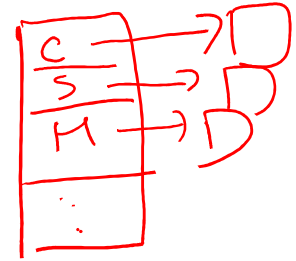CS 347M (Operating Systems Minor)

Spring 2022

# Lecture 11:
# Memory management in xv6

Mythili Vutukuru

CSE, IIT Bombay
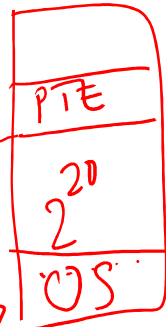
# Virtual memory in xv6

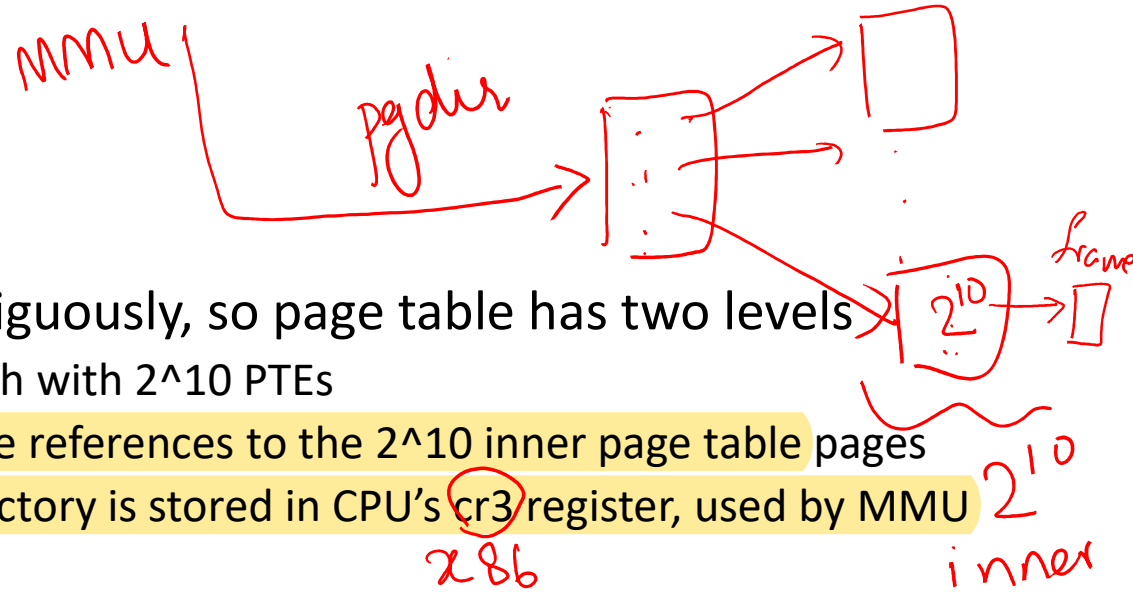- 32-bit OS, so 2^32=4GB virtual address space for every process

- 4KB pages, so 32 bit VA = 20 bit page number + 12 bit offset

- No demand paging, all pages assigned frames always

- Each PTE has 20 bit physical frame number, and some flags
  - PTE_P indicates if page is present (if not set, access will cause page fault)
  - PTE_W indicates if writeable (if not set, only reading is permitted)
  - PTE_U indicates if user page (if not set, only kernel can access the page)

- Address translation: use page number (top 20 bits of virtual address) to index into page table, find physical frame number, add 12-bit offset

# Two level page table

- 2^20 PTEs cannot be stored contiguously, so page table has two levels
  - 2^10 "inner" page table pages, each with 2^10 PTEs
  - Outer page directory stores PTE-like references to the 2^10 inner page table pages
  - Physical address of outer page directory is stored in CPU's cr3 register, used by MMU during address translation
- 32 bit virtual address = 10 bits index into page directory, next 10 bits index into inner page table, last 12 bits are offset within page
  - PFN from PTE + offset = physical address

```
0773 // A virtual address 'la' has a three-part structure as follows:
0774 //
0775 // +--------10--------+-------10--------+---------12-----------+
0776 // | Page Directory   |   Page Table    | Offset within Page   |
0777 // |     Index        |     Index       |                      |
0778 // +------------------+-----------------+----------------------+
0779 //  \--- PDX(va) ---/ \--- PTX(va) ---/
0780
```
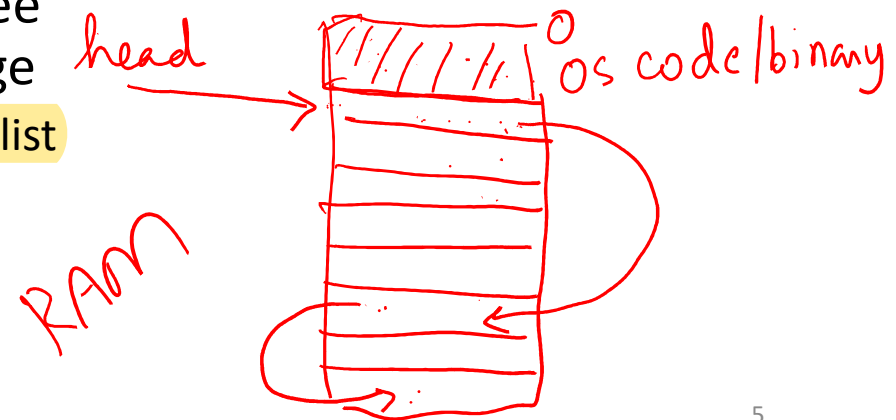
# Process virtual address space in xv6

- Memory image of a process starting at address 0
  - Contains code/data from executable, one page stack (with extra, guard page), expandable heap

- Kernel code/data is mapped beginning at address KERNBASE (2GB)
  - Contains kernel code/data, and free pages maintained by kernel

- Page table of a process contains two sets of PTEs
  - User entries map low virtual addresses to physical memory used by the process for its code/data/stack/heap
  - Kernel entries map high virtual addresses (KERNBASE to KERNBASE+PHYSTOP) to physical memory containing OS code/data and free memory (0 to PHYSTOP)

*(handwritten annotations: "O", "exec", "Stack", "heap", "2GB", "PT", "user", "OS", "Virtual", "0", "2GB", "4GB", "RAM")*
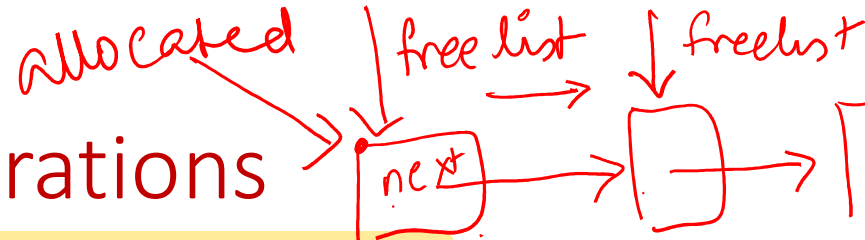
# Maintaining free memory

- After boot up, RAM contains OS code/data and free pages

- OS collects all free pages into a free list, so that it can be assigned to user processes
  - Used for user memory (code/data/stack/heap) and page tables of user processes

- Free list is a linked list, pointer to next free page embedded within previous free page
  - Kernel maintains pointer to first page in the list

```
3115 struct run {
3116   struct run *next;
3117 };
3118
3119 struct {
3120   struct spinlock lock;
3121   int use_lock;
3122   struct run *freelist;
3123 } kmem;
```

head

RAM

OS code/binary
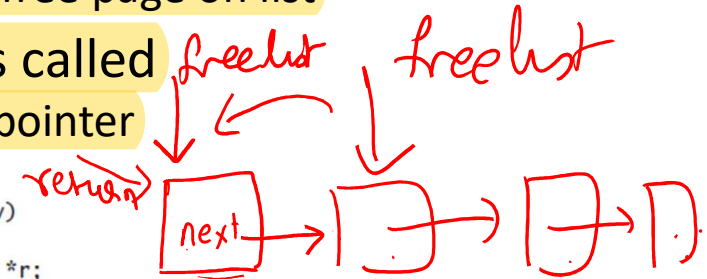
5

# alloc and free operations

- Anyone who needs a free page calls kalloc()
  - Sets free list pointer to next page and returns first free page on list
- When memory needs to be freed up, kfree() is called
  - Add free page to head of free list, update free list pointer

```
3186 char*
3187 kalloc(void)
3188 {
3189   struct run *r;
3190
3191   if(kmem.use_lock)
3192     acquire(&kmem.lock);
3193   r = kmem.freelist;
3194   if(r)
3195     kmem.freelist = r->next;
3196   if(kmem.use_lock)
3197     release(&kmem.lock);
3198   return (char*)r;
3199 }
```

```
3163 void
3164 kfree(char *v)
3165 {
3166   struct run *r;
3167
3168   if((uint)v % PGSIZE || v < end || V2P(v) >= PHYSTOP)
3169     panic("kfree");
3170
3171   // Fill with junk to catch dangling refs.
3172   memset(v, 1, PGSIZE);
3173
3174   if(kmem.use_lock)
3175     acquire(&kmem.lock);
3176   r = (struct run*)v;
3177   r->next = kmem.freelist;
3178   kmem.freelist = r;
3179   if(kmem.use_lock)
3180     release(&kmem.lock);
3181 }
```

# Memory management of user processes (1)

- User process needs memory pages to build its address space
  - User part of memory image (user code/data/stack/heap)
  - Page table (mappings to user memory image, as well as to kernel code/data)
- Note: only one copy of OS in memory, copied during boot time
- New virtual address space for a process is created during:
  - init process creation
  - fork system call
  - exec system call
- Existing virtual address space modified in sbrk system call
  - Invoked by malloc to expand heap

# Memory management of user processes (2)
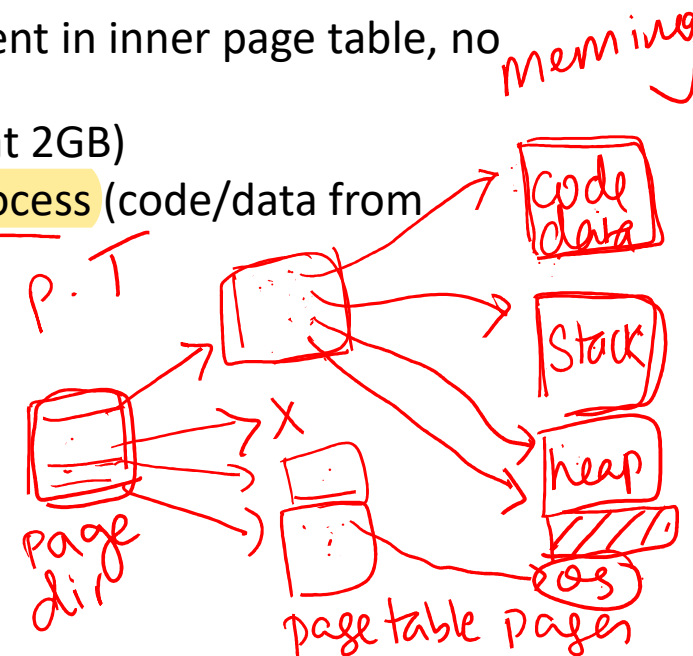
- Across all system calls, memory pages obtained from OS via kalloc()
- How is address space of process constructed in fork/exec?
    - Start with one page for the outer page directory
    - Allocate inner page tables on demand (if no entries present in inner page table, no need to allocate a page for it)
    - Add page table mappings for kernel code/data (starting at 2GB)
    - Allocate physical frames to store memory contents of process (code/data from executable, empty stack, ..)
    - Add mappings for user pages in page table
- How is address space of process expanded in sbrk?
    - Allocate physical frames for new virtual addresses
    - Add mappings for newly allocated pages in page table

# Recap: fork system call implementation

- Parent allocates new process in ptable, copies parent state to child
- Child process set to runnable, scheduler runs it at a later time
- Return value in parent is PID of child, return value in child is set to 0

*ret = fork();*

*P→C*

```
2579  int
2580  fork(void)
2581  {
2582    int i, pid;
2583    struct proc *np;
2584    struct proc *curproc = myproc();
2585
2586    // Allocate process.
2587    if((np = allocproc()) == 0){
2588      return -1;
2589    }
2590
2591    // Copy process state from proc.
2592    if((np->pgdir = copyuvm(curproc->pgdir, curproc->sz)) == 0){
2593      kfree(np->kstack);
2594      np->kstack = 0;
2595      np->state = UNUSED;
2596      return -1;
2597    }
2598    np->sz = curproc->sz;
2599    np->parent = curproc;
```

```
2600    *np->tf = *curproc->tf;
2601
2602    // Clear %eax so that fork returns 0 in the child.
2603    np->tf->eax = 0;
2604
2605    for(i = 0; i < NOFILE; i++)
2606      if(curproc->ofile[i])
2607        np->ofile[i] = filedup(curproc->ofile[i]);
2608    np->cwd = idup(curproc->cwd);
2609
2610    safestrcpy(np->name, curproc->name, sizeof(curproc->name));
2611
2612    pid = np->pid;
2613
2614    acquire(&ptable.lock);
2615
2616    np->state = RUNNABLE;
2617
2618    release(&ptable.lock);
2619
2620    return pid;
2621  }
```

*P  pid of child*

*Kernel mode*

*Kernel stack*
*Context*
*tf :*

*Memory copy*

*Copy trapframe*

*files*

*np*

*ptable*   *P  .  C*
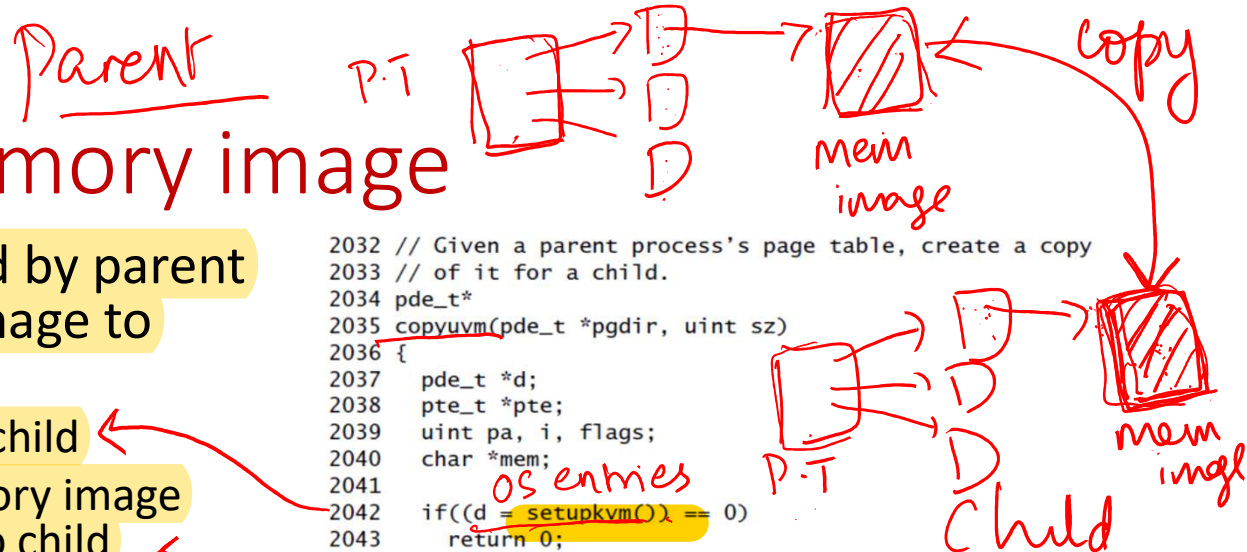
# Fork: copying memory image

- Function "copyuvm" called by parent to copy parent memory image to child
  - Create new page table for child
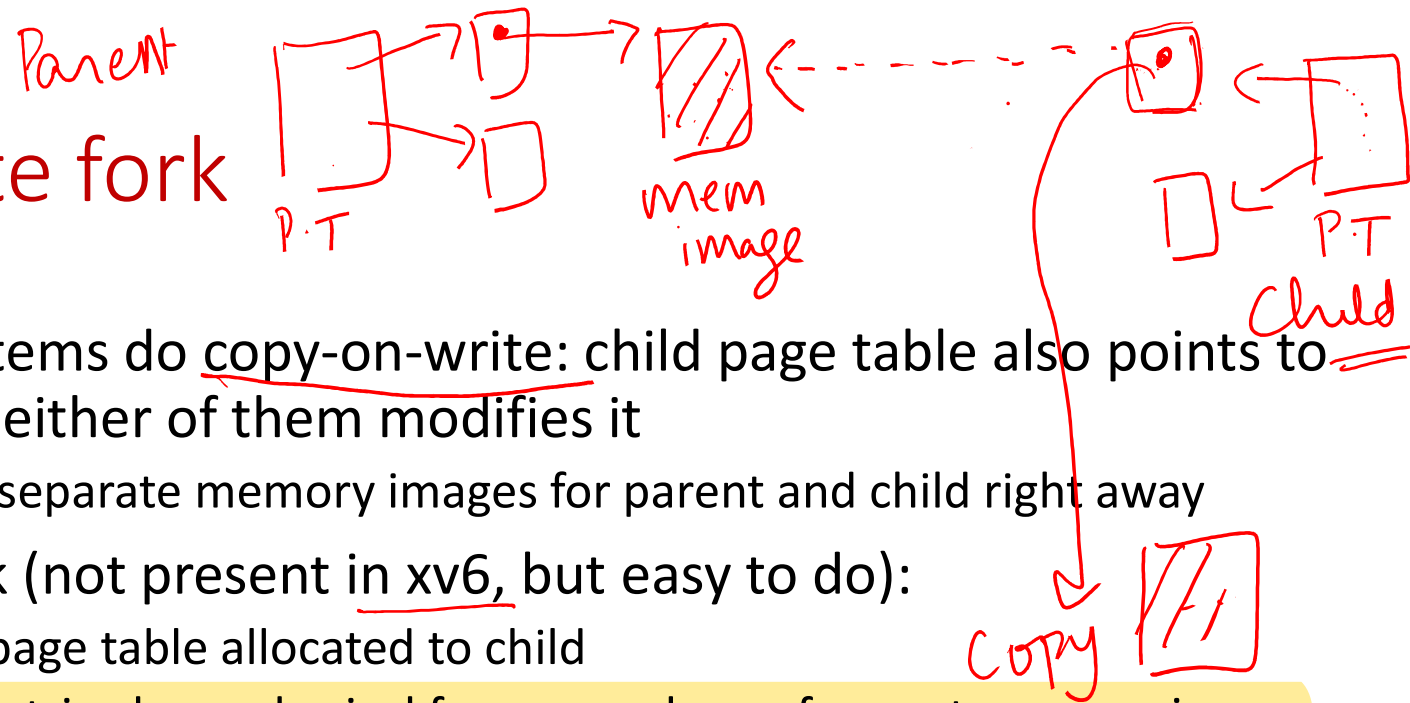  - Walk through parent memory image page by page and copy it to child
- For each page in parent
  - Fetch PTE, get physical address, permissions
  - Allocate new frame for child, copy contents of parent's page to new page of child
  - Add a PTE from virtual address to physical address of new page in child page table

```
2032  // Given a parent process's page table, create a copy
2033  // of it for a child.
2034  pde_t*
2035  copyuvm(pde_t *pgdir, uint sz)
2036  {
2037    pde_t *d;
2038    pte_t *pte;
2039    uint pa, i, flags;
2040    char *mem;
2041
2042    if((d = setupkvm()) == 0)
2043      return 0;
2044    for(i = 0; i < sz; i += PGSIZE){
2045      if((pte = walkpgdir(pgdir, (void *) i, 0)) == 0)
2046        panic("copyuvm: pte should exist");
2047      if(!(*pte & PTE_P))
2048        panic("copyuvm: page not present");
2049      pa = PTE_ADDR(*pte);
2050      flags = PTE_FLAGS(*pte);
2051      if((mem = kalloc()) == 0)
2052        goto bad;
2053      memmove(mem, (char*)P2V(pa), PGSIZE);
2054      if(mappages(d, (void*)i, PGSIZE, V2P(mem), flags) < 0) {
2055        kfree(mem);
2056        goto bad;
2057      }
2058    }
2059    return d;
2060
2061  bad:
2062    freevm(d);
2063    return 0;
2064  }
```
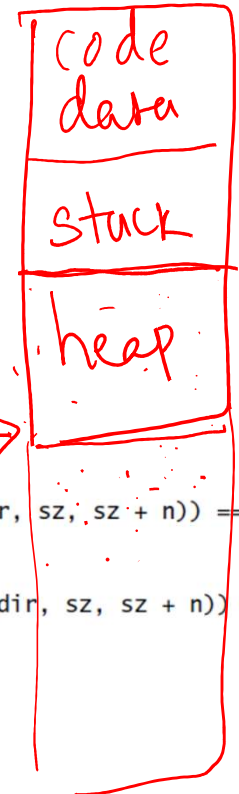
10

# Copy-on-write fork

- Real operating systems do copy-on-write: child page table also points to parent pages until either of them modifies it
  - Here, xv6 creates separate memory images for parent and child right away
- Copy-on-write fork (not present in xv6, but easy to do):
  - During fork, new page table allocated to child
  - Child page table entries have physical frame numbers of parent memory image pages only, no copy created for child
  - Parent's memory image is marked as read only
  - When parent or child tries to modify, MMU traps to OS
  - As part of trap handling, separate copy of memory image created
  - Finally, two separate copies of memory image for parent and child

Parent

P.T

mem image

P.T

Child

Copy

# Growing memory image: sbrk

- Initially heap is empty, program "break" (end of user memory) is at end of stack
  - sbrk() system call invoked by malloc to expand heap
- To grow memory, allocuvm allocates new pages, adds mappings into page table for new pages
- Whenever page table updated, must update cr3 register and TLB (done even during context switching)

```
2557 int
2558 growproc(int n)
2559 {
2560   uint sz;
2561   struct proc *curproc = myproc();
2562
2563   sz = curproc->sz;
2564   if(n > 0){
2565     if((sz = allocuvm(curproc->pgdir, sz, sz + n)) == 0)
2566       return -1;
2567   } else if(n < 0){
2568     if((sz = deallocuvm(curproc->pgdir, sz, sz + n)) == 0)
2569       return -1;
2570   }
2571   curproc->sz = sz;
2572   switchuvm(curproc);
2573   return 0;
2574 }
```
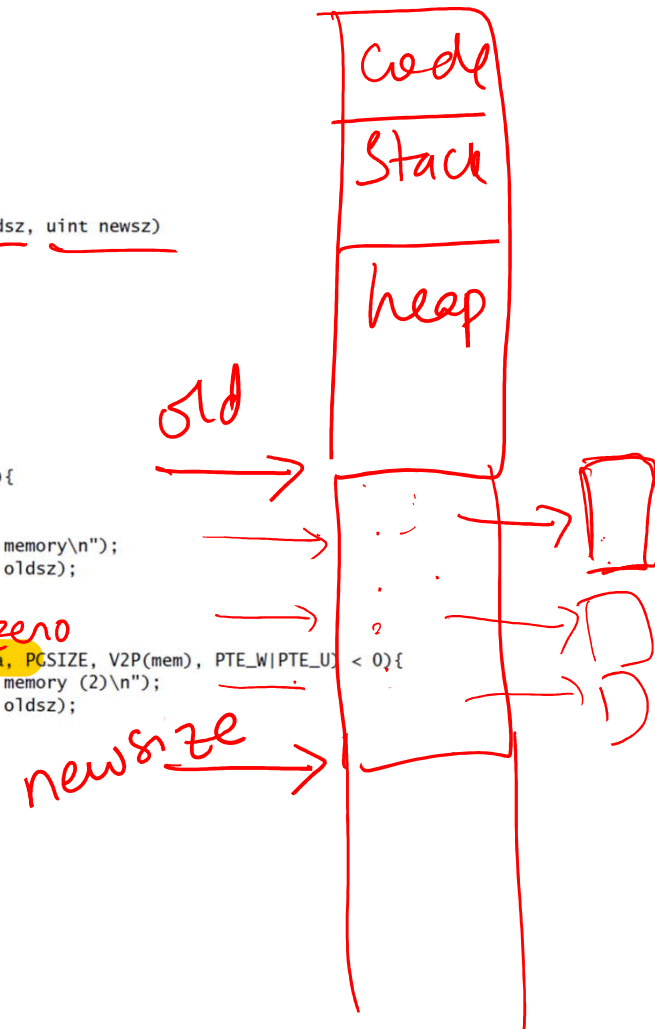
*(handwritten annotations: "break", "code", "data", "stuck", "heap", "update pagetable", "MMU")*

# allocuvm: grow address space

- Walk through new virtual addresses, page by page

- Allocate new frame, add mapping to page table with suitable user permissions

- Similarly deallocuvm shrinks memory image, frees up pages

```
1926 int
1927 allocuvm(pde_t *pgdir, uint oldsz, uint newsz)
1928 {
1929   char *mem;
1930   uint a;
1931
1932   if(newsz >= KERNBASE)
1933     return 0;
1934   if(newsz < oldsz)
1935     return oldsz;
1936
1937   a = PGROUNDUP(oldsz);
1938   for(; a < newsz; a += PGSIZE){
1939     mem = kalloc();
1940     if(mem == 0){
1941       cprintf("allocuvm out of memory\n");
1942       deallocuvm(pgdir, newsz, oldsz);
1943       return 0;
1944     }
1945     memset(mem, 0, PGSIZE);
1946     if(mappages(pgdir, (char*)a, PGSIZE, V2P(mem), PTE_W|PTE_U) < 0){
1947       cprintf("allocuvm out of memory (2)\n");
1948       deallocuvm(pgdir, newsz, oldsz);
1949       kfree(mem);
1950       return 0;
1951     }
1952   }
1953   return newsz;
1954 }
```

# Exec system call (1)

- Read ELF binary file from disk into memory
- Start with new page table, add mappings to new executable pages and grow virtual address space
  - Do not overwrite old page table yet

```
6609  int
6610  exec(char *path, char **argv)
6611  {
6612    char *s, *last;
6613    int i, off;
6614    uint argc, sz, sp, ustack[3+MAXARG+1];
6615    struct elfhdr elf;
6616    struct inode *ip;
6617    struct proghdr ph;
6618    pde_t *pgdir, *oldpgdir;
6619    struct proc *curproc = myproc();
6620
6621    begin_op();
6622
6623    if((ip = namei(path)) == 0){
6624      end_op();
6625      cprintf("exec: fail\n");
6626      return -1;
6627    }
6628    ilock(ip);
6629    pgdir = 0;
6630
6631    // Check ELF header
6632    if(readi(ip, (char*)&elf, 0, sizeof(elf)) != sizeof(elf))
6633      goto bad;
6634    if(elf.magic != ELF_MAGIC)
6635      goto bad;
6636
6637    if((pgdir = setupkvm()) == 0)
6638      goto bad;
```
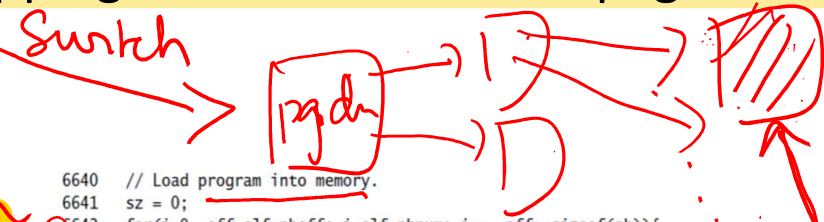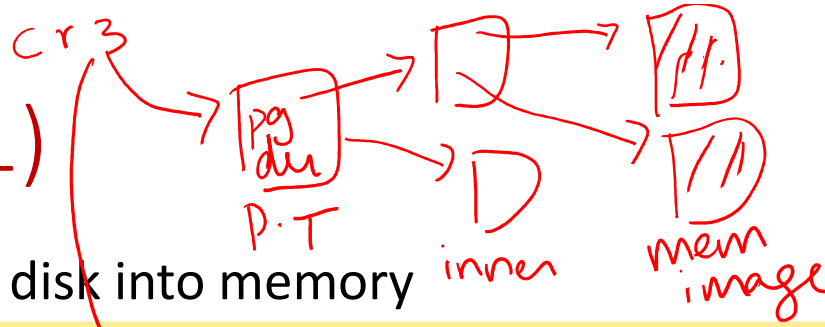
```
6640    // Load program into memory.
6641    sz = 0;
6642    for(i=0, off=elf.phoff; i<elf.phnum; i++, off+=sizeof(ph)){
6643      if(readi(ip, (char*)&ph, off, sizeof(ph)) != sizeof(ph))
6644        goto bad;
6645      if(ph.type != ELF_PROG_LOAD)
6646        continue;
6647      if(ph.memsz < ph.filesz)
6648        goto bad;
6649      if(ph.vaddr + ph.memsz < ph.vaddr)
6650        goto bad;
6651      if((sz = allocuvm(pgdir, sz, ph.vaddr + ph.memsz)) == 0)
6652        goto bad;
6653      if(ph.vaddr % PGSIZE != 0)
6654        goto bad;
6655      if(loaduvm(pgdir, (char*)ph.vaddr, ip, ph.off, ph.filesz) < 0)
6656        goto bad;
6657    }
6658    iunlockput(ip);
6659    end_op();
6660    ip = 0;
```
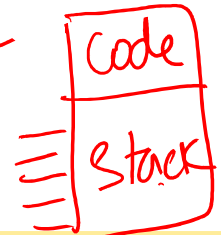
14

exec ( comm-d , " ... " , " ____ " )

a.out [Code] [Stack]

# Exec system call (2)

- After executable is copied to memory image, allocate 2 pages for stack (one is guard page, permissions cleared, access will trap)
- Push exec arguments onto user stack for main function of new program

main (argc, argv)
{ local ...

```
6662    // Allocate two pages at the next page boundary.
6663    // Make the first inaccessible.  Use the second as the user stack.
6664    sz = PGROUNDUP(sz);
6665    if((sz = allocuvm(pgdir, sz, sz + 2*PGSIZE)) == 0)
6666      goto bad;
6667    clearpteu(pgdir, (char*)(sz - 2*PGSIZE));
6668    sp = sz;
6669
6670    // Push argument strings, prepare rest of stack in ustack.
6671    for(argc = 0; argv[argc]; argc++) {
6672      if(argc >= MAXARG)
6673        goto bad;
6674      sp = (sp - (strlen(argv[argc]) + 1)) & ~3;
6675      if(copyout(pgdir, sp, argv[argc], strlen(argv[argc]) + 1) < 0)
6676        goto bad;
6677      ustack[3+argc] = sp;
6678    }
6679    ustack[3+argc] = 0;
6680
6681    ustack[0] = 0xffffffff;  // fake return PC
6682    ustack[1] = argc;
6683    ustack[2] = sp - (argc+1)*4;  // argv pointer
6684
6685    sp -= (3+argc+1) * 4;
6686    if(copyout(pgdir, sp, ustack, (3+argc+1)*4) < 0)
6687      goto bad;
6688
```
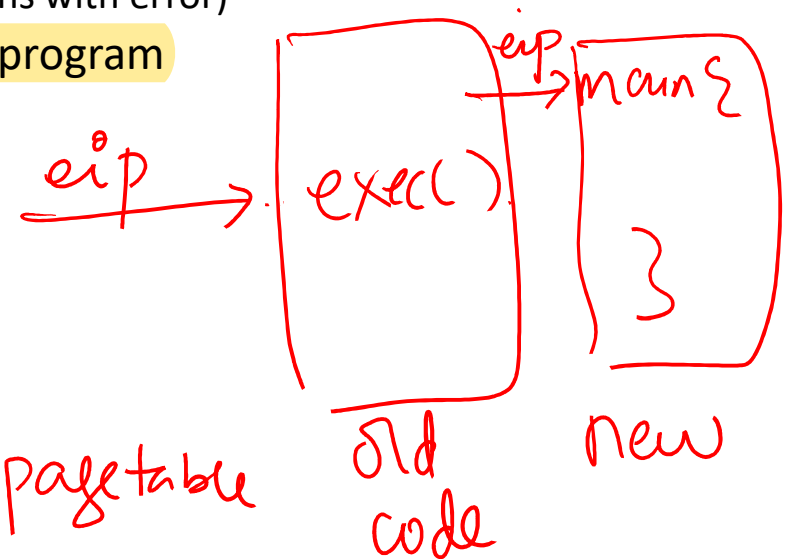
2 pages stack

command line argumen

15

# Exec system call (3)

- If no errors so far, switch to new page table that is pointing to new memory image
  - If any error, go back to old memory image (exec returns with error)
- Set eip in trapframe to start at entry point of new program
  - Returning from trap, process will run new executable

```
6689    // Save program name for debugging.
6690    for(last=s=path; *s; s++)
6691      if(*s == '/')
6692        last = s+1;
6693    safestrcpy(curproc->name, last, sizeof(curproc->name));
6694
6695    // Commit to the user image.
6696    oldpgdir = curproc->pgdir;
6697    curproc->pgdir = pgdir;
6698    curproc->sz = sz;
6699    curproc->tf->eip = elf.entry;   // main
6700    curproc->tf->esp = sp;
6701    switchuvm(curproc);
6702    freevm(oldpgdir);
6703    return 0;
6704
6705  bad:
6706    if(pgdir)
6707      freevm(pgdir);
6708    if(ip){
6709      iunlockput(ip);
6710      end_op();
6711    }
6712    return -1;
6713 }
```

*(handwritten annotations)* eip → exec() → main { } ; eip → main ; old code ; new ; Switch to new pagetable ; free up old memory