

slides 3

IPC

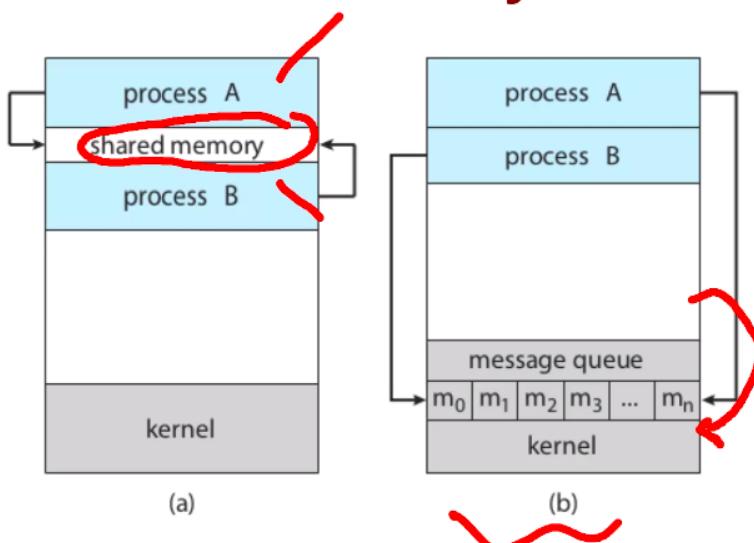
Revision of process related concepts

- PCB, struct proc
- Process lifecycle – different states
- Memory layout
- Memory management
- Interrupts handling, system call handling, code from xv6
- Scheduler, code of scheduler in xv6

IPC: Inter Process Communication

- Processes within a system may be independent or cooperating
 - Cooperating process can affect or be affected by other processes, including sharing data
 - Reasons for cooperating processes:
 - Information sharing, e.g. copy paste
 - Computation speedup, e.g. matrix multiplication
 - Modularity, e.g. chrome – separate process for display, separate for fetching data
 - Convenience,
 - Cooperating processes need interprocess communication (IPC)
 - Two models of IPC
 - Shared memory
 - Message passing
-

Shared Memory Vs Message Passing



- Each requires OS to provide system calls for
- Creating the IPC mechanism
 - To read/write using the IPC mechanism
 - Delete the IPC mechanism
- Note:** processes communicating with each other with the help of OS!

Figure 3.11 Communications models. (a) Shared memory. (b) Message passing.



Example of co-operating processes: Producer Consumer Problem

- Shared data

```
#define BUFFER_SIZE 10
typedef struct {
    ...
} item;

item buffer[BUFFER_SIZE];
int in = 0;
int out = 0;
```

- Can only use BUFFER_SIZE-1 elements



Example of co-operating processes: Producer Consumer Problem

- Code of Producer

```
while (true) {
    /* Produce an item */
        while (((in = (in + 1) % BUFFER_SIZE) == out)
            ; /* do nothing -- no free buffers */
    buffer[in] = item;
    in = (in + 1) % BUFFER_SIZE;
}
```

Example of co-operating processes: Producer Consumer Problem

- **Code of Consumer**

```
while (true) {  
    while (in == out)  
        ; // do nothing -- nothing to consume  
  
    { // remove an item from the buffer  
        item = buffer[out];  
        out = (out + 1) % BUFFER SIZE;  
        return item;  
    }  
}
```

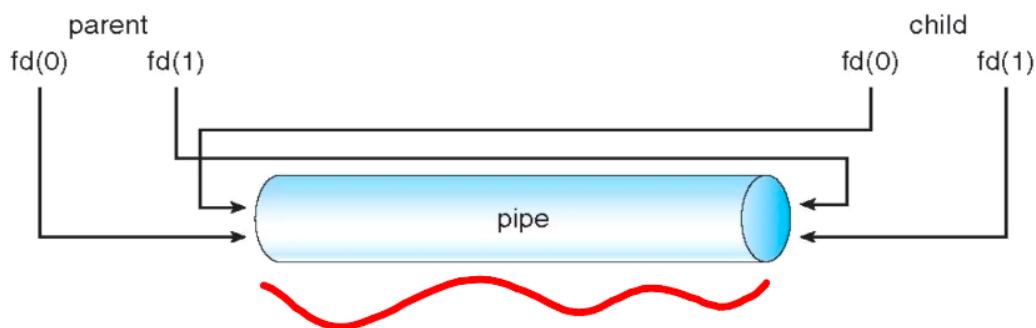
Pipes for IPC

- **Two types**

- **Unnamed Pipes or ordinary pipes**
- **Named Pipe**

Ordinary pipes

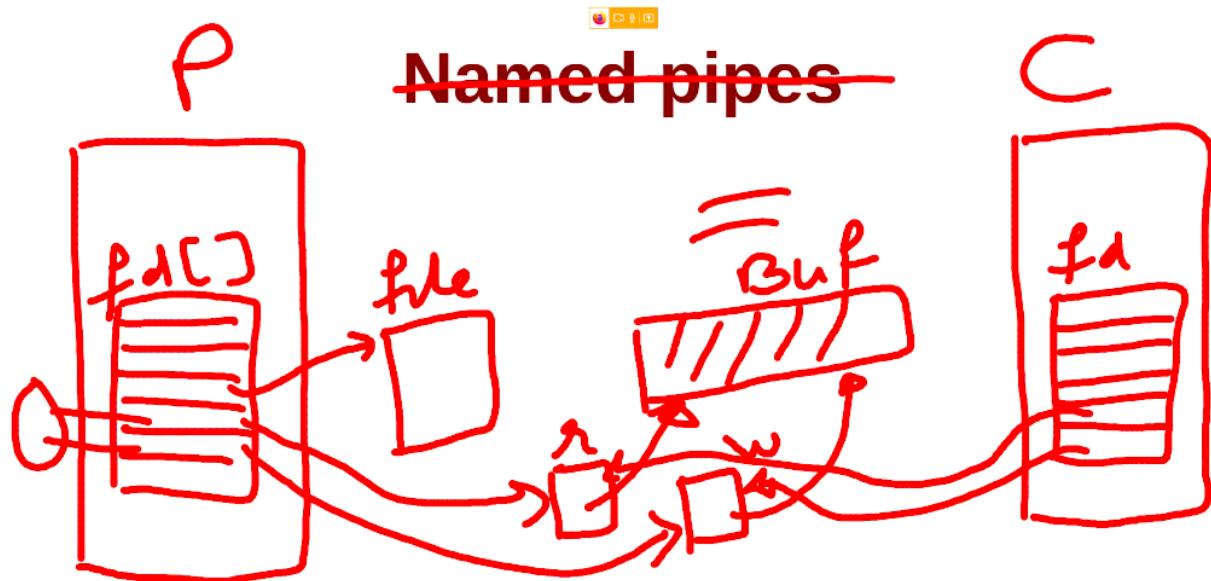
- Ordinary Pipes allow communication in standard producer-consumer style
- Producer writes to one end (the write-end of the pipe)
- Consumer reads from the other end (the read-end of the pipe)
- Ordinary pipes are therefore unidirectional
- Requires a parent-child (or sibling, etc) kind of relationship between communicating processes



```

Terminal Feb 8 14:47 abhijit@abhijit-laptop: /courses/os-2022/programs/pipe abhijit@abhijit-laptop: -/abhijit/coep/courses/os-2022/programs/pipe
1 #include <sys/types.h>
2 #include <sys/stat.h>
3 #include <fcntl.h>
4 #include <stdio.h>
5 #include <errno.h>
6 #include <unistd.h>
7
8 int main(int argc, char *argv[]) {
9     int fd, fd2;
10    char str[1600];
11    int pfd[2];
12
13    pipe(pfd);
14    /*pfd[0] - read file descriptor
15    pfd[1] - write file descriptor */
16    write(pfd[1], "hello\n", 6);
17    read(pfd[0], str, 6);
18    str[6] = '\0';
19    printf("%s", str);
20 }
~"pipe.c" 20L, 378C 1,1 All

```



pipe is unidirectional, so all unused ends should be closed

A screenshot of a Linux desktop environment showing a terminal window. The terminal window has two tabs: the first tab shows a C program for pipes, and the second tab shows the same program with some output. The C code includes #include directives for unistd.h, stdlib.h, sys/types.h, sys/wait.h, and stdio.h. It defines a main() function that creates a pipe, forks, and then reads from or writes to the pipe based on its process ID (pid). The terminal window also shows the date and time (Feb 8 14:54) and the user's name (abhijit@abhijit-laptop).

```
1 #include <unistd.h>
2 #include <stdlib.h>
3 #include <sys/types.h>
4 #include <sys/wait.h>
5 #include <stdio.h>
6 int main() {
7     char cmd[128], string[16];
8     int cpid, pfd[2];
9     pipe(pfd);
10    cpid = fork();
11    if(cpid == 0) {
12        close(pfd[0]);
13        write(pfd[1], "hello\n", 6);
14        exit(0);
15    } else {
16        close(pfd[1]);
17        read(pfd[0], string, 6);
18        string[6] = '\0';
19        printf("Parent read %s\n", string);
20    }
21    return 0;
22 }
23
```

all the writes are delayed by kernel, and it will do it on its own convinience.

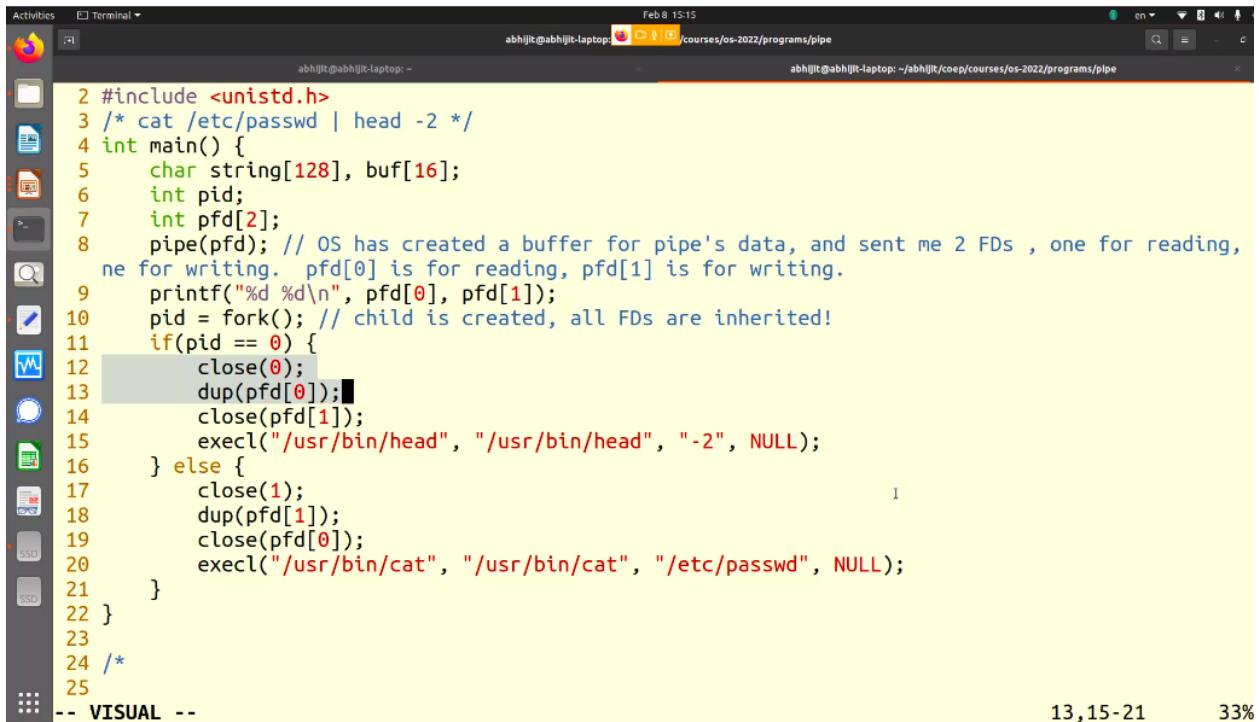
A screenshot of a Linux desktop environment showing a terminal window. The terminal window has two tabs: the first tab shows a C program for pipes, and the second tab shows the same program with some output. The C code includes #include directives for sys/types.h, sys/stat.h, fcntl.h, errno.h, stdio.h, and unistd.h. It defines a main() function that opens a file descriptor (fd) for reading, scans for input, and then prints it back. The terminal window also shows the date and time (Feb 8 15:03) and the user's name (abhijit@abhijit-laptop).

```
1 #include <sys/types.h>
2 #include <sys/stat.h>
3 #include <fcntl.h>
4 #include <errno.h>
5 #include <stdio.h>
6 #include <unistd.h>
7
8 int main(int argc, char *argv[]) {
9     int fd, n, i;
10    char buf[128];
11
12    close(0);
13    fd = open(argv[1], O_RDONLY);
14    if(fd == -1) {
15        perror("open failed");
16        return errno;
17    }
18    scanf("%s", buf);
19    printf("%s\n", buf);
20    return 0;
21 }
```

for this program, open will return file descriptor 0, and hence will read from file

see strace

exec, copies the file descriptors



```
2 #include <unistd.h>
3 /* cat /etc/passwd | head -2 */
4 int main() {
5     char string[128], buf[16];
6     int pid;
7     int pfd[2];
8     pipe(pfd); // OS has created a buffer for pipe's data, and sent me 2 FDs , one for reading,
ne for writing. pfd[0] is for reading, pfd[1] is for writing.
9     printf("%d %d\n", pfd[0], pfd[1]);
10    pid = fork(); // child is created, all FDs are inherited!
11    if(pid == 0) {
12        close(0);
13        dup(pfd[0]);
14        close(pfd[1]);
15        execl("/usr/bin/head", "/usr/bin/head", "-2", NULL);
16    } else {
17        close(1);
18        dup(pfd[1]);
19        close(pfd[0]);
20        execl("/usr/bin/cat", "/usr/bin/cat", "/etc/passwd", NULL);
21    }
22 }
23
24 /*
25
-- VISUAL --
```

13,15-21

33%

The screenshot shows a Linux desktop environment with a terminal window open. The terminal window has two tabs: one with the command 'abhi@abhi-laptop: ~' and another with the command 'abhi@abhi-laptop: ~/coep/courses/os-2022/programs/pipe'. The code in the terminal is a C program that creates a pipe, forks, and execs to demonstrate file descriptor inheritance. The code is as follows:

```
6 int pid;
7 int pfd[2];
8 printf("%d %d\n", pfd[0], pfd[1]);
9 while(1) {
10     pipe(pfd);
11     pid = fork(); // child is created, all FDs are inherited!
12     if(pid == 0) {
13         close(0);
14         dup(pfd[0]);
15         close(pfd[1]);
16         execl("/usr/bin/head", "/usr/bin/head", "-2", NULL);
17     } else {
18         pid = fork();
19         if(pid == 0) {
20             close(1);
21             dup(pfd[1]);
22             close(pfd[0]);
23             execl("/usr/bin/cat", "/usr/bin/cat", "/etc/passwd", NULL);
24         } else {
25             close(pfd[0]);
26             close(pfd[1]);
27             exit(0);
28         }
29     }
30 }
```

At the bottom right of the terminal window, there are status indicators: '25,18-30' and '45%'. The desktop interface includes a dock with icons for various applications like a file manager, terminal, and browser.

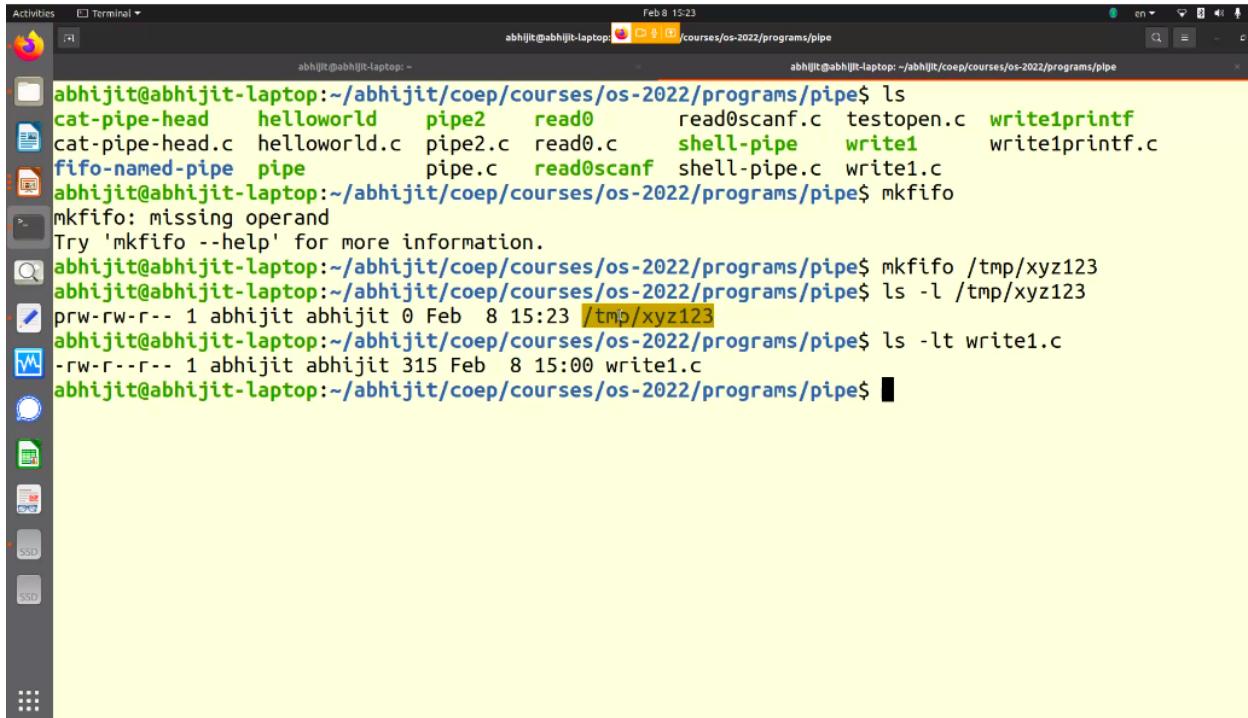
when we use process1 | process2

process 1 is continuously writing and process2 is continuously reading from pipe

Named pipes

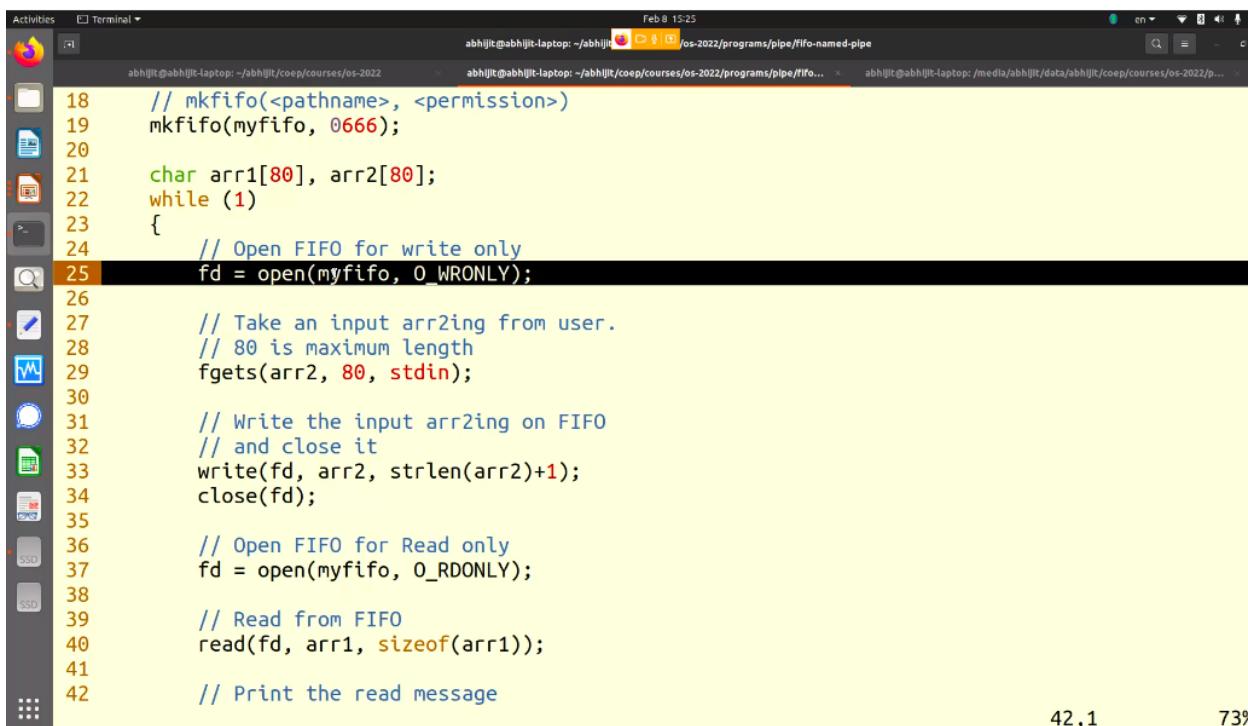
- Also called FIFO
- Processes can create a “file” that acts as pipe. Multiple processes can share the file to read/write as a FIFO
- Named Pipes are more powerful than ordinary pipes
- Communication is bidirectional
- No parent-child relationship is necessary between the communicating processes
- Several processes can use the named pipe for communication
- Provided on both UNIX and Windows systems
- Is not deleted automatically by OS

named pipe files are prefixed by p in file permissions



A screenshot of a Linux desktop environment showing a terminal window. The terminal window has two tabs: 'abhijit@abhijit-laptop: ~/abhijit/coep/courses/os-2022/programs/pipe' and 'abhijit@abhijit-laptop: ~/abhijit/coep/courses/os-2022/programs/pipe'. The terminal shows the following command-line session:

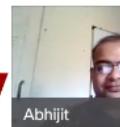
```
abhijit@abhijit-laptop:~/abhijit/coep/courses/os-2022/programs/pipe$ ls
cat-pipe-head  helloworld  pipe2  read0      read0scanf.c  testopen.c  write1printf
cat-pipe-head.c helloworld.c pipe2.c  read0.c    shell-pipe   write1    write1printf.c
fifo-named-pipe  pipe      pipe.c  read0scanf shell-pipe.c write1.c
abhijit@abhijit-laptop:~/abhijit/coep/courses/os-2022/programs/pipe$ mkfifo
mkfifo: missing operand
Try 'mkfifo --help' for more information.
abhijit@abhijit-laptop:~/abhijit/coep/courses/os-2022/programs/pipe$ mkfifo /tmp/xyz123
abhijit@abhijit-laptop:~/abhijit/coep/courses/os-2022/programs/pipe$ ls -l /tmp/xyz123
prw-rw-r-- 1 abhijit abhijit 0 Feb  8 15:23 /tmp/xyz123
abhijit@abhijit-laptop:~/abhijit/coep/courses/os-2022/programs/pipe$ ls -lt write1.c
-rw-r--r-- 1 abhijit abhijit 315 Feb  8 15:00 write1.c
abhijit@abhijit-laptop:~/abhijit/coep/courses/os-2022/programs/pipe$ █
```



A screenshot of a Linux desktop environment showing a terminal window. The terminal window has three tabs: 'abhijit@abhijit-laptop: ~/abhijit/coep/courses/os-2022', 'abhijit@abhijit-laptop: ~/abhijit/coep/courses/os-2022/programs/pipe/fifo-named-pipe', and 'abhijit@abhijit-laptop: ~/media/abhijit/data/abhijit/coep/courses/os-2022/p...'. The terminal shows the following command-line session:

```
18 // mkfifo(<pathname>, <permission>)
19 mkfifo(myfifo, 0666);
20
21 char arr1[80], arr2[80];
22 while (1)
23 {
24     // Open FIFO for write only
25     fd = open(myfifo, O_WRONLY);
26
27     // Take an input arr2ing from user.
28     // 80 is maximum length
29     fgets(arr2, 80, stdin);
30
31     // Write the input arr2ing on FIFO
32     // and close it
33     write(fd, arr2, strlen(arr2)+1);
34     close(fd);
35
36     // Open FIFO for Read only
37     fd = open(myfifo, O_RDONLY);
38
39     // Read from FIFO
40     read(fd, arr1, sizeof(arr1));
41
42     // Print the read message
```

Shared memory

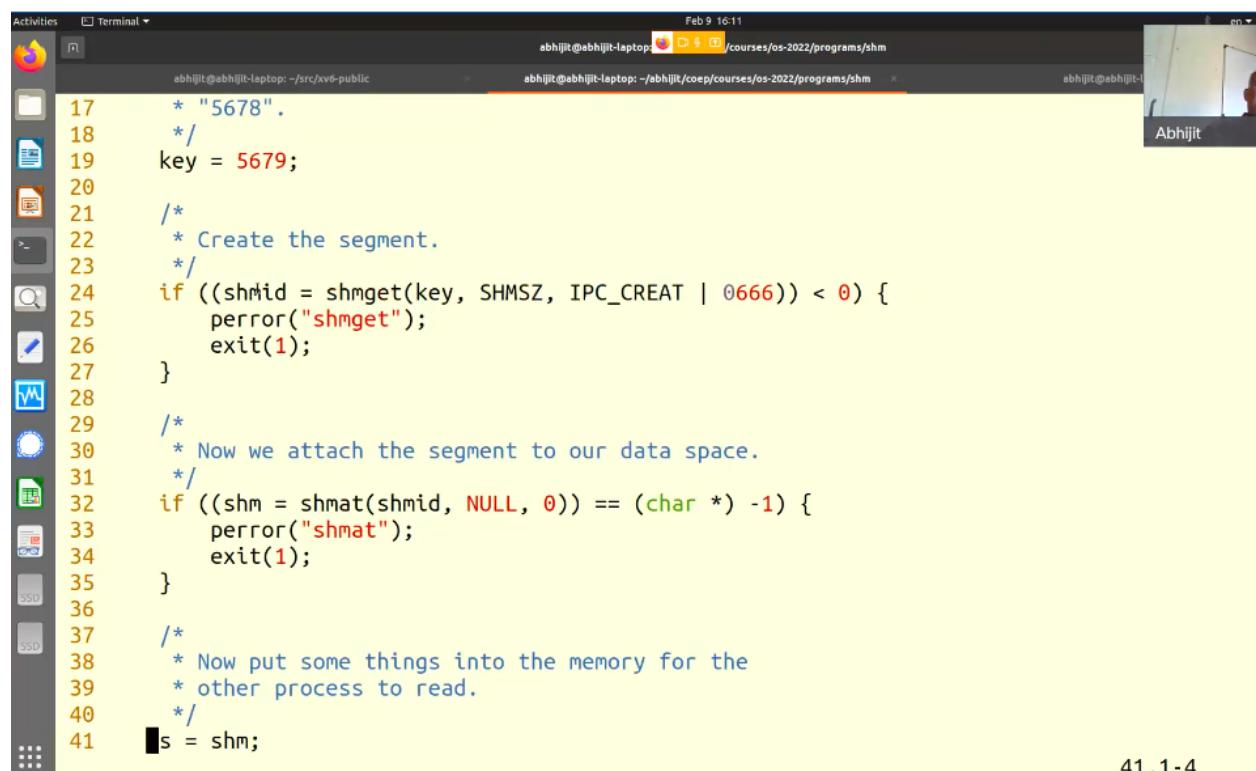


System V shared memory

- Process first creates shared memory segment
`segment id = shmget(IPC_PRIVATE, size, S_IRUSR | S_IWUSR);`
- Process wanting access to that shared memory must attach to it
`shared_memory = (char *) shmat(id, NULL, 0);`
- Now the process could write to the shared memory
`sprintf(shared_memory, "Writing to shared memory");`
- When done, a process can detach the shared memory from its address space

`shmdt()` detaches

`shm` is a `char*`, `key` should be common with all process sharing the shared memory

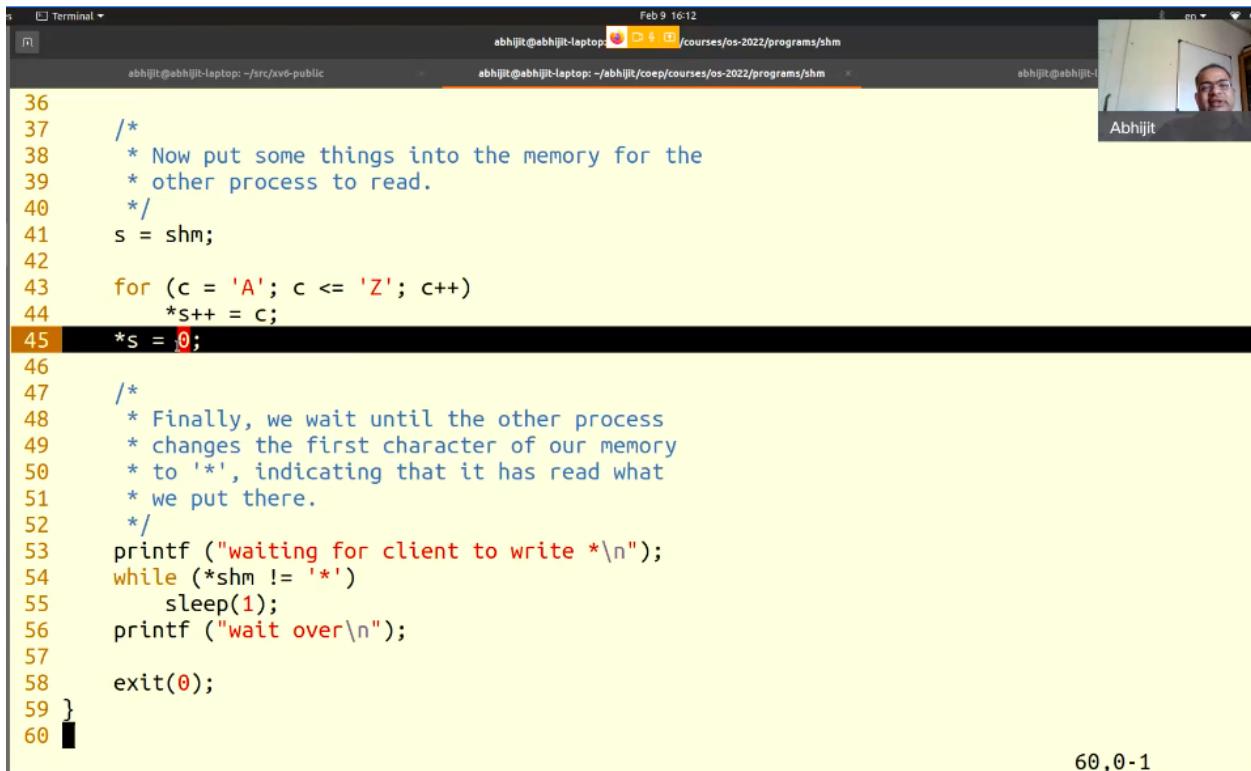


```

17     * "5678".
18     */
19 key = 5679;
20
21 /*
22  * Create the segment.
23 */
24 if ((shmid = shmget(key, SHMSZ, IPC_CREAT | 0666)) < 0) {
25     perror("shmget");
26     exit(1);
27 }
28
29 /*
30  * Now we attach the segment to our data space.
31 */
32 if ((shm = shmat(shmid, NULL, 0)) == (char *) -1) {
33     perror("shmat");
34     exit(1);
35 }
36
37 /*
38  * Now put some things into the memory for the
39  * other process to read.
40 */
41 s = shm;

```

41.1-4



```
36
37     /*
38      * Now put some things into the memory for the
39      * other process to read.
40     */
41     s = shm;
42
43     for (c = 'A'; c <= 'Z'; c++)
44         *s++ = c;
45     *s = 0;
46
47     /*
48      * Finally, we wait until the other process
49      * changes the first character of our memory
50      * to '*', indicating that it has read what
51      * we put there.
52     */
53     printf ("waiting for client to write *\n");
54     while (*shm != '*')
55         sleep(1);
56     printf ("wait over\n");
57
58     exit(0);
59 }
60
```

60,0-1

the above shared memory systemcall were system 5 compliant,
linux gives POSIX(portable operating system interfaces) systemcalls too

POSIX Shared Memory

- `shm_open`
- `ftruncate`
- `Mmap`
- See the example in Textbook

ordinary pipes are unidirectional, whereas named pipes are bidirectional

Message Passing

indirect message passing example

receiver



abhi@abhi-laptop: ~/src/xv6-public

```

11 int main()
12 {
13     key_t key;
14     int msgid;
15
16     // ftok to generate unique key
17     key = ftok("progfile", 65);
18
19     // msgget creates a message queue
20     // and returns identifier
21     msgid = msgget(key, 0666 | IPC_CREAT);
22
23     // msgrcv to receive message
24     msgrcv(msgid, &message, sizeof(message), 1, 0);
25
26     // display the message
27     printf("Data Received is : %s \n",
28           message.msg_text);
29
30     // to destroy the message queue
31     msgctl(msgid, IPC_RMID, NULL);
32
33

```

sender



abhi@abhi-laptop: ~/src/xv6-public

```

13 int main()
14 {
15     key_t key;
16     int msgid;
17
18     // ftok to generate unique key
19     key = ftok("progfile", 65);
20
21     // msgget creates a message queue
22     // and returns identifier
23     msgid = msgget(key, 0666 | IPC_CREAT);
24     message.msg_type = 1;
25
26     printf("Write Data : ");
27     fgets(message.msg_text, MAX, stdin);
28
29     // msgsnd to send message
30     msgsnd(msgid, &message, sizeof(message), 0);
31
32     // display the message
33     printf("Data send is : %s \n", message.msg_text);
34
35     return 0;

```

message passing using naming



Message passing using direct communication

- Processes must name each other explicitly:
 - **send (P, message)** – send a message to process P
 - receive(Q, message) – receive a message from process Q
- Properties of communication link
 - Links are established automatically
 - A link is associated with exactly one pair of communicating processes
 - Between each pair there exists exactly one link
 - The link may be unidirectional, but is usually bi-directional



Message passing using IN-direct communication

- Operations
 - create a new mailbox
 - send and receive messages through mailbox
 - destroy a mailbox
- Primitives are defined as:
 - **send(A, message)** – send a message to mailbox A
 - **receive(A, message)** – receive a message from mailbox A

here A is the mailbox, not a sender or receiver

Message passing using IN-direct communication



- **Mailbox sharing**
 - P1, P2, and P3 share mailbox A
 - P1 sends; P2 and P3 receive
 - Who gets the message?
- **Solutions**
 - Allow a link to be associated with at most two processes
 - Allow only one process at a time to execute a receive operation
 - Allow the system to select arbitrarily the receiver. Sender is notified who the receiver was.

Message Passing implementation: Synchronization issues



- **Message passing may be either blocking or non-blocking**
- **Blocking is considered synchronous**
 - Blocking send has the sender block until the message is received
 - Blocking receive has the receiver block until a message is available
- **Non-blocking is considered asynchronous**
 - Non-blocking send has the sender send the message and continue
 - Non-blocking receive has the receiver receive a valid message or null

Producer consumer using blocking send receive



Producer

```
message next_produced;  
while (true) {  
    /* produce an item in  
    next_produced */  
    send(next_produced);  
}
```

Consumer

```
message  
next_consumed;  
while (true) {  
    receive(next_consumed);  
}
```

writing blocking code is much easier.

Message Passing implementation: choice Buffering



- Queue of messages attached to the link; implemented in one of three ways
 1. Zero capacity – 0 messages
 - Sender must wait for receiver (rendezvous)
 2. Bounded capacity – finite length of n messages
 - Sender must wait if link full
 3. Unbounded capacity – infinite length
 - Sender never waits

```

20     return exec(path, args);
21 }
22 int
23 sys_pipe(void)
24 {
25     int *fd;
26     struct file *rf, *wf;
27     int fd0, fd1;
28
29     if(argptr(0, (void*)&fd, 2*sizeof(fd[0])) < 0)
30         return -1;
31     if(pipealloc(&rf, &wf) < 0)
32         return -1;
33     fd0 = -1;
34     if((fd0 = fdalloc(rf)) < 0 || (fd1 = fdalloc(wf)) < 0){
35         if(fd0 >= 0)
36             myproc()->ofile[fd0] = 0;
37         fileclose(rf);
38         fileclose(wf);
39         return -1;
40     }
41     fd[0] = fd0;
42     fd[1] = fd1;
43     return 0;
44 }
sysfile.c" 444L, 7377C

```

sys_pipe is called for pipe systemcall, which creates 2 file pointers and then calls pipealloc, does error checking and then allocate fd array and returns;

```

12
13 struct pipe {
14     struct spinlock lock;
15     char data[PIPESIZE];
16     uint nread;      // number of bytes read
17     uint nwrite;     // number of bytes written
18     int readopen;   // read fd is still open
19     int writeopen;  // write fd is still open
20 };
21
22 int

```

PIPESIZE is 512, data is the buffer used to store data,

inside pipealloc, it calls filealloc, which gives a struct file, then there is a call to kalloc() which will give you a 4kb page, and it is typecasted in pipe pointer.

then variables are initialized, then file structure, f0 and f1 are being initialized, one for writing and another for reading and then just return

```
18 int readopen; // read fd is still open
19 int writeopen; // write fd is still open
20 };
21
22 int
23 pipealloc(struct file **f0, struct file **f1)
24 {
25     struct pipe *p;
26
27     p = 0;
28     *f0 = *f1 = 0;
29     if((*f0 = filealloc()) == 0 || (*f1 = filealloc()) == 0)
30         goto bad;
31     if((p = (struct pipe*)kalloc()) == 0)
32         goto bad;
33     p->readopen = 1;
34     p->writeopen = 1;
35     p->nwrite = 0;
36     p->nread = 0;
37     initlock(&p->lock, "pipe");
38     (*f0)->type = FD_PIPE;
39     (*f0)->readable = 1;
40     (*f0)->writable = 0;
41     (*f0)->pipe = p;
42     (*f1)->type = FD_PIPE;
43     (*f1)->readable = 0;
44     (*f1)->writable = 1;
45     (*f1)->pipe = p;
'pipe.c' 121L, 2411C
```

30,2

1 +43

struct file: → which either points to a actual file(inode pointer is used) else a pipe pointer is used

```

1 struct file {
2     enum { FD_NONE, FD_PIPE, FD_INODE } type;
3     int ref; // reference count
4     char readable;
5     char writable;
6     struct pipe *pipe;
7     struct inode *ip;
8     uint off;
9 };
0
1
2 // in-memory copy of an inode
3 struct inode {
4     uint dev;           // Device number
5     uint inum;          // Inode number
6     int ref;            // Reference count
7     struct sleeplock lock; // protects everything below here
8     int valid;          // inode has been read from disk?
9
0     short type;         // copy of disk inode
1     short major;
2     short minor;
3     short nlink;
4     uint size;
5     uint addrs[NDIRECT+1];
6 };

```

when you read from a file, read is done through sys_read, which calls fileread

```

int
sys_read(void)
{
    struct file *f;
    int n;
    char *p;

    if(argfd(0, 0, &f) < 0 || argint(2, &n) < 0 || argptr(1, &p, n) < 0)
        return -1;
    return fileread(f, p, n);
}

int
sys_write(void)
{
    struct file *f;
    int n;

```

fileread calls piperead, if f->type was FD_PIPE

```
int
fileread(struct file *f, char *addr, int n)
{
    int r;

    if(f->readable == 0)
        return -1;
    if(f->type == FD_PIPE)
        return piperead(f->pipe, addr, n);
    if(f->type == FD_INODE){
        ilock(f->ip);
        if((r = readi(f->ip, addr, f->off, n)) > 0)
            f->off += r;
        iunlock(f->ip);
        return r;
    }
    panic("fileread");
}

//PAGEBREAK!
// Write to file f.
int
filewrite(struct file *f, char *addr, int n)
{
```

this piperead is like a queue with 2 indices that are nread and nwwrite

```

}

int
piperead(struct pipe *p, char *addr, int n)
{
    int i;

    acquire(&p->lock);
    while(p->nread == p->nwrite && p->writeopen){ //DOC: pipe-empty
        if(myproc()->killed){
            release(&p->lock);
            return -1;
        }
        sleep(&p->nread, &p->lock); //DOC: piperead-sleep
    }
    for(i = 0; i < n; i++){ //DOC: piperead-copy
        if(p->nread == p->nwrite)
            break;
        addr[i] = p->data[p->nread++ % PIPESIZE];
    }
    wakeup(&p->nwrite); //DOC: piperead-wakeup
    release(&p->lock);
    return i;
}

```

similar code for pipewrite, which is called by filewrite, which is called by sys_write

```

int
pipewrite(struct pipe *p, char *addr, int n)
{
    int i;

    acquire(&p->lock);
    for(i = 0; i < n; i++){
        while(p->nwrite == p->nread + PIPESIZE){ //DOC: pipewrite-full
            if(p->readopen == 0 || myproc()->killed){
                release(&p->lock);
                return -1;
            }
            wakeup(&p->nread);
            sleep(&p->nwrite, &p->lock); //DOC: pipewrite-sleep
        }
        p->data[p->nwrite++ % PIPESIZE] = addr[i];
    }
    wakeup(&p->nread); //DOC: pipewrite-wakeup1
    release(&p->lock);
    return n;
}

int
piperead(struct pipe *p, char *addr, int n)

```

paging

Viewing Abhijit's application

Review of last class

- MMU : Hardware features for MM
- OS: Sets up MMU for a process, then schedules process
- Compiler : Generates object code for a particular OS + MMU architecture
- MMU: Detects memory violations and raises interrupt --> Effectively passing control to OS

Review of last class

- Different types of MMUs
 - No translation (no MMU)!
 - Base + limit (also called relocation + offset) scheme
 - Leads to what type of OS ?
 - Multiple base + limit scheme --> segmentation
 - Offers which benefits?
 - Problem: fragmentation
 - Paging – logical division of the physical memory into equal sized pages, in memory page table + PTBR in hardware, OS's job and compiler's job in this case

Review of last class

- Different types of MMUs
 - No translation (no MMU)!
 - Base + limit (also called relocation + offset) scheme
 - Leads to what type of OS ?
 - Multiple base + limit scheme --> segmentation
 - Offers which benefits?
 - Problem: fragmentation
 - Paging – logical division of the physical memory into equal sized pages, in memory page table + PTBR in hardware, OS's job and compiler's job in this case

linux kernel is dynamically linked and dynamically loader

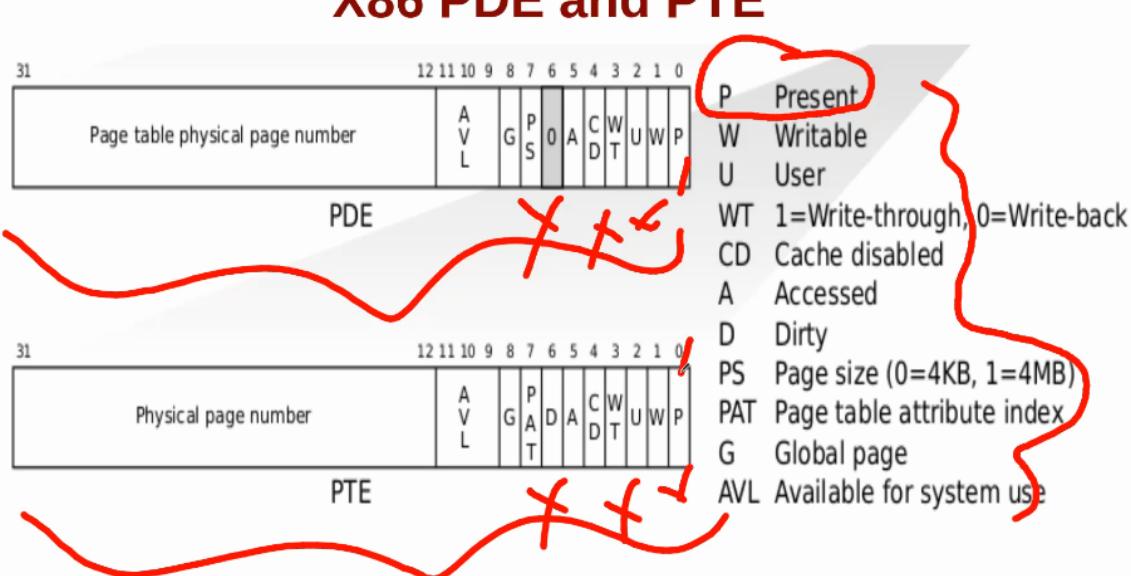
xv6 is static loading and linking

Dynamic Linking, Loading

- Dynamic linking necessarily demands an advanced type of loader that understands dynamic linking
 - Hence called ‘link-loader’
 - Static or dynamic loading is still a choice
- Question: which of the MMU options will allow for which type of linking, loading ?

process is discontinuous is the main distinction between segmentation and paging
 paging still has internal fragmentation, but the process is no longer continuous
 every process has a separate page table

X86 PDE and PTE



problem with paging is page table is continuous

Viewing Abhijit's application

Swapping

- Standard swapping
 - Entire process swapped in or swapped out
 - With continuous memory management
- Swapping with paging
 - Some pages are “paged out” and some “paged in”
 - Term “paging” refers to paging with swapping now

Virtual memory

Introduction

- Virtual memory != Virtual address

Virtual address is address issued by CPU's execution unit, later converted by MMU to physical address

Virtual memory is a memory management technique employed by OS (with hardware support, of course)

Unused parts of program

```
int a[4096][4096]
int f(int m[][4096]) {
    int i, j;
    for(i = 0; i < 1024; i++)
        m[0][i] = 200;
}
int main() {
    int i, j;
    for(i = 0; i < 1024; i++)
        a[1][i] = 200;
    if(random() == 10)
        f(a);
}
```

All parts of array a[] not accessed

Function f() may not be called

What is virtual memory?

Virtual memory - separation of user logical memory from physical memory

Only part of the program needs to be in memory for execution

Logical address space can therefore be much larger than physical address space

Allows address spaces to be shared by several processes

Allows for more efficient process creation

More programs running concurrently

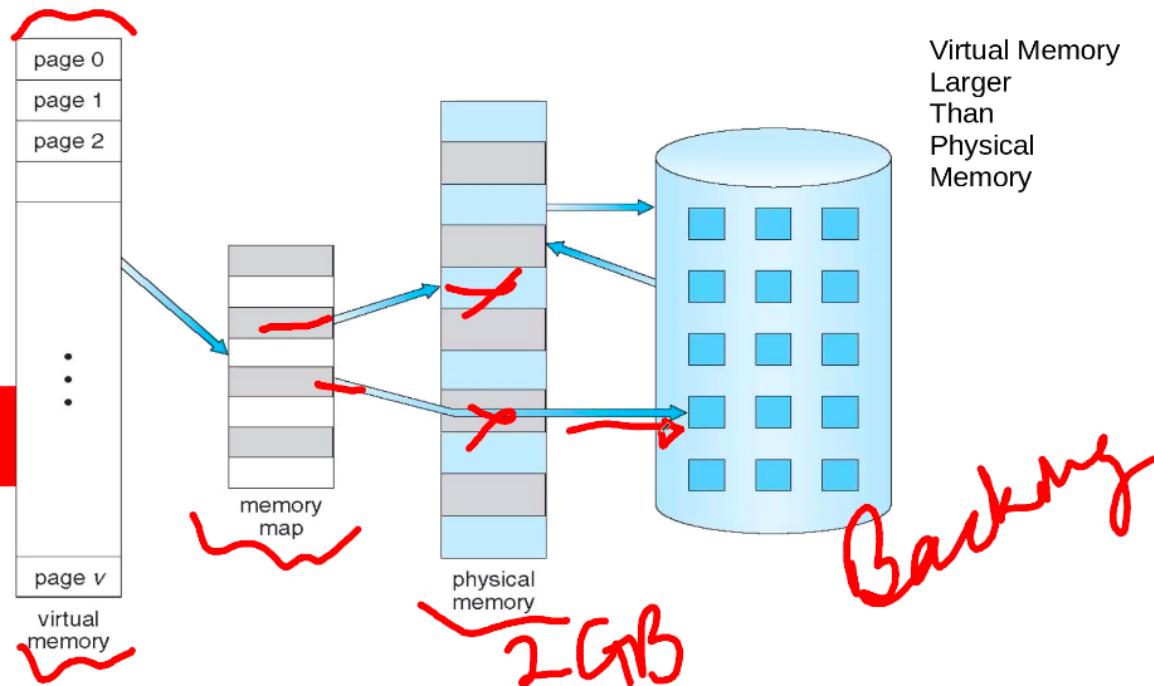
Less I/O needed to load or swap processes

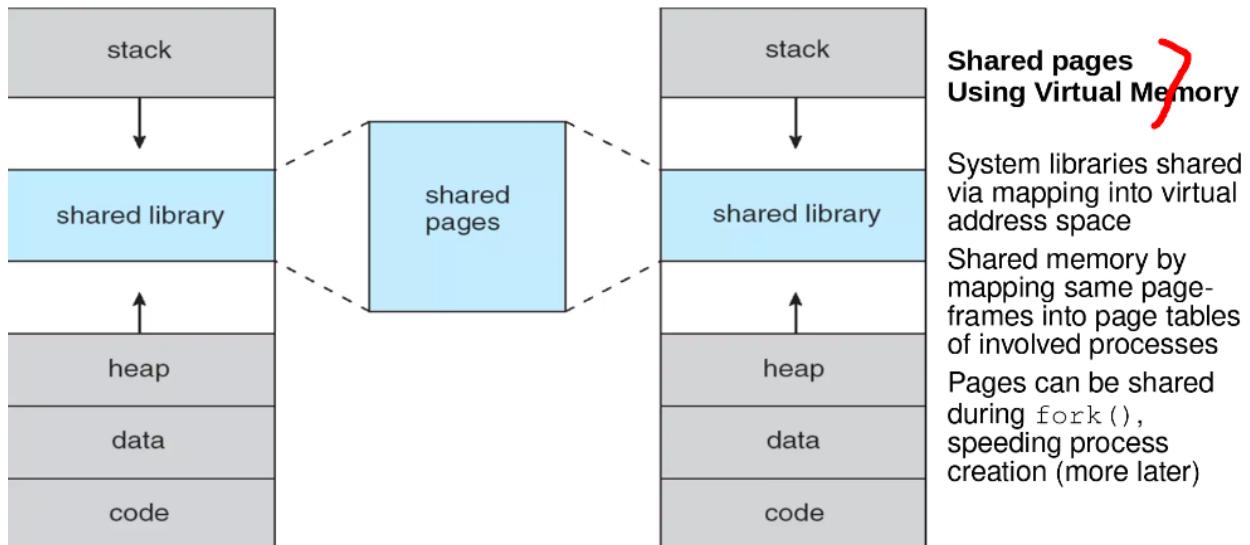
Virtual memory can be implemented via:

Demand paging

Demand segmentation

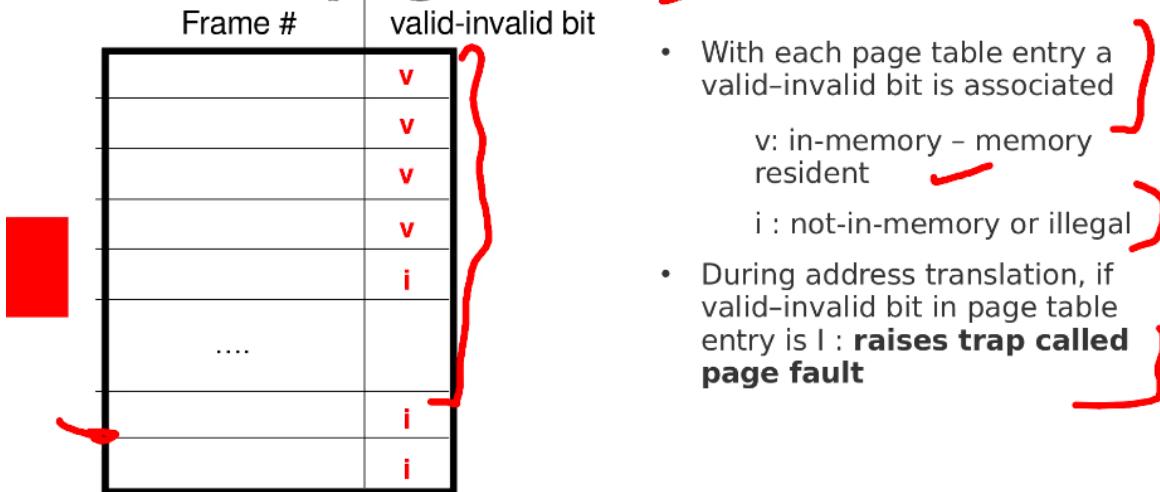
this is how Virtual memory is created, cylinder is backing store





syscall called `sbrk`, which asks the kernel to allocate some region and returns a pointer to it.

New meaning for valid/invalid bits in page table



1 instruction could cause 3 page faults, 1 for instruction, and then for mem1,mem2

Additional problems

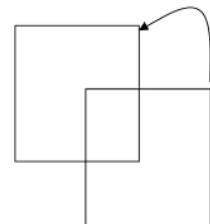
- Extreme case – start process with *no* pages in memory
OS sets instruction pointer to first instruction of process, non-memory-resident -> page fault
And for every other process pages on first access
add m₁, m₂
- Pure demand paging
- Actually, a given instruction could access multiple pages -> multiple page faults
Pain decreased because of **locality of reference**
- Hardware support needed for demand paging
Page table with valid / invalid bit
Secondary memory (swap device with **swap space**)
Instruction restart

Instruction restart

- Consider an instruction that could access several different locations

```
movarray 0x100, 0x200, 20
```

```
# copy 20 bytes from address 0x100 to  
address 0x200
```



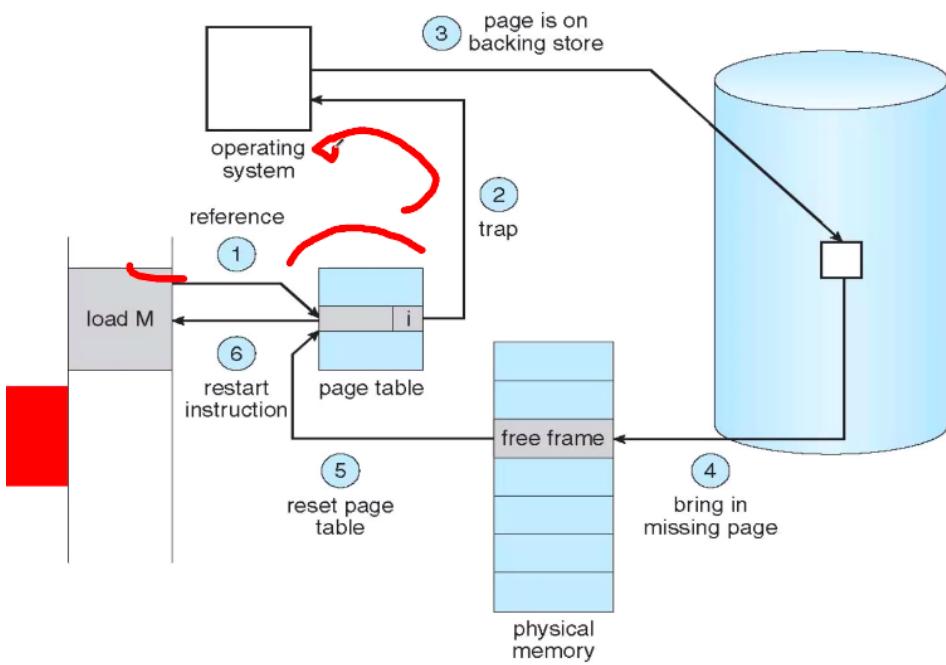
Hardware support needed for demand paging

Page table with valid / invalid bit

Secondary memory (swap device with **swap space**)

Instruction restart

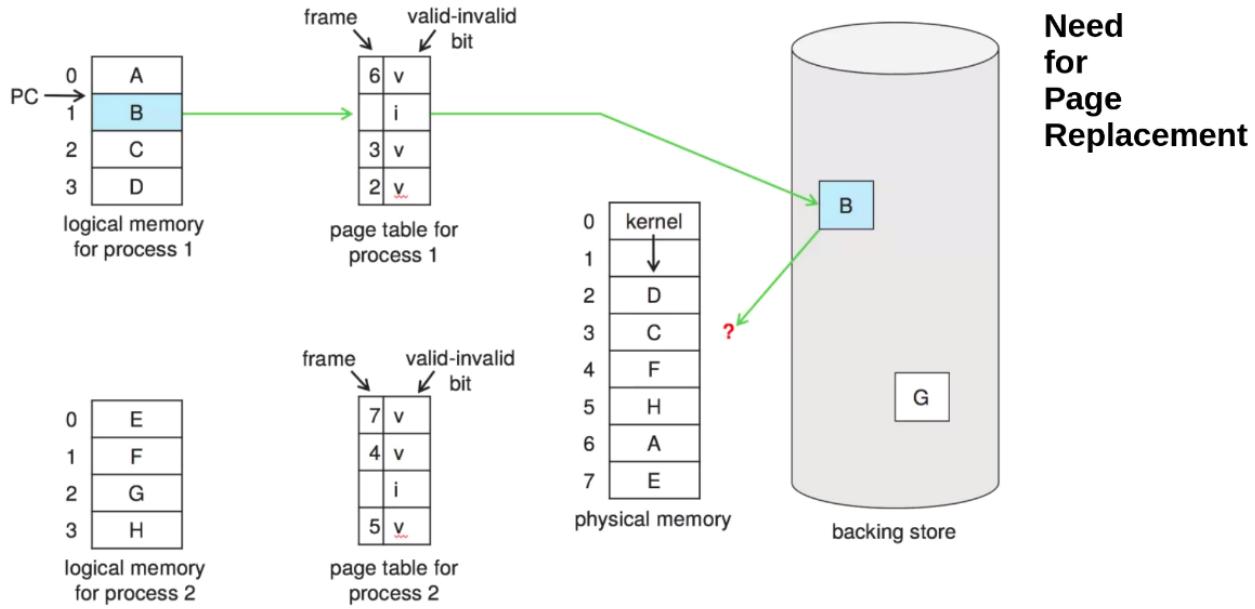
Handling A Page Fault



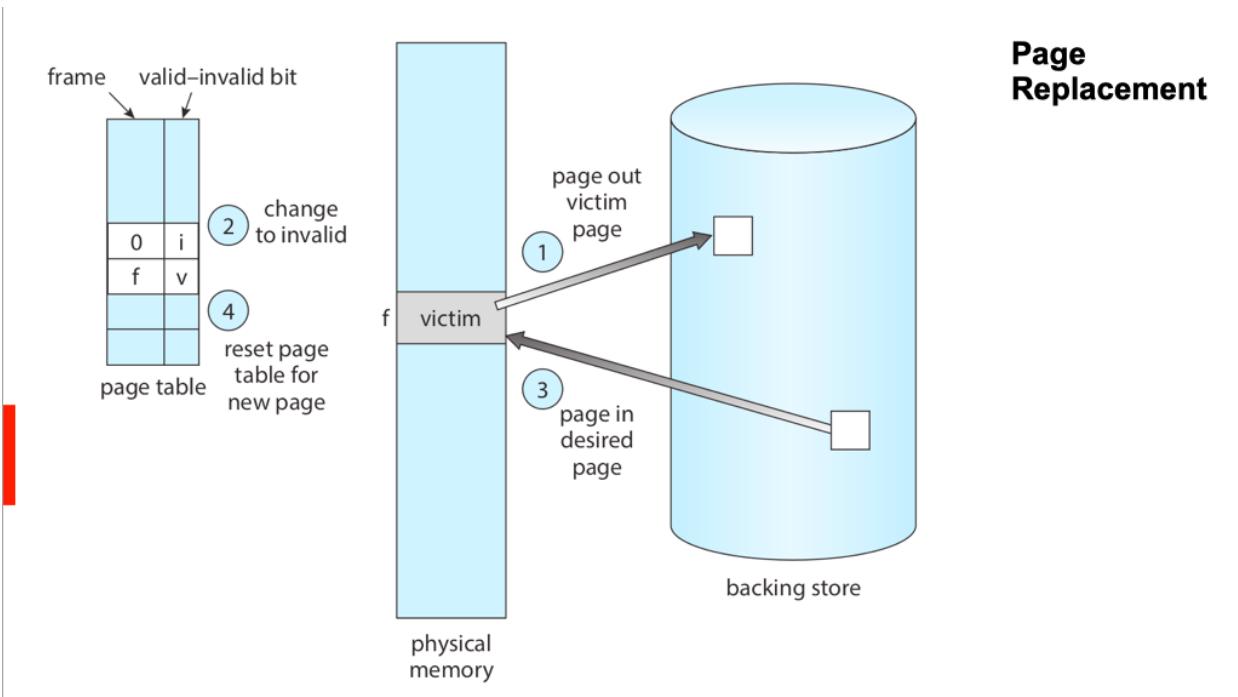
NUMBER 5 IS THE KERNEL EXECUTING,

What if no free frame found on page fault?

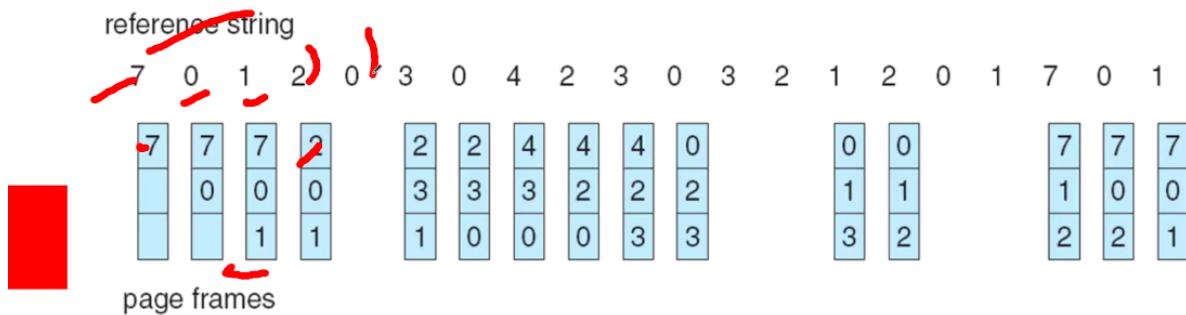
- Page frames in use depends on “Degree of multiprogramming”
 - More multiprogramming -> overallocation of frames
 - Also in demand from the kernel, I/O buffers, etc
 - How much to allocate to each process? How many processes to allow?
- Page replacement – find some page(frame) in memory, but not really in use, page it out
 - Questions : terminate process? Page out process? replace the page?
 - For performance, need an algorithm which will result in minimum number of page faults



strategies for page replacement



FIFO Algorithm



balady's anomaly

FIFO suffers with belady's anomaly

optimal and LRU doesn't

Thrashing

- If a process does not have “enough” pages, the page-fault rate is very high
 - Page fault to get page
 - Replace existing frame
 - But quickly need replaced frame back
- This leads to:
 - Low CPU utilization
 - Operating system thinking that it needs to increase the degree of multiprogramming

paging in xv6

kinit1() takes end, virtual address of 4mb [range of address i.e virtual address-end<4mb]

why virtual address? because code for this will be made by compiler which uses virtual addresses

end is defined in kernel symbol table which is the last address of ELF file

P2V will add KERNBASE,

why is this needed? because once the paging is set up, you need to use virtual address

kinit1 calls freerange, which runs a for loop on range of virtual address as shown and calls kfree for each page

```
36     freerange(vstart, vend);
37 }
38
39 void
40 kinit2(void *vstart, void *vend)
41 {
42     freerange(vstart, vend);
43     kmem.use_lock = 1;
44 }
45
46 void
47 freerange(void *vstart, void *vend)
48 {
49     char *p;
50     p = (char*)PGROUNDUP((uint)vstart);
51     for(; p + PGSIZE <= (char*)vend; p += PGSIZE)
52         kfree(p);
53 }
54 //PAGEBREAK: 21
55 // Free the page of physical memory pointed at by v,
56 // which normally should have been returned by a
57 // call to kalloc(). (The exception is when
58 // initializing the allocator; see kinit above.)
59 void
60 kfree(char *v)
```

PGROUNDUP will align address to page boundary like 4k,8k,12k

kfree() is called, v is the address of the chunk that will be part of linkedlist now

```
struct run{ struct run *next;}
```

entire xv6 code uses physical memory 0 to 224mb only

panic() runs a infinite loop after disabling interrupts, hence nothing will work after calling panic.

```
53 }
54 //PAGEBREAK: 21
55 // Free the page of physical memory pointed at by v,
56 // which normally should have been returned by a
57 // call to kalloc(). (The exception is when
58 // initializing the allocator; see kinit above.)
59 void
60 kfree(char *v)
61 {
62     struct run *r;
63
64     if((uint)v % PGSIZE || v < end || V2P(v) >= PHYSTOP)
65         panic("kfree");
66
67     // Fill with junk to catch dangling refs.
68     memset(v, 1, PGSIZE);
69
70     if(kmem.use_lock)
71         acquire(&kmem.lock);
72     r = (struct run*)v;
73     r->next = kmem.freelist;
74     kmem.freelist = r;
75     if(kmem.use_lock)
76         release(&kmem.lock);
77 }
```

memset fills the entire page with '1', virtual address v is typecasted into struct run*, then here is kmem,which has pointer call freelist which is global pointer, so line 72,73,74 creates a linkedlist

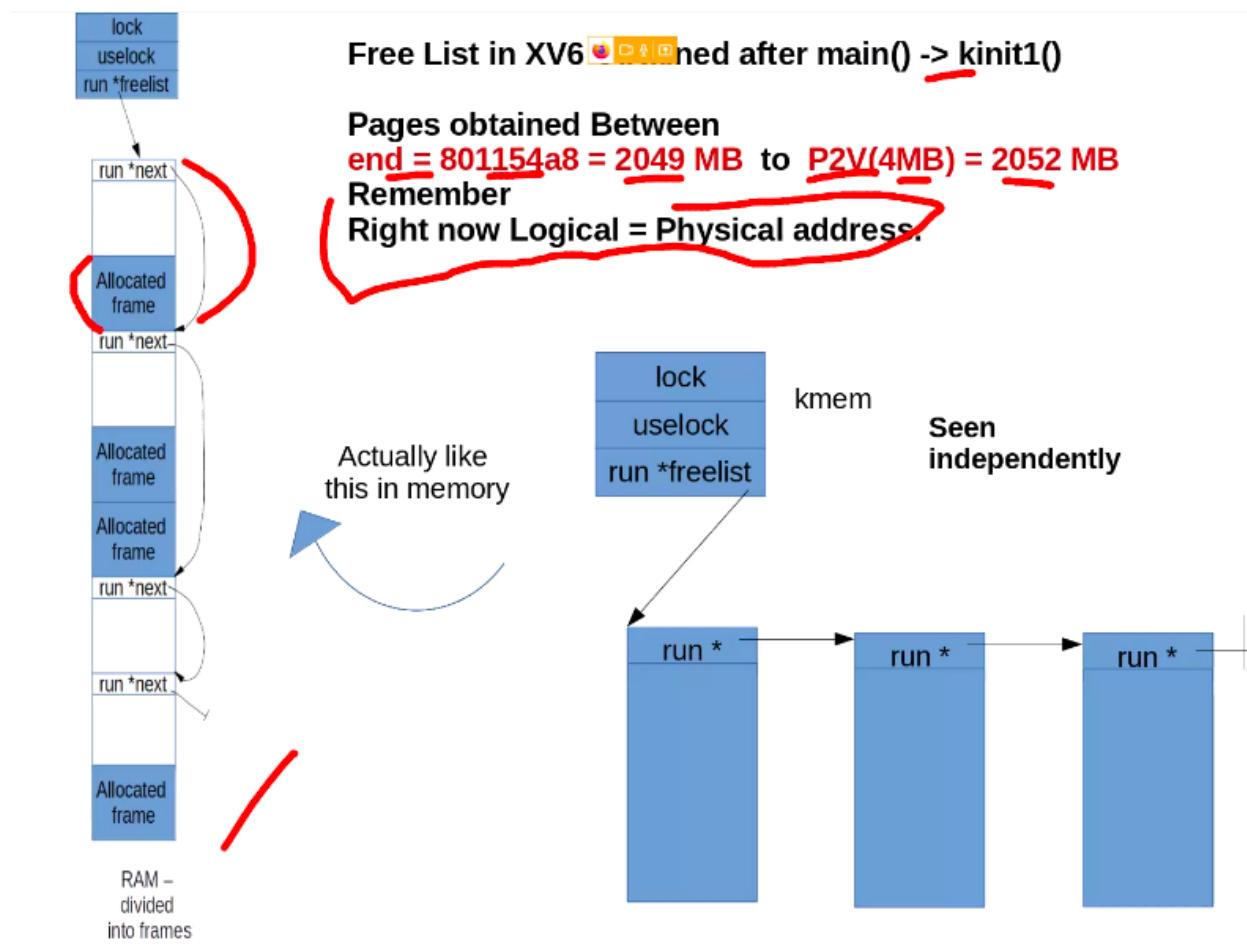
```

struct run {
    struct run *next;
};

struct {
    struct spinlock lock;
    int use_lock;
    struct run *freelist;
} kmem;

```

so $2052 - 2049 = 3\text{mb}$, which will have $3\text{mb}/4\text{kb} = 1536$ pages



we have created a free list out of 3mb.

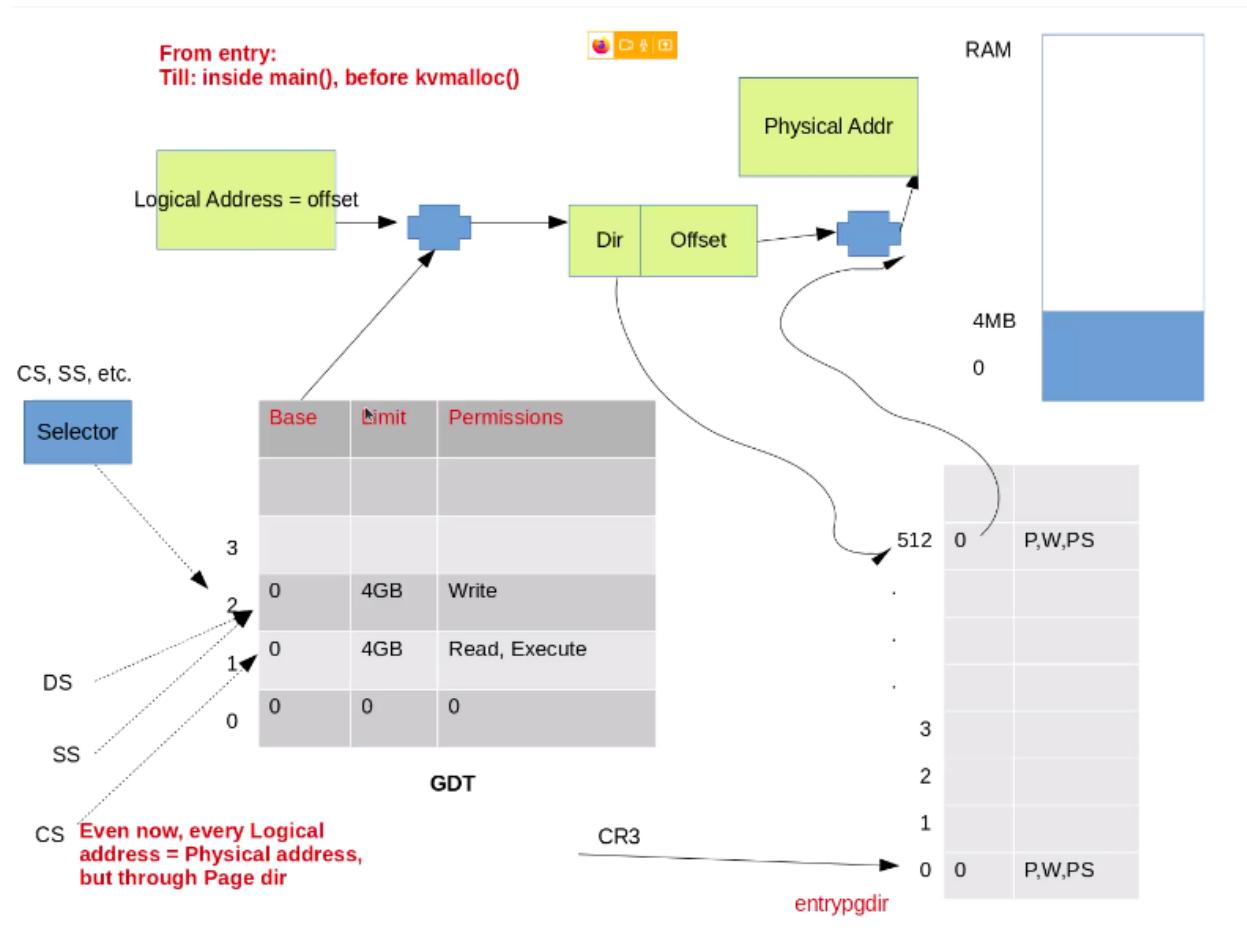
all the addresses 0 to 4mb, 2gb to 2gb+4mb will be mapped to 0 to 4mb

we also call kinit2() in bootmain()

here beginning address is 4mb, ending address is PHYSTOP

kvmalloc() will setup kernel's page table

this is the mmu before kvmalloc()



As of now only kernel is running, no applications are running.

```

7
3 // Allocate one page table for the machine for the kernel address
9 // space for scheduler processes.
9 void
1 kvmalloc(void)
2 {
3   kpgdir = setupkvm();
4   switchkvm();
5 }
5
7 // Switch h/w page table register to the kernel-only page table,
3 // for when no process is running.
9 void
9 switchkvm(void)
1 {
2   lcr3(V2P(kpgdir));    // switch to the kernel page table
3 }
n.c" 394L, 9918C

```

kvmalloc creates the mmu setup for scheduler, it calls switchkvm() lcr3() is load into cr3
 setupkvm is called from many functions.

setupkvm() will setup the page directory for the kernel and it's address is stored in
 kpgdir

kpgdir is a global uint*

```

1 // Set up kernel part of a page table.
2 pde_t*
3 setupkvm(void)
4 {
5     pde_t *pgdir;
6     struct kmap *k;
7
8     if((pgdir = (pde_t*)kalloc()) == 0)
9         return 0;
10    memset(pgdir, 0, PGSIZE);
11    if (P2V(PHYSTOP) > (void*)DEVSPACE)
12        panic("PHYSTOP too high");
13    for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
14        if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
15                    (uint)k->phys_start, k->perm) < 0) {
16            freevm(pgdir);
17            return 0;
18        }
19    return pgdir;
20 }

```

k is a pointer to struct kmap

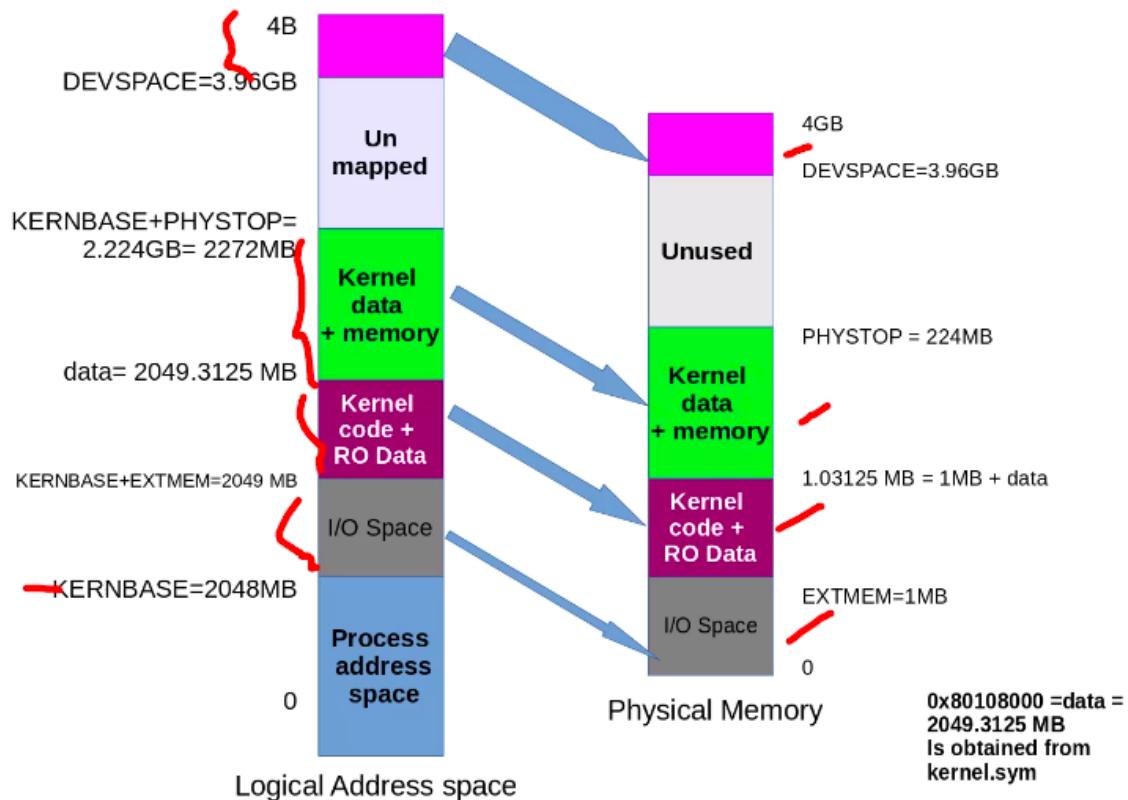
```

//                                     rw data + free physical memory
// 0xfe000000..0: mapped direct (devices such as ioapic)
//
// The kernel allocates physical memory for its heap and for user memory
// between V2P(end) and the end of physical memory (PHYSTOP)
// (directly addressable from end..P2V(PHYSTOP)).
//
// This table defines the kernel's mappings, which are present in
// every process's page table.
static struct kmap {
    void *virt;
    uint phys_start;
    uint phys_end;
    int perm;
} kmap[] = {
{ (void*)KERNBASE, 0,           EXTMEM,      PTE_W}, // I/O space
{ (void*)KERNLINK, V2P(KERNLINK), V2P(data), 0}, // kern text+rodata
{ (void*)data,     V2P(data),    PHYSTOP,    PTE_W}, // kern data+mem
{ (void*)DEVSPACE, DEVSPACE,    0,           PTE_W}, // more devices
};

```

it is array of structures

kmap[] mappings done in kvmalloc(). This shows segmentwise, entries are done in page directory and page table for corresponding VA > PA mappings



EXTMEM is 1MB, which is mapped to 0 to 1MB in physical space

0 as permission is read permission

kalloc() is counterpart of kfree

```

// Allocate one 4096-byte page of physical memory.
// Returns a pointer that the kernel can use.
// Returns 0 if the memory cannot be allocated.
char*
kalloc(void)
{
    struct run *r;

    if(kmem.use_lock)
        acquire(&kmem.lock);
    r = kmem.freelist;
    if(r)
        kmem.freelist = r->next;
    if(kmem.use_lock)
        release(&kmem.lock);
    return (char*)r;
}

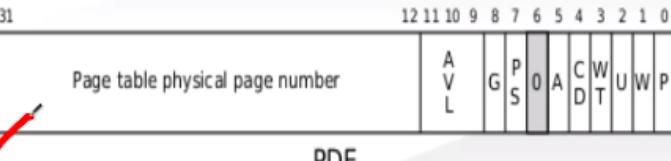
```

it picks up the first frame on freelist and returns a pointer to it. kalloc() will give you a 4K frame

pgdir points to frame returned by kalloc()

we call memset() which sets 0 in all pgdir, so we make P=0,PS=0 i.e present entry=0, PS=0 i.e. 4KB page size

Page Directory Entry (PDE) Page Table Entry (PTE)



P Present

W Writable

U User

WT 1=Write-through, 0=Write-back

CD Cache disabled

A Accessed

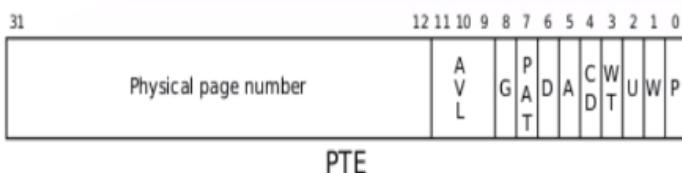
D Dirty

PS Page size (0=4KB, 1=4MB)

PAT Page table attribute index

G Global page

AVL Available for system use



for loop iterates through all kmap entries and calls mappages() which creates PDE and PTE entries, allocate page tables, and setup entries in PT and PD

```

    return &pgtab[PTX(va)];
}

// Create PTEs for virtual addresses starting at va that refer to
// physical addresses starting at pa. va and size might not
// be page-aligned.
static int
mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
{
    char *a, *last;
    pte_t *pte;

    a = (char*)PGROUNDDOWN((uint)va);
    last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
    for(;;){
        if((pte = walkpgdir(pgdir, a, 1)) == 0)
            return -1;
        if(*pte & PTE_P)
            panic("remap");
        *pte = pa | perm | PTE_P;
        if(a == last)
            break;
        a += PGSIZE;
        pa += PGSIZE;
    }
.c" 394L. 9918C

```

66.1

inside mappages() it calls walkpgdir() for all pagedir entries

```

// Return the address of the PTE in page table pgdir
// that corresponds to virtual address va. If alloc!=0,
// create any required page table pages.
static pte_t *
walkpgdir(pde_t *pgdir, const void *va, int alloc)
{
    pde_t *pde;
    pte_t *pgtab;

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

```

41.17

```

// A virtual address 'la' has a three-part structure as follows:
//
// +-----+-----+-----+
// | Page Directory | Page Table | Offset within Page |
// |     Index      |     Index   |                         |
// +-----+-----+-----+
// \--- PDX(va) ---/ \--- PTX(va) ---/

// page directory index
#define PDX(va) ((uint)(va) >> PDXSHIFT) & 0x3FF

// page table index
#define PTX(va) ((uint)(va) >> PTXSHIFT) & 0x3FF

// construct virtual address from indexes and offset

```

PDX will give page directory index and PTX will give page table index

if (Present bit is 1) which is false as we memset to 0

if this was true, then you will get a page number

but in else part we allocate a 4kb page [using kalloc]

we set up a page directory entry with physical address of pgtable, with flags returns pointer to page table entry

then we return in mappages, which fills a entry in the page table [as *pde] with physical address pa passed as argument[page number]

which returns in setupkvm, and thus 4 entries of kernel are setup.

summary

walkpgdir checks if va is mapped inside the pgdir entry, if yes return entry, no then allocate a page table then return address to the entry in pgdir of that page table

kvmalloc switches from 1 level paging to 2 level paging.

kvmalloc → setupkvm → mappages → walkpgdir

if alloc is 1, walkpgdir allocates a page table and returns its address

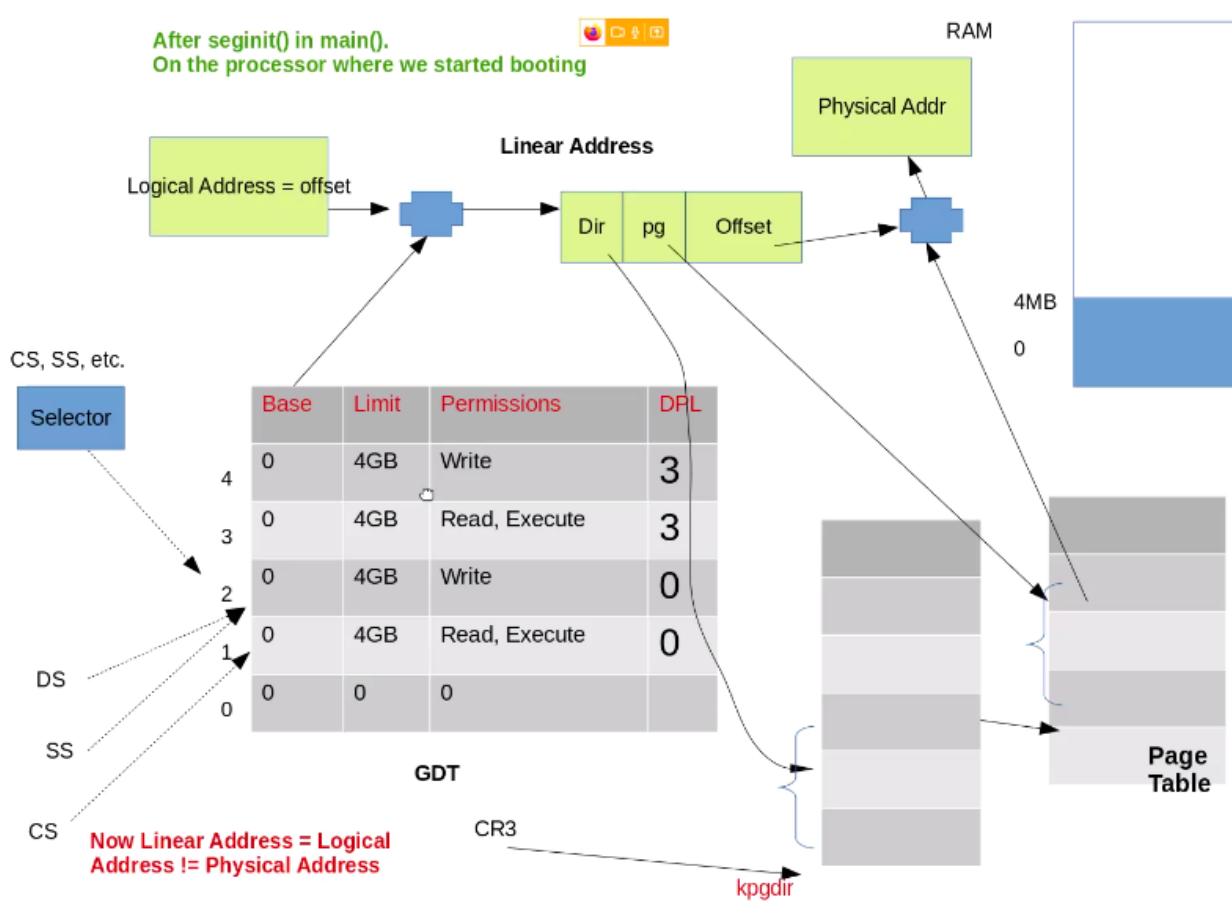
mappages fills the entry returned by walkpgdir with physical address pa ,permissions and sets Present bit.

setupkvm runs mappages for all the 4 entries of kmap.

kvmalloc sets the address of pgdir returned by setupkvm into global variable called kpgdir and loads the same address into the cr3 register as well.

```
// Set up CPU's kernel segment descriptors.  
// Run once on entry on each CPU.  
void  
seginit(void)  
{  
    struct cpu *c;  
  
    // Map "logical" addresses to virtual addresses using identity map.  
    // Cannot share a CODE descriptor for both kernel and user  
    // because it would have to have DPL_USER, but the CPU forbids  
    // an interrupt from CPL=0 to CPL=3.  
    c = &cpus[cpid()];  
    c->gdt[SEG_KCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, 0);  
    c->gdt[SEG_KDATA] = SEG(STA_W, 0, 0xffffffff, 0);  
    c->gdt[SEG_UCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, DPL_USER);  
    c->gdt[SEG_UDATA] = SEG(STA_W, 0, 0xffffffff, DPL_USER);  
    lgdt(c->gdt, sizeof(c->gdt));  
}
```

seginit will change segmentation table, now there are 4 entries



when application is running you ensure CS,DS is 3 and 4, for kernel is 1 and 2

- **Thread-local storage (TLS)** allows each thread to have its own copy of data
- Useful when you do not have control over the thread creation process (i.e., when using a thread pool)
- Different from local variables
 - Local variables visible only during single function invocation
 - TLS visible across function invocations
- Similar to **static** data
 - TLS is unique to each thread

Scheduler activations for threads

- Both M:M and Two-level models require communication to maintain the appropriate number of kernel threads allocated to the application
- Typically use an intermediate data structure between user and kernel threads – **lightweight process (LWP)**
 - Appears to be a virtual processor on which process can schedule user thread to run
 - Each LWP attached to kernel thread
 - How many LWPs to create?
- Scheduler activations provide **upcalls** - a communication mechanism from the kernel to the **upcall handler** in the thread library
- This communication allows an application to maintain the correct number kernel threads

Operating systems provide system calls like kill() and signal() to enable processes to deliver and receive signals

Signal() - is used by a process to specify a “signal handler” – a code that should run on receiving a signal

Kill() is used by a process to send another process a signal