifdef COW
147     if(pmem.use_lock)
148       acquire(&pmem.lock);
149
150     pmem.numfreepages++;
151
152     if(pmem.use_lock)
153       release(&pmem.lock);
154 #endif
155   }
156   if(kmem.use_lock)
157     release(&kmem.lock);
158
159   return (char*)r;
160 }
161
```

# Modification – Reference Counter

- Add reference counter and associated functions
  - kalloc.c

```c
70 #ifdef COW
71 // reference counter APIs
72 int freemem()
73 {
74     return 0;
75 }
76
77 uint
78 get_ref(uint pa)
79 {
80     return 0;
81 }
82
83 void
84 inc_ref(void)
85 {
86     return;
87 }
88
89 void
90 dec_ref(uint pa)
91 {
92     return;
93 }
94 #endif
```

# Modification – Reference Counter

- Add reference counter and associated functions
  - `mmu.h`

```
87 #define PGSHIFT          12        // log2(PGSIZE)
88 #define PTXSHIFT         12        // offset of PTX in a linear address
89 #define PDXSHIFT         22        // offset of PDX in a linear address
```

  - `memlayout.h`

```
1 // Memory layout
2
3 #define EXTMEM   0x100000          // Start of extended memory
4 #define PHYSTOP  0xE000000         // Top physical memory
5 #define DEVSPACE 0xFE000000        // Other devices are at high addresses
6
```

# Modification – Reference Counter

- Add reference counter and associated functions in kalloc.c
  - `ref[PHYSTOP >> PGSHIFT]`: reference counter for physical memory pages

  - `uint get_ref(uint pa)`: read the reference count of pa
  - `void inc_ref(uint pa)`: increase the reference count of pa
  - `void dec_ref(uint pa)`: decrease the reference count of pa
- Increase reference counter
  - When allocating the physical page
  - When referencing the physical page
- Decrease reference counter
  - When de-allocating the physical page
  - When de-referencing the physical page
- Increase / decrease reference counter appropriately
  - kalloc / kfree

# Modification – Sharing pages

vm.c

- copyuvm() in vm.c
  - DO NOT allocate a new page for child's address space
  - Install parent's address spaces to child's page table
  - Set the page NOT writable (flags) for both child and parent
  - Increase reference counter for shared page
  - TLB flush: `lcr3(V2P(pgdir))`

```c
313 // Given a parent process's page table, create a copy
314 // of it for a child.
315 pde_t*
316 copyuvm(pde_t *pgdir, uint sz)
317 {
318   pde_t *d;
319   pte_t *pte;
320   uint pa, i, flags;
321   char *mem;
322
323   if((d = setupkvm()) == 0)
324     return 0;
325   for(i = 0; i < sz; i += PGSIZE){
326     if((pte = walkpgdir(pgdir, (void *) i, 0)) == 0)
327       panic("copyuvm: pte should exist");
328     if(!(*pte & PTE_P))
329       panic("copyuvm: page not present");
330     pa = PTE_ADDR(*pte);
331     flags = PTE_FLAGS(*pte);
332     if((mem = kalloc()) == 0)
333       goto bad;
334     memmove(mem, (char*)P2V(pa), PGSIZE);
335     if(mappages(d, (void*)i, PGSIZE, V2P(mem), flags) < 0) {
336       kfree(mem);
337       goto bad;
338     }
339   }
340   return d;
341
342 bad:
343   freevm(d);
344   return 0;
345 }
```

# Modification – Page fault handler

- Add page fault handler in trap.c

- T_PGFLT occurs when memory access is invalid

trap.c

```
35 //PAGEBREAK: 41
36 void
37 trap(struct trapframe *tf)
38 {
39   if(tf->trapno == T_SYSCALL){
40     if(myproc()->killed)
41       exit();
42     myproc()->tf = tf;
43     syscall();
44     if(myproc()->killed)
45       exit();
46     return;
47   }
48
49   switch(tf->trapno){
50 #ifdef COW
51   case T_PGFLT:
52     page_fault();
53     break;
54 #endif
55   case T_IRQ0 + IRQ_TIMER:
56     if(cpuid() == 0){
57       acquire(&tickslock);
58       ticks++;
59       wakeup(&ticks);
60       release(&tickslock);
61     }
```

# Modification – Page fault handler

- void page_fault(void)
  - `rcr2()`: return virtual address incurring page fault
  - check if the virtual address is valid
  - locate page table entry for the virtual address
    - `pte_t`
    - `walkpgdir`()
  - check the reference count of the physical address correspondin the virtual address
    - `get_refcounter()`
  - Perform copy-on-write
    - if reference counter > 2: allocate new page & copy
      - `kalloc`(), `memmove`()
      - decrease reference count for original page
    - Make pages writeable (update page table entry)
  - `lcr3(V2P(pgdir))`: TLB flush

vm.c

```
387
388 #ifdef COW
389 void
390 page_fault(void)
391 {
392     uint va = rcr2();
393
394     // fill this part
395
396
397     return 0;
398 }
399 #endif
400
401 //PAGEBREAK!
402 // Blank page.
403 //PAGEBREAK!
404 // Blank page.
405 //PAGEBREAK!
406 // Blank page.
407
```

# COW Test

- Three tests

```c
1 #include "types.h"
2 #include "stat.h"
3 #include "user.h"
4
5 int
6 main(int argc, char **argv)
7 {
8   int before, after;
9   int pid;
10
11  printf(1, "TEST1: ");
12
13  before = freemem();
14
15  pid = fork();
16  if(pid == 0){
17    after = freemem();
18    if(before - after == 68)
19      printf(1, "OK\n");
20    else
21      printf(1, "WRONG\n");
22    exit();
23  }
24  else{
25    wait();
26  }
27
28  exit();
29 }
~
~
~
cowtest1.c                  1,1           All
```

```c
1 #include "types.h"
2 #include "stat.h"
3 #include "user.h"
4
5 int
6 main(int argc, char **argv)
7 {
8   int before, after;
9   int pid;
10
11  printf(1, "TEST2: ");
12
13  before = freemem();
14
15  pid = fork();
16  if(pid == 0){
17    exit();
18  }
19  else{
20    wait();
21  }
22
23  after = freemem();
24  if(before == after)
25    printf(1, "OK\n");
26  else
27    printf(1, "WRONG\n");
28
29  exit();
30 }
~
cowtest2.c                  1,1           All
```

```c
1 #include "types.h"
2 #include "stat.h"
3 #include "user.h"
4
5 int global = 3;
6
7 int
8 main(int argc, char **argv)
9 {
10   int before, after;
11   int pid;
12
13   printf(1, "TEST3: ");
14
15   pid = fork();
16   if(pid == 0){
17     before = freemem();
18     global = 4;
19     after = freemem();
20     if(before - after == 1)
21       printf(1, "OK\n");
22     else
23       printf(1, "WRONG\n");
24     exit();
25   }
26   else{
27     wait();
28   }
29
30   exit();
31 }
~
cowtest3.c                  1,1           All
```
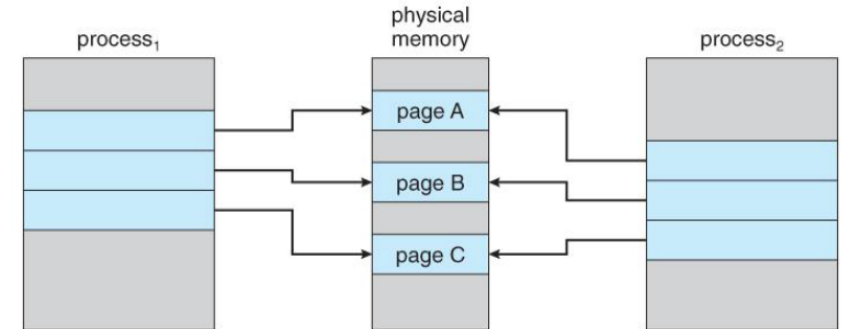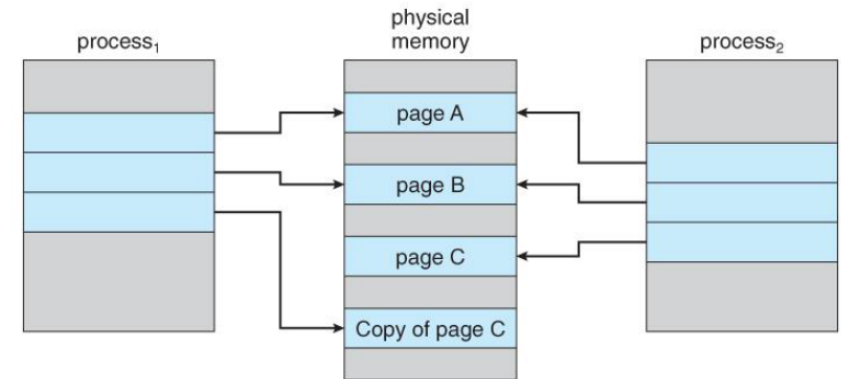
# COW Test

- All should be "OK"

# Copy-on-Write

- When a process forks
  - Create shared mappings to the same page frames in physical page
  - Shared pages are protected as read-only

- When data is written to shared pages
  - Protection fault is generated
  - OS allocates new space in physical memory and directs the write to it

- Reference counter for physical pages is needed



Before process 1 modifies page C



After process 1 modifies page C