

---

# 2'nd Lab Project

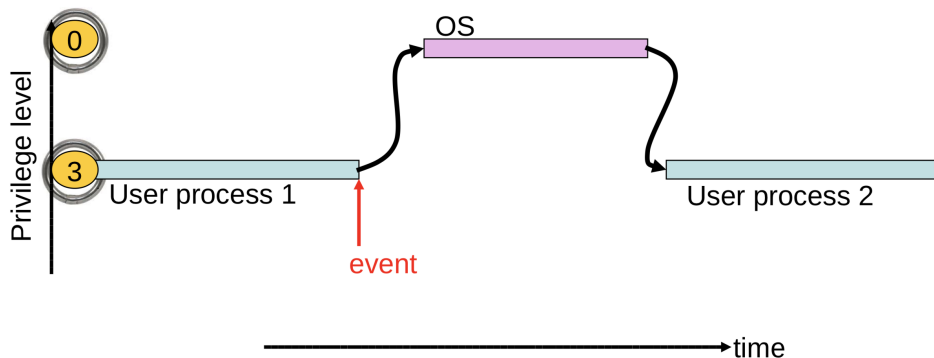
## System Calls and Processes

Omid Bodaghi - Saber EbrahimNejad

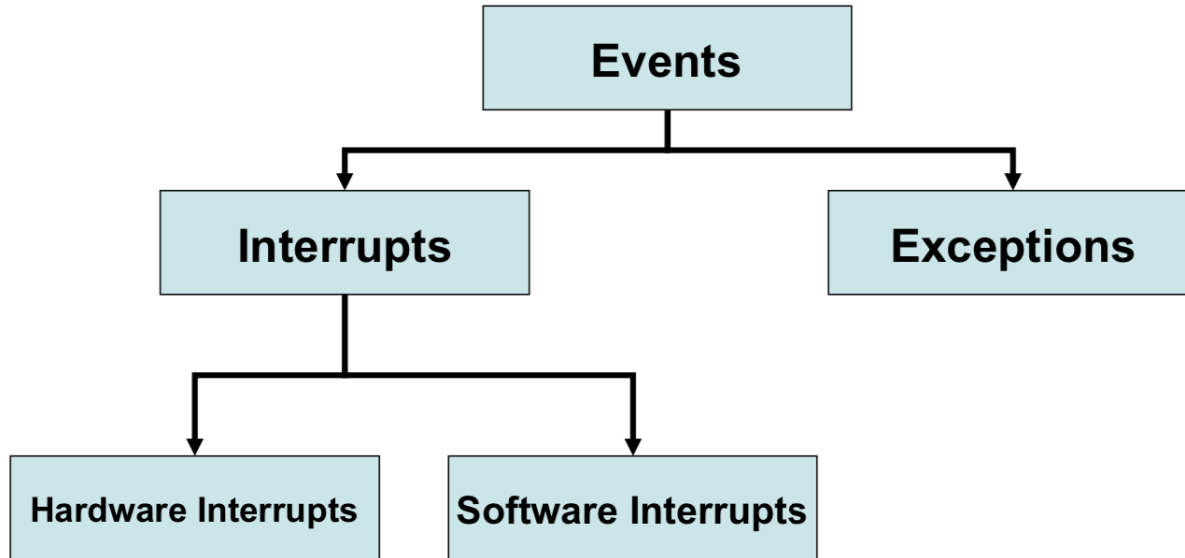
---

# Why event driven design?

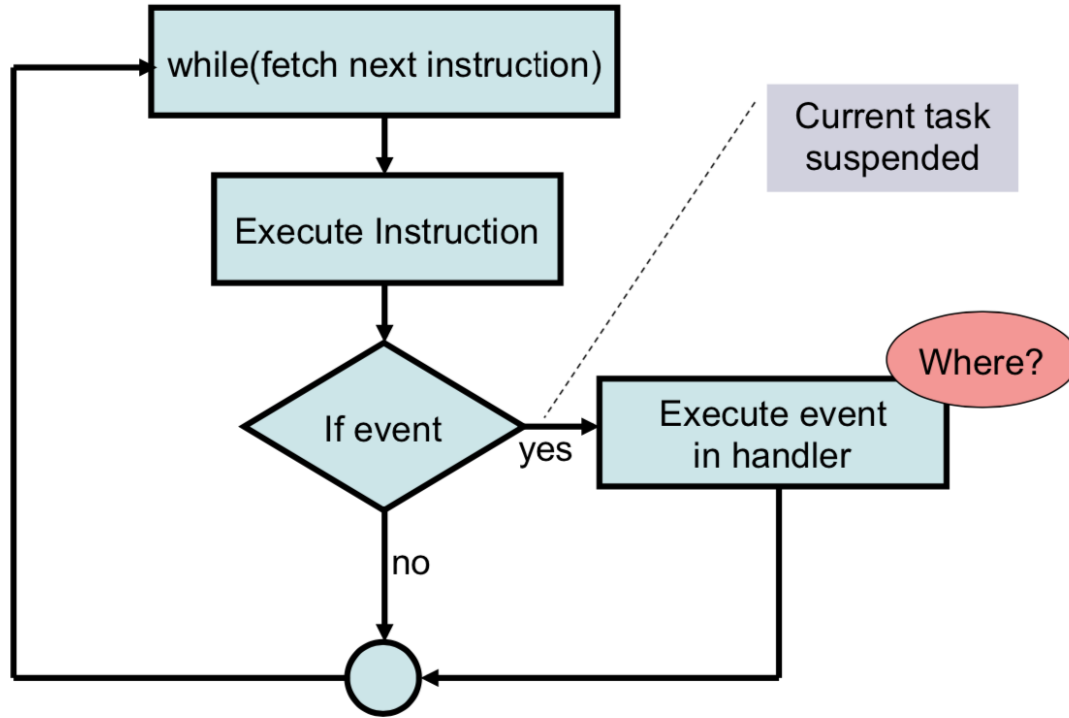
- OS cannot **trust** user processes
  - User processes may be buggy or malicious
  - User process crash should not affect OS
- OS needs to guarantee **fairness** to all user processes
  - One process cannot 'hog' CPU time
  - Timer interrupts



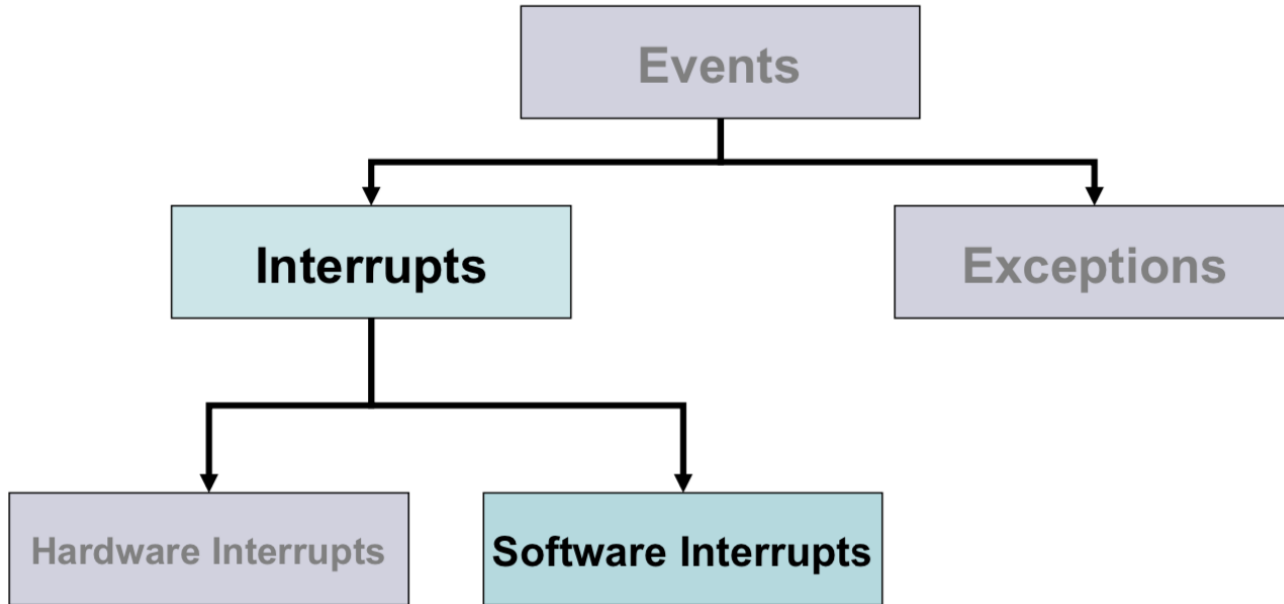
# Event Types



# Event View of CPU

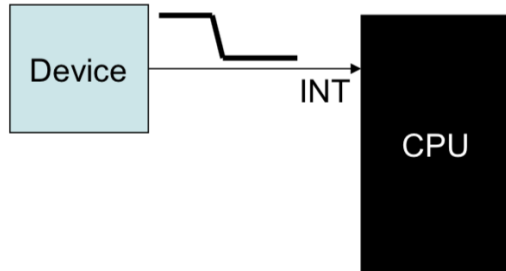


# System Calls



# Hardware vs Software Interrupt

## Hardware Interrupt



A device (like PIC)  
asserts a pin in the CPU

## Software Interrupt

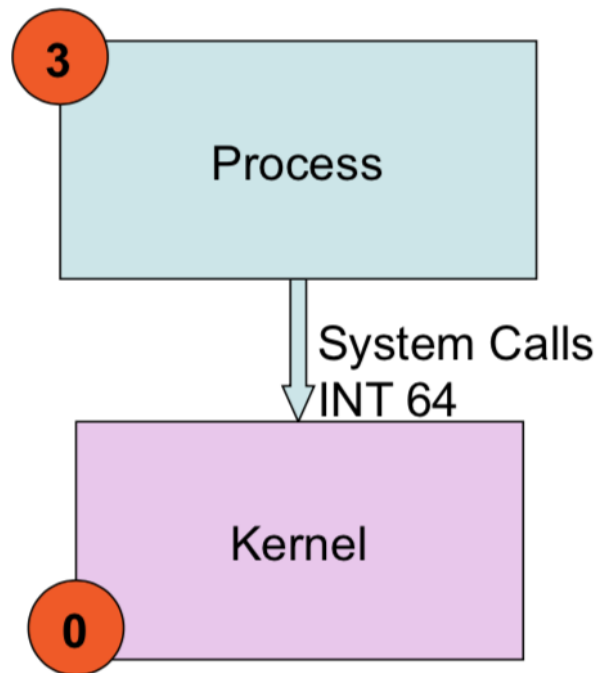


An instruction which when executed  
causes an interrupt

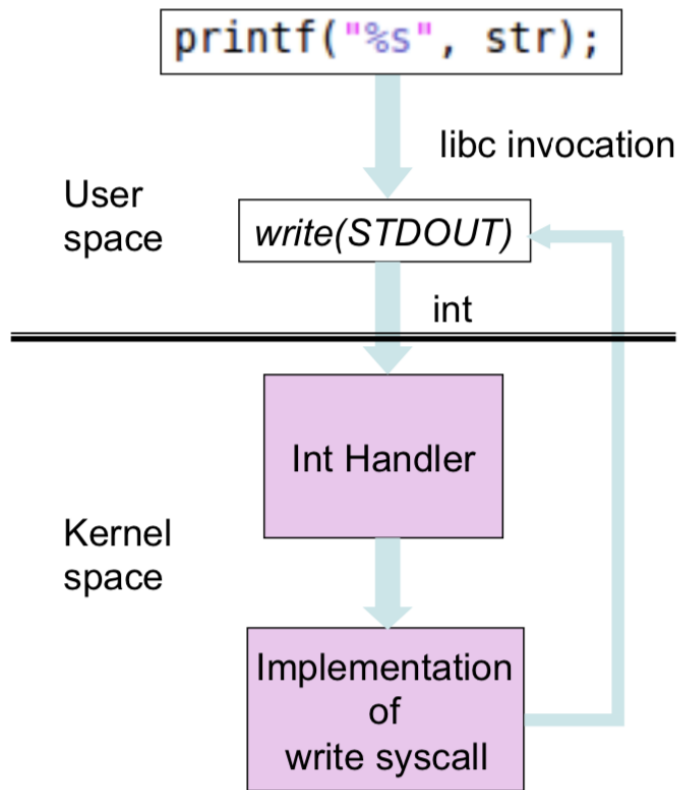
# Software Interrupt

Software interrupt used for implementing system calls

- Previously in Linux INT 128, was used for system calls
- In xv6, INT 64 is used for system calls



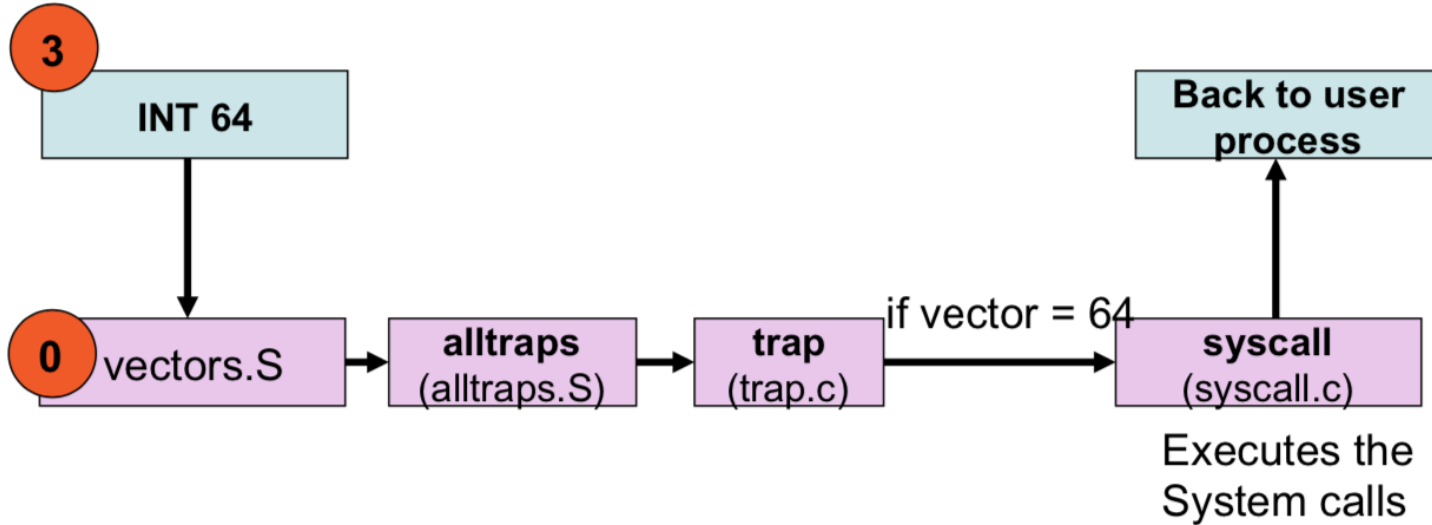
## Example (write system call)



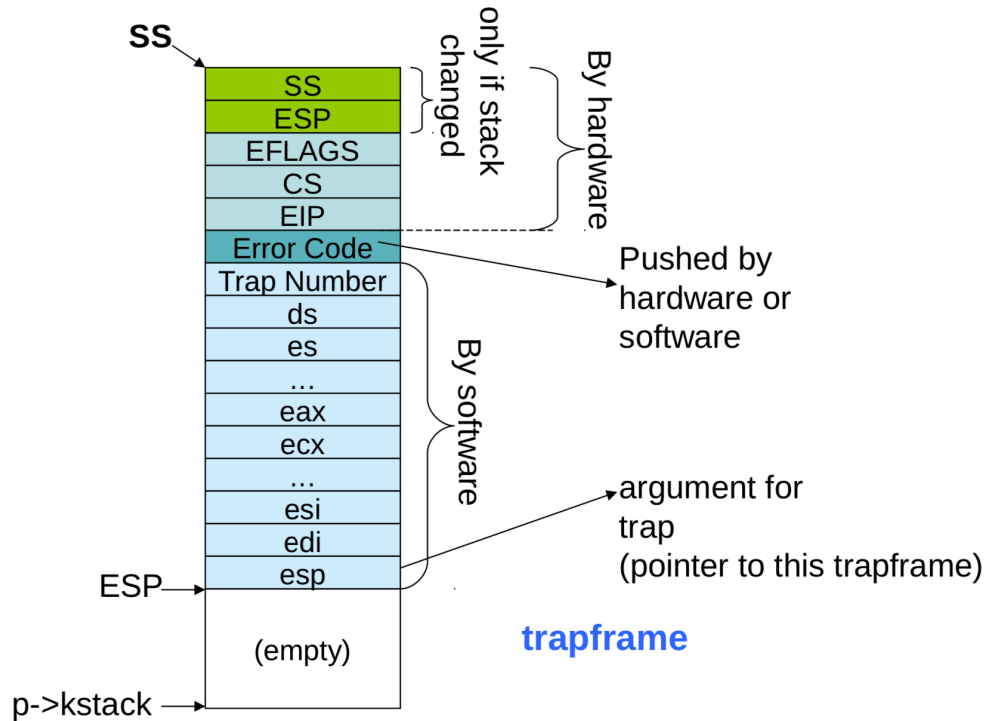


# System Call Processing in kernel

Almost similar to hardware interrupts



# Trapframe



ref : struct trapframe in x86.h (0602 [06])

# Trapframe Struct

```
0602 struct trapframe {
0603     // registers as pushed by pusha
0604     uint edi;
0605     uint esi;
0606     uint ebp;
0607     uint oesp;    // useless & ignored
0608     uint ebx;
0609     uint edx;
0610     uint ecx;
0611     uint eax;
0612
0613     // rest of trap frame
0614     ushort gs;
0615     ushort padding1;
0616     ushort fs;
0617     ushort padding2;
0618     ushort es;
0619     ushort padding3;
0620     ushort ds;
0621     ushort padding4;
0622     uint trapno;
0623
0624     // below here defined by x86 hardware
0625     uint err;
0626     uint eip;
0627     ushort cs;
0628     ushort padding5;
0629     uint eflags;
0630
0631     // below here only when crossing rings, such as from user to kernel
0632     uint esp;
0633     ushort ss;
0634     ushort padding6;
0635 };
```

SS
ESP
EFLAGS
CS
EIP
Error Code
Trap Number
ds
es
...
eax
ecx
...
esi
edi
esp
(empty)

# System Calls in xv6

System call	Description
fork()	Create process
exit()	Terminate current process
wait()	Wait for a child process to exit
kill(pid)	Terminate process pid
getpid()	Return current process's id
sleep(n)	Sleep for n seconds
exec(filename, *argv)	Load a file and execute it
sbrk(n)	Grow process's memory by n bytes
open(filename, flags)	Open a file; flags indicate read/write
read(fd, buf, n)	Read n bytes from an open file into buf
write(fd, buf, n)	Write n bytes to an open file
close(fd)	Release open file fd
dup(fd)	Duplicate fd
pipe(p)	Create a pipe and return fd's in p
chdir(dirname)	Change the current directory
mkdir(dirname)	Create a new directory
mknod(name, major, minor)	Create a device file
fstat(fd)	Return info about an open file
link(f1, f2)	Create another name (f2) for the file f1
unlink(filename)	Remove a file

**How does the  
OS distinguish  
between the system  
calls?**

# System Call Number

System  
call number

```
mov x, %eax  
INT 64
```

Based on the system call number  
function syscall invokes the  
corresponding syscall handler

## System call numbers

```
#define SYS_fork      1  
#define SYS_exit      2  
#define SYS_wait      3  
#define SYS_pipe      4  
#define SYS_read      5  
#define SYS_kill      6  
#define SYS_exec      7  
#define SYS_fstat     8  
#define SYS_chdir     9  
#define SYS_dup      10  
#define SYS_getpid    11  
#define SYS_sbrk      12  
#define SYS_sleep     13  
#define SYS_uptime    14  
#define SYS_open      15  
#define SYS_write     16  
#define SYS_mknod     17  
#define SYS_unlink    18  
#define SYS_link      19  
#define SYS_mkdir     20  
#define SYS_close     21
```

## System call handlers

```
[SYS_fork]    sys_fork,  
[SYS_exit]    sys_exit,  
[SYS_wait]    sys_wait,  
[SYS_pipe]    sys_pipe,  
[SYS_read]    sys_read,  
[SYS_kill]    sys_kill,  
[SYS_exec]    sys_exec,  
[SYS_fstat]   sys_fstat,  
[SYS_chdir]   sys_chdir,  
[SYS_dup]     sys_dup,  
[SYS_getpid]  sys_getpid,  
[SYS_sbrk]    sys_sbrk,  
[SYS_sleep]   sys_sleep,  
[SYS_uptime]  sys_uptime,  
[SYS_open]    sys_open,  
[SYS_write]   sys_write,  
[SYS_mknod]   sys_mknod,  
[SYS_unlink]  sys_unlink,  
[SYS_link]    sys_link,  
[SYS_mkdir]   sys_mkdir,  
[SYS_close]   sys_close,
```

ref : syscall.h, syscall() in syscall.c

# xv6 System Call Naming Convention

- Usually a library function `foo()` will do some work and then call a system call `sys_foo()`
  - `sys_foo()` implemented in `sys*.c` (`sysfile.c`, `sysproc.c`)
- System call number for `foo()` is `SYS_foo`
  - `syscalls.h`
- All system calls begin with `sys_`

## Syscall(void)

All system calls are handled in this function.

The sys num which is saved in eax register is retrieved and system call is read from table.

```
void
syscall(void)
{
    int num;

    num = proc->tf->eax;
    if(num >= 0 && num < NELEM(syscalls) && syscalls[num])
        proc->tf->eax = syscalls[num]();
    else {
        cprintf("%d %s: unknown sys call %d\n",
                proc->pid, proc->name, num);
        proc->tf->eax = -1;
    }
}
```

ref : syscall.h, syscall() in syscall.c

# Prototype of a Typical System Call

```
int system_call( resource_descriptor, parameters)
```



```
graph TD; A["int system_call( resource_descriptor, parameters)"] --> B["return is generally 'int' (or equivalent) sometimes 'void'"]; A --> C["What OS resource is the target here? For example a file, device, etc. If not specified, generally means the current process"]; A --> D["System call specific parameters passed. How are they passed?"]
```

return is generally  
'int' (or equivalent)  
sometimes 'void'

int used to denote completion  
status of system call sometimes  
also has additional information  
like number of bytes written to  
file

What OS resource is the target  
here?  
For example a file, device, etc.

If not specified, generally means  
the current process

System call specific parameters  
passed.  
How are they passed?



# Adding New System Call

- A system call body is defined in sysproc.c or sysfile.c
- Multiple files are needed to be altered to add a new syscall.
- Desired new system calls in this project:
  - void calculate\_biggest\_perfect\_square (int n)
  - void set\_sleep(int seconds)
  - void process\_start\_time(void)
  - void get\_ancestors(int pid)
  - void get\_descendant(int pid)

# Processes

A process is the running of a program, including the program's state and data. The state includes such things as:

- Memory the program occupied
- Memory contents
- Register values
- Files
- Kernel structures

# Managing Processes

The kernel has a simple data structure for each process, organized in some list. The kernel juggles between the processes, using a context switch, which:

- Saves state of old process to memory.
- Loads state of new process from memory

The processes are held in a struct named `ptable`, who has a vector of processes named `proc`.

# First Process

Processes are created by the kernel, after another process asks it to. Therefore, the kernel needs to run the first process itself, in order to create someone who will ask for new processes to be created.

```
// Per-process state
struct proc {
    uint sz;                // Size of process memory (bytes)
    pde_t* pgdir;           // Page table
    char *kstack;           // Bottom of kernel stack for this process
    enum procstate state;   // Process state
    volatile int pid;       // Process ID
    struct proc *parent;    // Parent process
    struct trapframe *tf;   // Trap frame for current syscall
    struct context *context; // switch() here to run process
    void *chan;             // If non-zero, sleeping on chan
    int killed;             // If non-zero, have been killed
    struct file *ofile[NOFILE]; // Open files
    struct inode *cwd;       // Current directory
    char name[16];          // Process name (debugging)
};
```

Proc structure Ref: proc.h

# Passing Parameters in System Calls

- Passing parameters to system calls not similar to passing parameters in function calls.
  - Recall stack changes from user mode stack to kernelstack.
- Typical Methods
  - Pass by Registers (eg. Linux)
  - Pass via user mode stack (eg. xv6)
    - Complex
  - Pass via a designated memory region
    - Address passed through registers

## Pass By Registers(Linux)

- System calls with fewer than 6 parameters passed in registers
  - %eax (sys call number), %ebx, %ecx,, %esi, %edi, %ebp
- If 6 or more arguments
  - Pass pointer to block structure containing argument list
- Max size of argument is the register size (eg. 32 bit)
  - Larger pointers passed through pointers

# Pass via User Mode Stack(xv6)

## User process

```
push param1
push param2
push param3
mov sysnum, %eax
int 64
```

## User stack

param1
param2
param3

## trapframe

SS
ESP
EFLAGS
CS
EIP
Error Code
Trap Number
ds
es
...
eax
ecx
...
esi
edi
ESP
(empty)

ESP pushed by hardware contains user mode stack pointer

proc entry  
for process

Points to trapframe

ref : sys\_open (sysfile.c), argint, fetchint (syscall.c)

# Returns from System Calls

## User process

```
push param1
push param2
push param3
mov sysnum, %eax
int 64
.....
```

## Return value

register EAX

## trapframe

SS
ESP
EFLAGS
CS
EIP
Error Code
Trap Number
ds
es
...
eax
ecx
...
esi
edi
ESP
(empty)

## in system call

move result to eax in  
trap frame

Automatically restored  
by hardware while returning  
to user process



# Calculate biggest perfect square

- Int calculate\_biggest\_perfect\_square(int num)
  - Get a number as argument
  - Return biggest perfect square which is less than num
  - Use a register to store the argument

# Sleep system call

- `void set_sleep(int seconds)`
  - Get a number as argument
  - Wait until the seconds pass
  - Busy waiting
  - (Hint: each 100 amount of “ticks” variable = 1 second)

# Process start time

- `int process_start_time(void)`
  - Returns the start time of the current process based on **ticks**
  - knowledge about the pcb and process creation in XV6 is required
  - You should save the process creation time in pcb

# Print Ancestors

- `void get_ancestors(pid)`
  - Get pid as an argument
  - Print ancestors of pid
  - Write user program

# Print Descendants

- `void get_descendants(pid)`
  - Get pid as an argument
  - Print descendants of pid
  - Select youngest child
  - Write user program

# Deadline

- 29 Farvardin
- Provide Report
- Resources:
  - [http://www.cse.iitm.ac.in/~chester/courses/16o\\_os/slides/6\\_Interrupts.pdf](http://www.cse.iitm.ac.in/~chester/courses/16o_os/slides/6_Interrupts.pdf)