# Lecture 28: Memory management of user processes in xv6
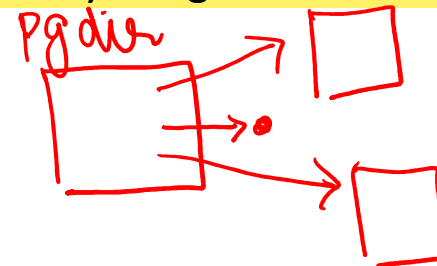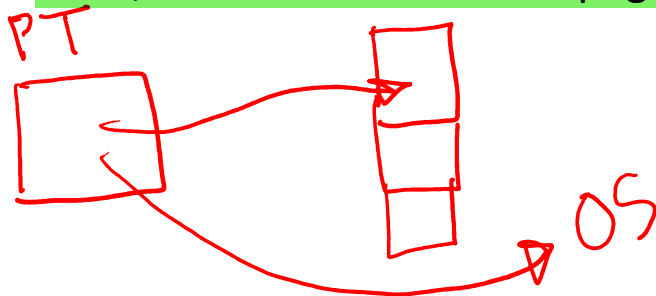
Mythili Vutukuru

IIT Bombay

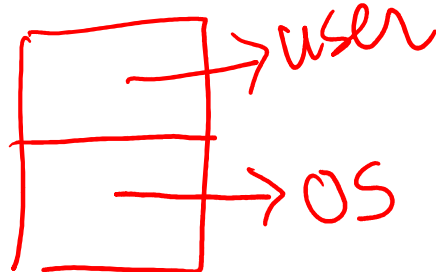https://www.cse.iitb.ac.in/~mythili/os/

# Memory management of user processes

- 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)
- Free list of kernel used to allocate memory for user processes via kalloc()
- 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 (expand heap)
- How is page table of a process constructed?
  - 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) as memory image created or updated
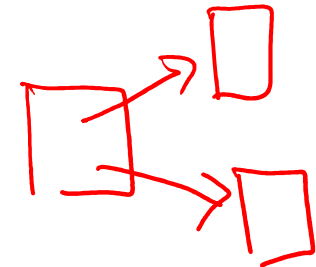
# Functions to build page table (1)

- Every page table begins with setting up kernel mappings in setupkvm()
- Outer pgdir allocated
- Kernel mappings defined in "kmap" added to page table by calling "mappages"
- After setupkvm(), user page table mappings added

```
1802 // This table defines the kernel's mappings, which are present in
1803 // every process's page table.
1804 static struct kmap {
1805   void *virt;
1806   uint phys_start;
1807   uint phys_end;
1808   int perm;
1809 } kmap[] = {
1810   { (void*)KERNBASE, 0,             EXTMEM,   PTE_W}, // I/O space
1811   { (void*)KERNLINK, V2P(KERNLINK), V2P(data), 0},    // kern text+rodata
1812   { (void*)data,     V2P(data),     PHYSTOP,  PTE_W}, // kern data+memory
1813   { (void*)DEVSPACE, DEVSPACE,      0,        PTE_W}, // more devices
1814 };
1815
1816 // Set up kernel part of a page table.
1817 pde_t*
1818 setupkvm(void)
1819 {
1820   pde_t *pgdir;
1821   struct kmap *k;
1822
1823   if((pgdir = (pde_t*)kalloc()) == 0)
1824     return 0;
1825   memset(pgdir, 0, PGSIZE);
1826   if (P2V(PHYSTOP) > (void*)DEVSPACE)
1827     panic("PHYSTOP too high");
1828   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1829     if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
1830                 (uint)k->phys_start, k->perm) < 0) {
1831       freevm(pgdir);
1832       return 0;
1833     }
1834   return pgdir;
1835 }
```
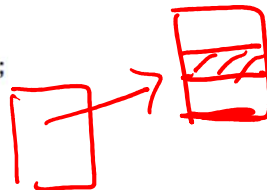
3

# Functions to build page table (2)

- Page table entries added by "mappages" actually assigns physical page frame to the pte.
  - Arguments: page directory, range of virtual addresses, physical addresses to map to, permissions of the pages
  - For each page, walks page table, get pointer to PTE via function "walkpgdir", fills it with physical address and permissions
- Function "walkpgdir" walks page table, returns PTE of a virtual address
  - Can allocate inner page table if it doesn't exist

VA → PA

2 level page table, pde: page directory entry
pte: page table entry.

```
1756 // Create PTEs for virtual addresses starting at va that refer to
1757 // physical addresses starting at pa. va and size might not
1758 // be page-aligned.
1759 static int
1760 mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
1761 {
1762   char *a, *last;
1763   pte_t *pte;
1764
1765   a = (char*)PGROUNDDOWN((uint)va);
1766   last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
1767   for(;;){
1768     if((pte = walkpgdir(pgdir, a, 1)) == 0)
1769       return -1;
1770     if(*pte & PTE_P)
1771       panic("remap");
1772     *pte = pa | perm | PTE_P;
1773     if(a == last)
1774       break;
1775     a += PGSIZE;
1776     pa += PGSIZE;
1777   }
1778   return 0;
1779 }
```

```
1731 // Return the address of the PTE in page table pgdir
1732 // that corresponds to virtual address va.  If alloc!=0,
1733 // create any required page table pages.
1734 static pte_t *
1735 walkpgdir(pde_t *pgdir, const void *va, int alloc)
1736 {
1737   pde_t *pde;
1738   pte_t *pgtab;
1739
1740   pde = &pgdir[PDX(va)];
1741   if(*pde & PTE_P){
1742     pgtab = (pte_t*)P2V(PTE_ADDR(*pde));
1743   } else {
1744     if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
1745       return 0;
1746     // Make sure all those PTE_P bits are zero.
1747     memset(pgtab, 0, PGSIZE);
1748     // The permissions here are overly generous, but they can
1749     // be further restricted by the permissions in the page table
1750     // entries, if necessary.
1751     *pde = V2P(pgtab) | PTE_P | PTE_W | PTE_U;
1752   }
1753   return &pgtab[PTX(va)];
1754 }
```

allocate

# Fork: copying memory image

```
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   }
```
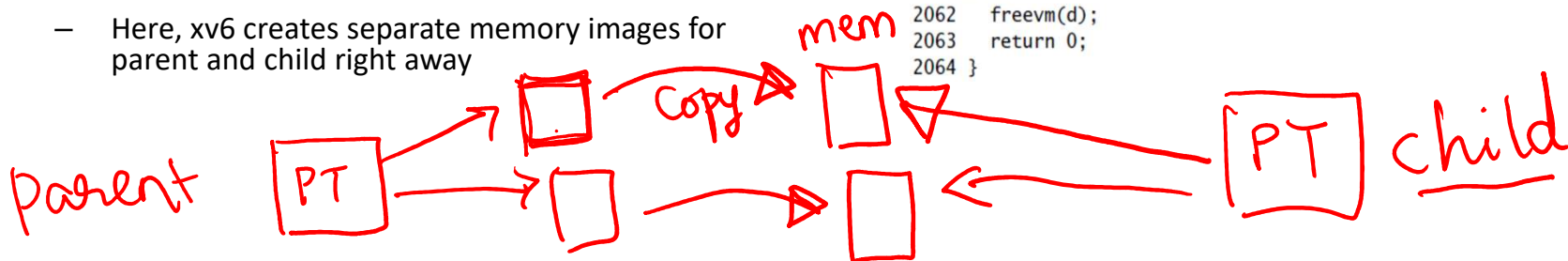
- Function "copyuvm" called by parent to copy parent memory image to child
  - Create new page table for child
  - Walk through parent memory image page by page and copy it to child, while adding child page table mappings
- For each page in parent
  - fetch PTE, get physical address, permissions
  - Allocate new page for child, and 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
- 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

```
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   }
```

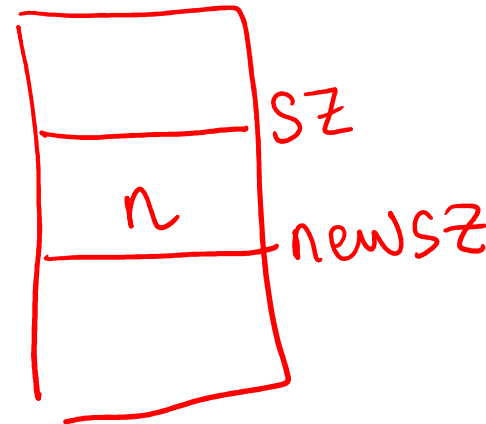copy the contents from parent's phy frame to child's phy frame.

# 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) This step is important.

```
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 }
```

# allocuvm: grow address space

- Walk through new virtual addresses to be added in page size chunks
- Allocate new page, add it 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;
```
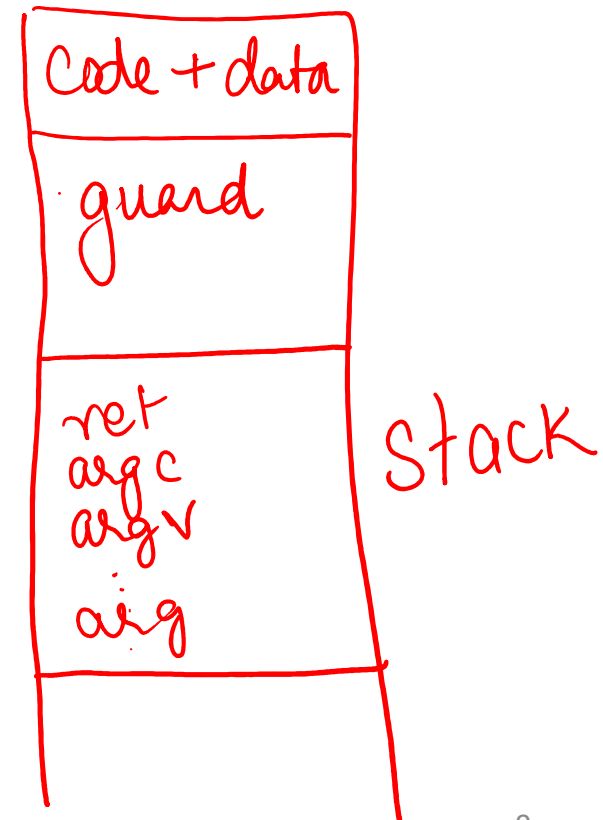
```
6640   // Load program into memory.
6641   sz = 0;
6642   for(i=0, off=elf.phoff; i<elf.phnum; i++, off+=sizeof(ph)){
6643     if(readi(ip, (char*)&ph, off, sizeof(ph)) != sizeof(ph))
6644       goto bad;
6645     if(ph.type != ELF_PROG_LOAD)
6646       continue;
6647     if(ph.memsz < ph.filesz)
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;
```

8

# 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
  - Stack has return address, argc, argv array (pointers to variable sized arguments), and the arguments themselves

```
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
```
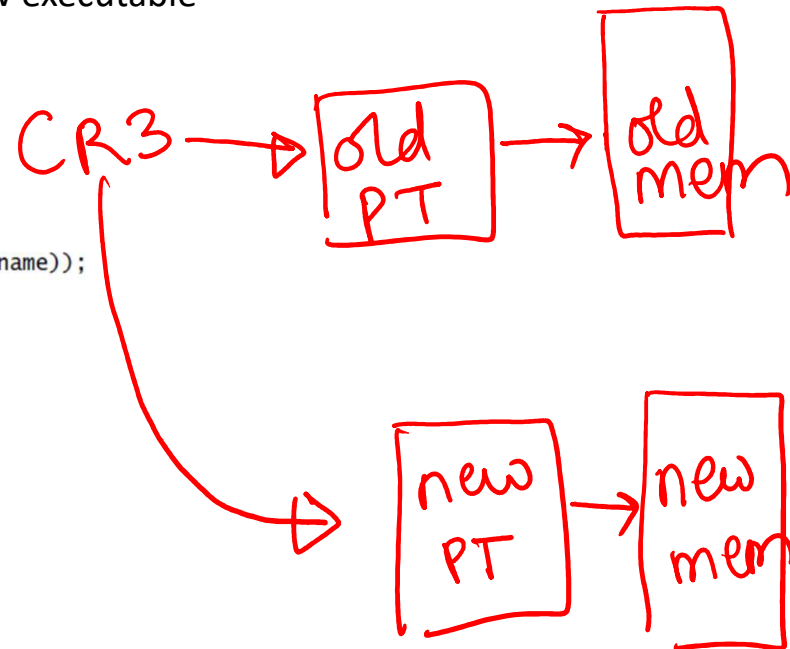
argc
argv

Code + data

guard

ret
argc
argv

arg

Stack

# 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  }
```

# Summary

- Memory management for user processes
    - Build page table: start with kernel mappings, add user entries to build virtual address space
    - Memory management code in fork, exec, sbrk