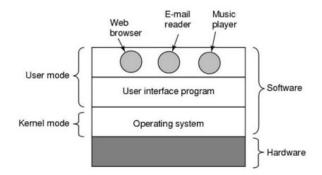# January 2025
# CSE 314 - OS Sessional

Xv6 Introduction
Nafis Tahmid
May 07, 2025

# Dual Mode Operation

CPU executes in 2 Modes

- **Kernel mode:** CPU can execute all machine instructions CPU can use every hardware feature
- **User mode:** permits only a subset of the instructions and a subset of the hardware features. Allows OS to protect itself and other system components.

# Kernel Mode Execution

- When does CPU start executing in kernel mode?
    A. System Boot - starting a computer
    B. Hardware Interrupt - generated by hardware devices to signal that they need some attention from the OS.
    C. Trap ←————————————
        A. A software-generated interrupt caused either by an error (i.e. division by 0 or invalid memory access) or
        B. By a specific request from a user program that an operating-system service needs to be performed.

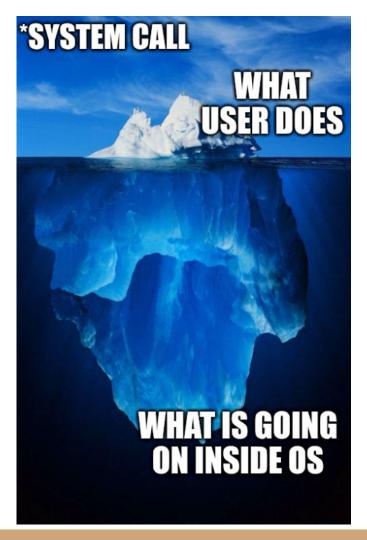# Switching Modes

- To obtain services from the operating system,
  - an user program must make a *system call*, which *traps* into the kernel and invokes the operating system.

- The **TRAP** instruction switches from user mode to kernel mode and starts the operating system.

- When the work has been completed, control is returned to the user program at the instruction following the system call.

# System Calls

- Programming interface to the services provided by the OS

- Typically used from a high-level language (C or C++) program.

- Typically a number is associated with each system call, and the OS maintains a table indexed according to these numbers.

# Let's see xv6

Prerequisite tools:
https://pdos.csail.mit.edu/6.828/2022/tools.html

Press Ctrl+A then X to quit xv6

# Challenges you will face in this assignment

Too many files

Use

- grep <pattern> *
- grep <pattern> kernel/*
- grep <pattern> user/*

When you are the only
language people choose
to write a kernel*..

Signature look of
superiority

*Windows, Linux & Unix(Mac)

C

yuva.krishna.memes

# Warnings

me

xv6

# Adding a User program

- Create a file under user folder
- These three lines must be added on top strictly following order

```
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"
```

# Adding a User program

- Create a file under user folder
- These three lines must be added on top **strictly following order**

```
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"
```

User can call these functions **"declared"** here

# Adding a User program

- Create a file under user folder
- These three lines must be added on top strictly following order

```
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"
```

- Then add prog_name in Makefile.
- Then run make clean; make qemu

```
UPROGS=\
    $U/_cat\
    $U/_echo\
    $U/_forktest\
    $U/_grep\
    $U/_init\
    $U/_kill\
    $U/_ln\
    $U/_ls\
    $U/_mkdir\
    $U/_myprog\ # add this line
    $U/_rm\
    $U/_sh\
    $U/_stressfs\
    $U/_usertests\
    $U/_grind\
    $U/_wc\
    $U/_zombie\
```

# Adding a system call

First let's understand what happens during a system call.

Suppose we have added **our own** system call with signature **int getuid()**; which simply returns an **integer value** stored in a variable named **uid** in **kernel.**

gibuid.c

```
int uid = getuid();

printf("%d\n", uid);
```

# A series of things happen

user/gibuid.c

```c
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"

int main()
{
    int uid = getuid();

    printf("%d\n", uid);

    return 0;
}
```

# A series of things happen

user/user.h

```c
C  user.h > ...
  struct stat;

  // system calls
  int fork(void);
  int exit(int) __attribute__((noreturn));
  int wait(int*);
  int pipe(int*);
  int write(int, const void*, int);
  int read(int, void*, int);
  int close(int);
  int kill(int);
  int exec(const char*, char**);
  int open(const char*, int);
  int mknod(const char*, short, short);
  int unlink(const char*);
  int fstat(int fd, struct stat*);
  int link(const char*, const char*);
  int mkdir(const char*);
  int chdir(const char*);
  int dup(int);
  int getpid(void);
  char* sbrk(int);
  int sleep(int);
  int uptime(void);
  int getuid(void);   // signature declaration
```

user/gibuid.c

```c
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"

int main()
{
    int uid = getuid();

    printf("%d\n", uid);

    return 0;
}
```

Signature is **declared** here

# A series of things happen

user/gibuid.c

```c
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"

int main()
{
    int uid = getuid();

    printf("%d\n", uid);

    return 0;
}
```

What about definition?

# A series of things happen

user/gibuid.c

```c
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"

int main()
{
    int uid = getuid();

    printf("%d\n", uid);

    return 0;
}
```

What about definition?

Available in user/usys.S (generated from user/usys.pl)

Some other files come into play also. Let's see...

# System Calls [Recap]

- Programming interface to the services provided by the OS

- Typically used from a high-level language (C or C++) program.

- Typically a number is associated with each system call, and the OS maintains a table indexed according to these numbers.

# A series of things happen

user/gibuid.c

```c
int uid = getuid();
```

kernel/syscall.h

```c
// 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
#define SYS_getuid   22
```

# A series of things happen

user/gibuid.c

```c
int uid = getuid();
```

user/usys.pl

```perl
 3    # Generate usys.S, the stubs for syscalls.
 4
 5    print "# generated by usys.pl - do not edit\n";
 6
 7    print "#include \"kernel/syscall.h\"\n";
 8
 9    sub entry {
10        my $name = shift;
11        print ".global $name\n";
12        print "${name}:\n";
13        print " li a7, SYS_${name}\n";
14        print " ecall\n";
15        print " ret\n";
16    }
34    entry("dup");
35    entry("getpid");
36    entry("sbrk");
37    entry("sleep");
38    entry("uptime");
39    entry("getuid");
```

kernel/syscall.h

```c
// 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
#define SYS_getuid 22
```

# A series of things happen

kernel/syscall.h

```
// 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
#define SYS_getuid  22
```

user/gibuid.c

```c
int uid = getuid();
```

user/usys.pl

```perl
3     # Generate usys.S, the stubs for syscalls.
4
5     print "# generated by usys.pl - do not edit\n";
6
7     print "#include \"kernel/syscall.h\"\n";
8
9     sub entry {
10        my $name = shift;
11        print ".global $name\n";
12        print "${name}:\n";
13        print " li a7, SYS_${name}\n";
14        print " ecall\n";
15        print " ret\n";
16    }
34    entry("dup");
35    entry("getpid");
36    entry("sbrk");
37    entry("sleep");
38    entry("uptime");
39    entry("getuid");
40
```

# A series of things happen

user/gibuid.c

```
int uid = getuid();
```

user/usys.pl

```
3    # Generate usys.S, the stubs for syscalls.
4
5    print "# generated by usys.pl - do not edit\n";
6
7    print "#include \"kernel/syscall.h\"\n";
8
9    sub entry {
10       my $name = shift;
11       print ".global $name\n";
12       print "${name}:\n";
13       print " li a7, SYS_${name}\n";
14       print " ecall\n";
15       print " ret\n";
16   }
34   entry("dup");
35   entry("getpid");
36   entry("sbrk");
37   entry("sleep");
38   entry("uptime");
39   entry("getuid");
```

user/usys.S

```
1    # generated by usys.pl - do not edit
2    #include "kernel/syscall.h"
108     .global getuid
109    getuid:
110      li a7, SYS_getuid
111      ecall
112      ret
```

kernel/syscall.h

```
// 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
#define SYS_getuid  22
```

# A series of things happen

kernel/syscall.h

```
// 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
#define SYS_getuid   22
```

user/gibuid.c

```c
int uid = getuid();
```

Definition we are looking for..

user/usys.pl

```perl
3    # Generate usys.S, the stubs for syscalls.
4
5    print "# generated by usys.pl - do not edit\n";
6
7    print "#include \"kernel/syscall.h\"\n";
8
9    sub entry {
10       my $name = shift;
11       print ".global $name\n";
12       print "${name}:\n";
13       print " li a7, SYS_${name}\n";
14       print " ecall\n";
15       print " ret\n";
16   }
34   entry("dup");
35   entry("getpid");
36   entry("sbrk");
37   entry("sleep");
38   entry("uptime");
39   entry("getuid");
```

user/usys.S

```
1    # generated by usys.pl - do not edit
2    #include "kernel/syscall.h"
108    .global getuid
109   getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

Loads 22 at a7 register

# A series of things happen

user/gibuid.c

```
int uid = getuid();
```

user/usys.pl

```
3    # Generate usys.S, the stubs for syscalls.
4
5    print "# generated by usys.pl - do not edit\n";
6
7    print "#include \"kernel/syscall.h\"\n";
8
9    sub entry {
10       my $name = shift;
11       print ".global $name\n";
12       print "${name}:\n";
13       print " li a7, SYS_${name}\n";
14       print " ecall\n";
15       print " ret\n";
16   }
34   entry("dup");
35   entry("getpid");
36   entry("sbrk");
37   entry("sleep");
38   entry("uptime");
39   entry("getuid");
```

Definition we are looking for..

user/usys.S

```
1    # generated by usys.pl - do not edit
2    #include "kernel/syscall.h"
108     .global getuid
109     getuid:
110      li a7, SYS_getuid
111      ecall
112      ret
```

Traps

Loads 22 at a7 register

kernel/syscall.h

```
// 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
#define SYS_getuid  22
```

# A series of things happen

user/gibuid.c ->      user/Usys.s ->      kernel/trampoline.S/uservec()

```c
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

```
19   trampoline:
20   .align 4
21   .globl uservec
22   uservec:
23        #
24        # trap.c sets stvec to point here, so
25        # traps from user space start here,
26        # in supervisor mode, but with a
27        # user page table.
28        #
29
30        # save user a0 in sscratch so
31        # a0 can be used to get at TRAPFRAME.
32        csrw sscratch, a0
33
34        # each process has a separate p->trapframe memory area,
35        # but it's mapped to the same virtual address
36        # (TRAPFRAME) in every process's user page table.
37        li a0, TRAPFRAME
38
39        # save the user registers in TRAPFRAME
40        sd ra, 40(a0)
41        sd sp, 48(a0)
42        sd gp, 56(a0)
43        sd tp, 64(a0)
44        sd t0, 72(a0)
45        sd t1, 80(a0)
46        sd t2, 88(a0)
47        sd s0, 96(a0)
48        sd s1, 104(a0)
49        sd a1, 120(a0)
50        sd a2, 128(a0)
51        sd a3, 136(a0)
52        sd a4, 144(a0)
53        sd a5, 152(a0)
54        sd a6, 160(a0)
55        sd a7, 168(a0)
56        sd s2, 176(a0)
57        sd s3, 184(a0)
```

# A series of things happen

user/gibuid.c ->          user/Usys.s ->          kernel/trampoline.S/uservec()

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

```
19    trampoline:
20    .align 4
21    .globl uservec
22    uservec:
23            #
24            # trap.c sets stvec to point here, so
25            # traps from user space start here,
26            # in supervisor mode, but with a
27            # user page table.
28            #
29
30            # save user a0 in sscratch so
31            # a0 can be used to get at TRAPFRAME.
32            csrw sscratch, a0
33
34            # each process has a separate p->trapframe memory area,
35            # but it's mapped to the same virtual address
36            # (TRAPFRAME) in every process's user page table.
37            li a0, TRAPFRAME
38
39            # save the user registers in TRAPFRAME
40            sd ra, 40(a0)
41            sd sp, 48(a0)
42            sd gp, 56(a0)
```

```
97            # jump to usertrap(), which does not return
98            jr t0
99
47            sd s0, 96(a0)
48            sd s1, 104(a0)
49            sd a1, 120(a0)
50            sd a2, 128(a0)
51            sd a3, 136(a0)
52            sd a4, 144(a0)
53            sd a5, 152(a0)
54            sd a6, 160(a0)
55            sd a7, 168(a0)
56            sd s2, 176(a0)
57            sd s3, 184(a0)
```

# A series of things happen

user/gibuid.c ->        user/Usys.s ->        kernel/trampoline.S/uservec()  -> kernel/trap.c/usertrap()

```c
int uid = getuid();
```

```asm
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

```c
33  // handle an interrupt, exception, or system call from user space.
34  // called from trampoline.S
35  //
36  void
37  usertrap(void)
38  {
39    int which_dev = 0;
40
41    if((r_sstatus() & SSTATUS_SPP) != 0)
42      panic("usertrap: not from user mode");
43
44    // send interrupts and exceptions to kerneltrap(),
45    // since we're now in the kernel.
46    w_stvec((uint64)kernelvec);
47
48    struct proc *p = myproc();
49
50    // save user program counter.
51    p->trapframe->epc = r_sepc();
52
53    if(r_scause() == 8){
54      // system call
55
56      if(killed(p))
57        exit(-1);
58
59      // sepc points to the ecall instruction,
60      // but we want to return to the next instruction.
61      p->trapframe->epc += 4;
62
63      // an interrupt will change sepc, scause, and sstatus,
64      // so enable only now that we're done with those registers.
65      intr_on();
66
67      syscall();
```

```asm
19    trampoline:
20    .align 4
21    .globl uservec
22    uservec:
23          #
24          # trap.c sets stvec to point here, so
25          # traps from user space start here,
26          # in supervisor mode, but with a
27          # user page table.
28          #
29
30          # save user a0 in sscratch so
31          # a0 can be used to get at TRAPFRAME.
32          csrw sscratch, a0
33
34          # each process has a separate p->trapframe memory area,
35          # but it's mapped to the same virtual address
36          # (TRAPFRAME) in every process's user page table.
37          li a0, TRAPFRAME
38
39          # save the user registers in TRAPFRAME
40          sd ra, 40(a0)
41          sd sp, 48(a0)
42          sd gp, 56(a0)

97          # jump to usertrap(), which does not return
98          jr t0
99

47          sd s0, 96(a0)
48          sd s1, 104(a0)
49          sd a1, 120(a0)
50          sd a2, 128(a0)
51          sd a3, 136(a0)
52          sd a4, 144(a0)
53          sd a5, 152(a0)
54          sd a6, 160(a0)
55          sd a7, 168(a0)
56          sd s2, 176(a0)
57          sd s3, 184(a0)
```

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/syscall.c

```
33  // handle an interrupt, exception, or system call from user space.
34  // called from trampoline.S
35  //
36  void
37  usertrap(void)
38  {
39    int which_dev = 0;
40
41    if((r_sstatus() & SSTATUS_SPP) != 0)
42      panic("usertrap: not from user mode");
43
44    // send interrupts and exceptions to kerneltrap(),
45    // since we're now in the kernel.
46    w_stvec((uint64)kernelvec);
47
48    struct proc *p = myproc();
49
50    // save user program counter.
51    p->trapframe->epc = r_sepc();
52
53    if(r_scause() == 8){
54      // system call
55
56      if(killed(p))
57        exit(-1);
58
59      // sepc points to the ecall instruction,
60      // but we want to return to the next instruction.
61      p->trapframe->epc += 4;
62
63      // an interrupt will change sepc, scause, and sstatus,
64      // so enable only now that we're done with those registers.
65      intr_on();
66
67      syscall();
```

```
135   void
136   syscall(void)
137   {
138     int num;
139     struct proc *p = myproc();
140
141     num = p->trapframe->a7;
142     if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143       // Use num to lookup the system call function for num, call it,
144       // and store its return value in p->trapframe->a0
145       p->trapframe->a0 = syscalls[num]();
146     } else {
147       printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149       p->trapframe->a0 = -1;
150     }
151   }
```

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()

```c
int uid = getuid();
```

```asm
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/syscall.c

```c
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

kernel/proc.h

```c
82    enum procstate { UNUSED, USED, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
83
84    // Per-process state
85    struct proc {
86      struct spinlock lock;
87
88      // p->lock must be held when using these:
89      enum procstate state;        // Process state
90      void *chan;                  // If non-zero, sleeping on chan
91      int killed;                  // If non-zero, have been killed
92      int xstate;                  // Exit status to be returned to parent's wait
93      int pid;                     // Process ID
94
95      // wait_lock must be held when using this:
96      struct proc *parent;         // Parent process
97
98      // these are private to the process, so p->lock need not be held.
99      uint64 kstack;               // Virtual address of kernel stack
100     uint64 sz;                   // Size of process memory (bytes)
101     pagetable_t pagetable;       // User page table
102     struct trapframe *trapframe; // data page for trampoline.S
103     struct context context;      // swtch() here to run process
104     struct file *ofile[NOFILE];  // Open files
105     struct inode *cwd;           // Current directory
106     char name[16];               // Process name (debugging)
107    };
108
```

# A series of things happen

user/gibuid.c->user/Usys.s ->**kernel/trampoline.S/uservec()**->kernel/t

```c
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/proc.h

```c
82  enum procstate { UNUSED, USED, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
83
84  // Per-process state
85  struct proc {
86    struct spinlock lock;
87
88    // p->lock must be held when using these:
89    enum procstate state;        // Process state
90    void *chan;                  // If non-zero, sleeping on chan
91    int killed;                  // If non-zero, have been killed
92    int xstate;                  // Exit status to be returned to parent's wait
93    int pid;                     // Process ID
94
95    // wait_lock must be held when using this:
96    struct proc *parent;         // Parent process
97
98    // these are private to the process, so p->lock need not be held.
99    uint64 kstack;               // Virtual address of kernel stack
100   uint64 sz;                   // Size of process memory (bytes)
101   pagetable_t pagetable;       // User page table
102   struct trapframe *trapframe; // data page for trampoline.S
103   struct context context;      // swtch() here to run process
104   struct file *ofile[NOFILE];  // Open files
105   struct inode *cwd;           // Current directory
106   char name[16];               // Process name (debugging)
107 };
108
```

```c
43  struct trapframe {
44    /*   0 */ uint64 kernel_satp;   // kernel page table
45    /*   8 */ uint64 kernel_sp;     // top of process's kernel stack
46    /*  16 */ uint64 kernel_trap;   // usertrap()
47    /*  24 */ uint64 epc;           // saved user program counter
48    /*  32 */ uint64 kernel_hartid; // saved kernel tp
49    /*  40 */ uint64 ra;
50    /*  48 */ uint64 sp;
51    /*  56 */ uint64 gp;
52    /*  64 */ uint64 tp;
53    /*  72 */ uint64 t0;
54    /*  80 */ uint64 t1;
55    /*  88 */ uint64 t2;
56    /*  96 */ uint64 s0;
57    /* 104 */ uint64 s1;
58    /* 112 */ uint64 a0;
59    /* 120 */ uint64 a1;
60    /* 128 */ uint64 a2;
61    /* 136 */ uint64 a3;
62    /* 144 */ uint64 a4;
63    /* 152 */ uint64 a5;
64    /* 160 */ uint64 a6;
65    /* 168 */ uint64 a7;
66    /* 176 */ uint64 s2;
67    /* 184 */ uint64 s3;
68    /* 192 */ uint64 s4;
69    /* 200 */ uint64 s5;
70    /* 208 */ uint64 s6;
71    /* 216 */ uint64 s7;
72    /* 224 */ uint64 s8;
73    /* 232 */ uint64 s9;
74    /* 240 */ uint64 s10;
75    /* 248 */ uint64 s11;
76    /* 256 */ uint64 t3;
77    /* 264 */ uint64 t4;
78    /* 272 */ uint64 t5;
79    /* 280 */ uint64 t6;
80  };
```

Trap Frames are used to store the registers of the current thread when an interrupt occurs, or when there is a system service call (the transfer from user mode to kernel mode). Loaded in trampoline.S
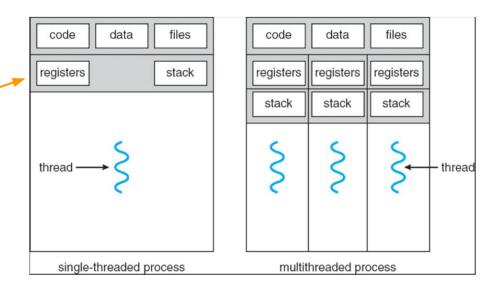
# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/proc.h

```
82   enum procstate { UNUSED, USED, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
83
84   // Per-process state
85   struct proc {
86     struct spinlock lock;
87
88     // p->lock must be held when using these:
89     enum procstate state;        // Process state
90     void *chan;                  // If non-zero, sleeping on chan
91     int killed;                  // If non-zero, have been killed
92     int xstate;                  // Exit status to be returned to parent's wait
93     int pid;                     // Process ID
94
95     // wait_lock must be held when using this:
96     struct proc *parent;         // Parent process
97
98     // these are private to the process, so p->lock need not be held.
99     uint64 kstack;               // Virtual address of kernel stack
100    uint64 sz;                   // Size of process memory (bytes)
101    pagetable_t pagetable;       // User page table
102    struct trapframe *trapframe; // data page for trampoline.S
103    struct context context;      // swtch() here to run process
104    struct file *ofile[NOFILE];  // Open files
105    struct inode *cwd;           // Current directory
106    char name[16];               // Process name (debugging)
107  };
108
```

## Process Control Block (PCB)

| Process management | Memory management | File management |
|---|---|---|
| Registers | Pointer to text segment | Root directory |
| Program counter | Pointer to data segment | Working directory |
| Program status word | Pointer to stack segment | File descriptors |
| Stack pointer | | User ID |
| Process state | | Group ID |
| Priority | | |
| Scheduling parameters | | |
| Process ID | | |
| Parent process | | |
| Process group | | |
| Signals | | |
| Time when process started | | |
| CPU time used | | |
| Children's CPU time | | |
| Time of next alarm | | |

Figure: Fields of a PCB

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/proc.h

```
82   enum procstate { UNUSED, USED, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
83
84   // Per-process state
85   struct proc {
86     struct spinlock lock;
87
88     // p->lock must be held when using these:
89     enum procstate state;        // Process state
90     void *chan;                  // If non-zero, sleeping on chan
91     int killed;                  // If non-zero, have been killed
92     int xstate;                  // Exit status to be returned to parent's wait
93     int pid;                     // Process ID
94
95     // wait_lock must be held when using this:
96     struct proc *parent;         // Parent process
97
98     // these are private to the process, so p->lock need not be held.
99     uint64 kstack;               // Virtual address of kernel stack
100    uint64 sz;                   // Size of process memory (bytes)
101    pagetable_t pagetable;       // User page table
102    struct trapframe *trapframe; // data page for trampoline.S
103    struct context context;      // swtch() here to run process
104    struct file *ofile[NOFILE];  // Open files
105    struct inode *cwd;           // Current directory
106    char name[16];               // Process name (debugging)
107  };
108
```

## Multithreaded Processes



single-threaded process     multithreaded process

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()

```c
int uid = getuid();
```

```asm
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/syscall.c

```c
107    // An array mapping syscall numbers from syscall.h
108    // to the function that handles the system call.
109    static uint64 (*syscalls[])(void) = {
110    [SYS_fork]    sys_fork,
111    [SYS_exit]    sys_exit,
112    [SYS_wait]    sys_wait,
113    [SYS_pipe]    sys_pipe,
114    [SYS_read]    sys_read,
115    [SYS_kill]    sys_kill,
116    [SYS_exec]    sys_exec,
117    [SYS_fstat]   sys_fstat,
118    [SYS_chdir]   sys_chdir,
119    [SYS_dup]     sys_dup,
120    [SYS_getpid]  sys_getpid,
121    [SYS_sbrk]    sys_sbrk,
122    [SYS_sleep]   sys_sleep,
123    [SYS_uptime]  sys_uptime,
124    [SYS_open]    sys_open,
125    [SYS_write]   sys_write,
126    [SYS_mknod]   sys_mknod,
127    [SYS_unlink]  sys_unlink,
128    [SYS_link]    sys_link,
129    [SYS_mkdir]   sys_mkdir,
130    [SYS_close]   sys_close,
131    [SYS_getuid]  sys_getuid,
132    };
```

kernel/syscall.c

```c
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()

```
int uid = getuid();
```

```
108     .global getuid
109     getuid:
110      li a7, SYS_getuid
111      ecall
112      ret
```

kernel/syscall.c

```
107    // An array mapping syscall numbers from syscall.h
108    // to the function that handles the system call.
109    static uint64 (*syscalls[])(void) = {
110    [SYS_fork]      sys_fork,
111    [SYS_exit]      sys_exit,
112    [SYS_wait]      sys_wait,
113    [SYS_pipe]      sys_pipe,
114    [SYS_read]      sys_read,
115    [SYS_kill]      sys_kill,
116    [SYS_exec]      sys_exec,
117    [SYS_fstat]     sys_fstat,
118    [SYS_chdir]     sys_chdir,
119    [SYS_dup]       sys_dup,
120    [SYS_getpid]    sys_getpid,
121    [SYS_sbrk]      sys_sbrk,
122    [SYS_sleep]     sys_sleep,
123    [SYS_uptime]    sys_uptime,
124    [SYS_open]      sys_open,
125    [SYS_write]     sys_write,
126    [SYS_mknod]     sys_mknod,
127    [SYS_unlink]    sys_unlink,
128    [SYS_link]      sys_link,
129    [SYS_mkdir]     sys_mkdir,
130    [SYS_close]     sys_close,
131    [SYS_getuid]    sys_getuid,
132    };
```

kernel/syscall.c

```
       // Prototypes for the functions that handle system calls.
       extern uint64 sys_fork(void);
       extern uint64 sys_exit(void);
       extern uint64 sys_wait(void);
       extern uint64 sys_pipe(void);
       extern uint64 sys_read(void);
       extern uint64 sys_kill(void);
       extern uint64 sys_exec(void);
       extern uint64 sys_fstat(void);
       extern uint64 sys_chdir(void);
       extern uint64 sys_dup(void);
       extern uint64 sys_getpid(void);
       extern uint64 sys_sbrk(void);
       extern uint64 sys_sleep(void);
       extern uint64 sys_uptime(void);
       extern uint64 sys_open(void);
       extern uint64 sys_write(void);
       extern uint64 sys_mknod(void);
       extern uint64 sys_unlink(void);
       extern uint64 sys_link(void);
       extern uint64 sys_mkdir(void);
       extern uint64 sys_close(void);
       extern uint64 sys_getuid(void);
```

Defined mostly in sysproc.c and sysfile.c

kernel/syscall.c

```
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

kernel/defs.h

```
190      // number of elements in fixed-size array
191      #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
192
```

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()

```
int uid = getuid();
```

```
108     .global getuid
109     getuid:
110      li a7, SYS_getuid
111      ecall
112      ret
```

These were defined in syscall.h (loaded in a7). So, line 145 will call sys_getuid() [defined in sysproc.c]

### kernel/syscall.c

```
107    // An array mapping syscall numbers from syscall.h
108    // to the function that handles the system call.
109    static uint64 (*syscalls[])(void) = {
110    [SYS_fork]     sys_fork,
111    [SYS_exit]     sys_exit,
112    [SYS_wait]     sys_wait,
113    [SYS_pipe]     sys_pipe,
114    [SYS_read]     sys_read,
115    [SYS_kill]     sys_kill,
116    [SYS_exec]     sys_exec,
117    [SYS_fstat]    sys_fstat,
118    [SYS_chdir]    sys_chdir,
119    [SYS_dup]      sys_dup,
120    [SYS_getpid]   sys_getpid,
121    [SYS_sbrk]     sys_sbrk,
122    [SYS_sleep]    sys_sleep,
123    [SYS_uptime]   sys_uptime,
124    [SYS_open]     sys_open,
125    [SYS_write]    sys_write,
126    [SYS_mknod]    sys_mknod,
127    [SYS_unlink]   sys_unlink,
128    [SYS_link]     sys_link,
129    [SYS_mkdir]    sys_mkdir,
130    [SYS_close]    sys_close,
131    [SYS_getuid]   sys_getuid,
132    };
```

kernel/syscall.c

```
      // Prototypes for the functions that handle system calls.
      extern uint64 sys_fork(void);
      extern uint64 sys_exit(void);
      extern uint64 sys_wait(void);
      extern uint64 sys_pipe(void);
      extern uint64 sys_read(void);
      extern uint64 sys_kill(void);
      extern uint64 sys_exec(void);
      extern uint64 sys_fstat(void);
      extern uint64 sys_chdir(void);
      extern uint64 sys_dup(void);
      extern uint64 sys_getpid(void);
      extern uint64 sys_sbrk(void);
      extern uint64 sys_sleep(void);
      extern uint64 sys_uptime(void);
      extern uint64 sys_open(void);
      extern uint64 sys_write(void);
      extern uint64 sys_mknod(void);
      extern uint64 sys_unlink(void);
      extern uint64 sys_link(void);
      extern uint64 sys_mkdir(void);
      extern uint64 sys_close(void);
      extern uint64 sys_getuid(void);
```

Defined mostly in sysproc.c and sysfile.c

### kernel/syscall.c

```
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

### kernel/defs.h

```
190    // number of elements in fixed-size array
191    #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
192
```

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()->kernel/sysproc.c

```c
int uid = getuid();
```

```asm
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

These were defined in syscall.h (loaded in a7). So, line 145 will call sys_getuid() [defined in sysproc.c]

kernel/syscall.c

```c
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

kernel/sysproc.c

```c
95    uint64
96    sys_getuid(void)
97    {
98      return getuid();
99    }
```

kernel/defs.h

```c
190    // number of elements in fixed-size array
191    #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
192
```

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()->kernel/sysproc.c

->kernel/proc.c

```
int uid = getuid();
```

```
108      .global getuid
109      getuid:
110       li a7, SYS_getuid
111       ecall
112       ret
```

kernel/sysproc.c

```
95    uint64
96    sys_getuid(void)
97    {
98      return getuid();
99    }
```

kernel/proc.c

```
697    int uid=123;
698
699    int getuid(void) {
700      return uid;
701    }
```

kernel/syscall.c

```
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

kernel/defs.h

```
190    // number of elements in fixed-size array
191    #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
192
```

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()->kernel/sysproc.c

->kernel/proc.c

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/syscall.c

```
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

kernel/sysproc.c

```
95    uint64
96    sys_getuid(void)
97    {
98      return getuid();
99    }
```

to make getuid() accessible from kernel/sysproc.c, add this **declaration** in **kernel/defs.h**. If no new function was defined, this line wouldn't have been required.

kernel/proc.c

```
697    int uid=123;
698
699    int getuid(void) {
700      return uid;
701    }
```

```
107    int              either_copyin(void  u
108    void             procdump(void);
109    int              getuid(void);
```

kernel/defs.h

```
190    // number of elements in fixed-size array
191    #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
192
```

# A series of things happen

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()->kernel/sysproc.c

->kernel/proc.c

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/syscall.c

```
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

kernel/sysproc.c

```
95     uint64
96     sys_getuid(void)
97     {
98       return getuid();
99     }
```

kernel/proc.c

```
697    int uid=123;
698
699    int getuid(void) {
700      return uid;
701    }
```

*Why bother adding another function in proc.c?*
**Ans:** It's the style of the coding in xv6, treating sysproc.c as a mediator having **sys call handlers** who verifies input and then delegates. We should adhere to the practice.

kernel/defs.h

```
190    // number of elements in fixed-size array
191    #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
192
```

```
107    int          either_copyin(void  d
108    void         procdump(void);
109    int          getuid(void);
```

# A series of things happened(now returning)

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()->kernel/sysproc.c

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/sysproc.c

```
95    uint64
96    sys_getuid(void)
97    {
98      return getuid();
99    }
```

kernel/syscall.c

```
135    void
136    syscall(void)
137    {
138      int num;
139      struct proc *p = myproc();
140
141      num = p->trapframe->a7;
142      if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143        // Use num to lookup the system call function for num, call it,
144        // and store its return value in p->trapframe->a0
145        p->trapframe->a0 = syscalls[num]();
146      } else {
147        printf("%d %s: unknown sys call %d\n",
148               p->pid, p->name, num);
149        p->trapframe->a0 = -1;
150      }
151    }
```

kernel/defs.h

```
190    // number of elements in fixed-size array
191    #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
192
```

# A series of things happened(now returning)

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->syscall.c/syscall()

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

kernel/syscall.c

```
135  void
136  syscall(void)
137  {
138    int num;
139    struct proc *p = myproc();
140
141    num = p->trapframe->a7;
142    if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143      // Use num to lookup the system call function for num, call it,
144      // and store its return value in p->trapframe->a0
145      p->trapframe->a0 = syscalls[num]();
146    } else {
147      printf("%d %s: unknown sys call %d\n",
148             p->pid, p->name, num);
149      p->trapframe->a0 = -1;
150    }
151  }
```

kernel/defs.h

```
190    // number of elements in fixed-size array
191    #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
192
```

# A series of things happenIng

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()

```
int uid = getuid();
```

```
108     .global getuid
109     getuid:
110      li a7, SYS_getuid
111      ecall
112      ret
```

```c
36   void
37   usertrap(void)
38   {
39     int which_dev = 0;
40
41     if((r_sstatus() & SSTATUS_SPP) != 0)
42       panic("usertrap: not from user mode");
43
44     // send interrupts and exceptions to kerneltrap(),
45     // since we're now in the kernel.
46     w_stvec((uint64)kernelvec);
47
48     struct proc *p = myproc();
49
50     // save user program counter.
51     p->trapframe->epc = r_sepc();
52
53     if(r_scause() == 8){
54       // system call
55
56       if(killed(p))
57         exit(-1);
58
59       // sepc points to the ecall instruction,
60       // but we want to return to the next instruction.
61       p->trapframe->epc += 4;
62
63       // an interrupt will change sepc, scause, and sstat
64       // so enable only now that we're done with those re
65       intr_on();
66
67       syscall();
68     } else if((which_dev = devintr()) != 0){
69       // ok
70     } else {
71       printf("usertrap(): unexpected scause 0x%lx pid=%d"
72       printf("            sepc=0x%lx stval=0x%lx\n", r_se
73       setkilled(p);
74     }
75
76     if(killed(p))
77       exit(-1);
78
79     // give up the CPU if this is a timer interrupt.
80     if(which_dev == 2)
81       yield();
82
83     usertrapret();
84   }
85
```

# A series of things happenIng

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->

kernel/trap.c/usertrapret()

```
int uid = getuid();
```

```
89   void
90   usertrapret(void)
91   {
92     struct proc *p = myproc();
93
94     // we're about to switch the destination of traps from
95     // kerneltrap() to usertrap(), so turn off interrupts until
96     // we're back in user space, where usertrap() is correct.
97     intr_off();
98
99     // send syscalls, interrupts, and exceptions to uservec in trampo
100    uint64 trampoline_uservec = TRAMPOLINE + (uservec - trampoline);
101    w_stvec(trampoline_uservec);
102
103    // set up trapframe values that uservec will need when
104    // the process next traps into the kernel.
105    p->trapframe->kernel_satp = r_satp();         // kernel page tabl
106    p->trapframe->kernel_sp = p->kstack + PGSIZE; // process's kernel
107    p->trapframe->kernel_trap = (uint64)usertrap;
108    p->trapframe->kernel_hartid = r_tp();         // hartid for cpui
109
110    // set up the registers that trampoline.S's sret will use
111    // to get to user space.
112
113    // set S Previous Privilege mode to User.
114    unsigned long x = r_sstatus();
115    x &= ~SSTATUS_SPP; // clear SPP to 0 for user mode
116    x |= SSTATUS_SPIE; // enable interrupts in user mode
117    w_sstatus(x);
118
119    // set S Exception Program Counter to the saved user pc.
120    w_sepc(p->trapframe->epc);
121
122    // tell trampoline.S the user page table to switch to.
123    uint64 satp = MAKE_SATP(p->pagetable);
124
125    // jump to userret in trampoline.S at the top of memory, which
126    // switches to the user page table, restores user registers,
127    // and switches to user mode with sret.
128    uint64 trampoline_userret = TRAMPOLINE + (userret - trampoline);
129    ((void (*)(uint64))trampoline_userret)(satp);
130  }
131
```

```
36   void
37   usertrap(void)
38   {
39     int which_dev = 0;
40
41     if((r_sstatus() & SSTATUS_SPP) != 0)
42       panic("usertrap: not from user mode");
43
44     // send interrupts and exceptions to kerneltrap(),
45     // since we're now in the kernel.
46     w_stvec((uint64)kernelvec);
47
48     struct proc *p = myproc();
49
50     // save user program counter.
51     p->trapframe->epc = r_sepc();
52
53     if(r_scause() == 8){
54       // system call
55
56       if(killed(p))
57         exit(-1);
58
59       // sepc points to the ecall instruction,
60       // but we want to return to the next instruction.
61       p->trapframe->epc += 4;
62
63       // an interrupt will change sepc, scause, and ssta
64       // so enable only now that we're done with those re
65       intr_on();
66
67       syscall();
68     } else if((which_dev = devintr()) != 0){
69       // ok
70     } else {
71       printf("usertrap(): unexpected scause 0x%lx pid=%d
72       printf("            sepc=0x%lx stval=0x%lx\n", r_se
73       setkilled(p);
74     }
75
76     if(killed(p))
77       exit(-1);
78
79     // give up the CPU if this is a timer interrupt.
80     if(which_dev == 2)
81       yield();
82
83     usertrapret();
84   }
85
```

# A series of things happenIng

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->kernel/trampoline.S/userret()

```
int uid = getuid();
```

```
 89   void
 90   usertrapret(void)
 91   {
 92       struct proc *p = myproc();
 93
 94       // we're about to switch the destination of traps from
 95       // kerneltrap() to usertrap(), so turn off interrupts until
 96       // we're back in user space, where usertrap() is correct.
 97       intr_off();
 98
 99       // send syscalls, interrupts, and exceptions to uservec in trampoline
100       uint64 trampoline_uservec = TRAMPOLINE + (uservec - trampoline);
101       w_stvec(trampoline_uservec);
102
103       // set up trapframe values that uservec will need when
104       // the process next traps into the kernel.
105       p->trapframe->kernel_satp = r_satp();         // kernel page table
106       p->trapframe->kernel_sp = p->kstack + PGSIZE; // process's kernel
107       p->trapframe->kernel_trap = (uint64)usertrap;
108       p->trapframe->kernel_hartid = r_tp();          // hartid for cpui
109
110       // set up the registers that trampoline.S's sret will use
111       // to get to user space.
112
113       // set S Previous Privilege mode to User.
114       unsigned long x = r_sstatus();
115       x &= ~SSTATUS_SPP; // clear SPP to 0 for user mode
116       x |= SSTATUS_SPIE; // enable interrupts in user mode
117       w_sstatus(x);
118
119       // set S Exception Program Counter to the saved user pc.
120       w_sepc(p->trapframe->epc);
121
122       // tell trampoline.S the user page table to switch to.
123       uint64 satp = MAKE_SATP(p->pagetable);
124
125       // jump to userret in trampoline.S at the top of memory, which
126       // switches to the user page table, restores user registers,
127       // and switches to user mode with sret.
128       uint64 trampoline_userret = TRAMPOLINE + (userret - trampoline);
129       ((void (*)(uint64))trampoline_userret)(satp);
130   }
131
```

kernel/trampoline.S

```
100       .globl userret
101   userret:
102           # userret(pagetable)
103           # called by usertrapret() in trap.c to
104           # switch from kernel to user.
105           # a0: user page table, for satp.
106
107           # switch to the user page table.
108           sfence.vma zero, zero
109           csrw satp, a0
110           sfence.vma zero, zero
111
112           li a0, TRAPFRAME
113
114           # restore all but a0 from TRAPFRAME
115           ld ra, 40(a0)
116           ld sp, 48(a0)
117           ld gp, 56(a0)
118           ld tp, 64(a0)
119           ld t0, 72(a0)
120           ld t1, 80(a0)
121           ld t2, 88(a0)
122           ld s0, 96(a0)
123           ld s1, 104(a0)
124           ld a1, 120(a0)
125           ld a2, 128(a0)
126           ld a3, 136(a0)
127           ld a4, 144(a0)
128           ld a5, 152(a0)
129           ld a6, 160(a0)
130           ld a7, 168(a0)
131           ld s2, 176(a0)
132           ld s3, 184(a0)
133           ld s4, 192(a0)
```

# A series of things happenIng

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->kernel/trampoline.S/userret()

```
int uid = getuid();
```

```
89   void
90   usertrapret(void)
91   {
92     struct proc *p = myproc();
93
94     // we're about to switch the destination of traps from
95     // kerneltrap() to usertrap(), so turn off interrupts until
96     // we're back in user space, where usertrap() is correct.
97     intr_off();
98
99     // send syscalls, interrupts, and exceptions to uservec in trampo
100    uint64 trampoline_uservec = TRAMPOLINE + (uservec - trampoline);
101    w_stvec(trampoline_uservec);
102
103    // set up trapframe values that uservec will need when
104    // the process next traps into the kernel.
105    p->trapframe->kernel_satp = r_satp();          // kernel page tab
106    p->trapframe->kernel_sp = p->kstack + PGSIZE; // process's kernel
107    p->trapframe->kernel_trap = (uint64)usertrap;
108    p->trapframe->kernel_hartid = r_tp();          // hartid for cpui
109
110    // set up the registers that trampoline.S's sret will use
111    // to get to user space.
112
113    // set S Previous Privilege mode to User.
114    unsigned long x = r_sstatus();
115    x &= ~SSTATUS_SPP; // clear SPP to 0 for user mode
116    x |= SSTATUS_SPIE; // enable interrupts in user mode
117    w_sstatus(x);
118
119    // set S Exception Program Counter to the saved user pc.
120    w_sepc(p->trapframe->epc);
121
122    // tell trampoline.S the user page table to switch to.
123    uint64 satp = MAKE_SATP(p->pagetable);
124
125    // jump to userret in trampoline.S at the top of memory, which
126    // switches to the user page table, restores user registers,
127    // and switches to user mode with sret.
128    uint64 trampoline_userret = TRAMPOLINE + (userret - trampoline);
129    ((void (*)(uint64))trampoline_userret)(satp);
130  }
131
```

## kernel/trampoline.S

```
100   .globl userret
101   userret:
102        # userret(pagetable)
103        # called by usertrapret() in trap.c to
104        # switch from kernel to user.
105        # a0: user page table, for satp.
106
107        # switch to the user page table.
108        sfence.vma zero, zero
109        csrw satp, a0
110        sfence.vma zero, zero
111
112        li a0, TRAPFRAME
113
114        # restore all but a0 from TRAPFRAME
115        ld ra, 40(a0)
116        ld sp, 48(a0)
117        ld gp, 56(a0)
118        ld tp, 64(a0)
119        ld t0, 72(a0)
120        ld t1, 80(a0)
121        ld t2, 88(a0)
122        ld s0, 96(a0)
123        ld s1, 104(a0)
124        ld a1, 120(a0)
125        ld a2, 128(a0)
126        ld a3, 136(a0)
127        ld a4, 144(a0)
128        ld a5, 152(a0)
129        ld a6, 160(a0)
130        ld a7, 168(a0)
131        ld s2, 176(a0)
132        ld s3, 184(a0)
133        ld s4, 192(a0)
```

trapframe->a0 contains the return value

# A series of things happenIng

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->kernel/trampoline.S/userret()

```c
int uid = getuid();
```

```c
89   void
90   usertrapret(void)
91   {
92     struct proc *p = myproc();
93
94     // we're about to switch the destination of traps from
95     // kerneltrap() to usertrap(), so turn off interrupts until
96     // we're back in user space, where usertrap() is correct.
97     intr_off();
98
99     // send syscalls, interrupts, and exceptions to uservec in trampoline
100    uint64 trampoline_uservec = TRAMPOLINE + (uservec - trampoline);
101    w_stvec(trampoline_uservec);
102
103    // set up trapframe values that uservec will need when
104    // the process next traps into the kernel.
105    p->trapframe->kernel_satp = r_satp();        // kernel page table
106    p->trapframe->kernel_sp = p->kstack + PGSIZE; // process's kernel
107    p->trapframe->kernel_trap = (uint64)usertrap;
108    p->trapframe->kernel_hartid = r_tp();        // hartid for cpuid
109
110    // set up the registers that trampoline.S's sret will use
111    // to get to user space.
112
113    // set S Previous Privilege mode to User.
114    unsigned long x = r_sstatus();
115    x &= ~SSTATUS_SPP; // clear SPP to 0 for user mode
116    x |= SSTATUS_SPIE; // enable interrupts in user mode
117    w_sstatus(x);
118
119    // set S Exception Program Counter to the saved user pc.
120    w_sepc(p->trapframe->epc);
121
122    // tell trampoline.S the user page table to switch to.
123    uint64 satp = MAKE_SATP(p->pagetable);
124
125    // jump to userret in trampoline.S at the top of memory, which
126    // switches to the user page table, restores user registers,
127    // and switches to user mode with sret.
128    uint64 trampoline_userret = TRAMPOLINE + (userret - trampoline);
129    ((void (*)(uint64))trampoline_userret)(satp);
130  }
131
```

## kernel/trampoline.S

```asm
100    .globl userret
101  userret:
102       # userret(pagetable)
103       # called by usertrapret() in trap.c to
104       # switch from kernel to user.
105       # a0: user page table, for satp.
106
107       # switch to the user page table.
108       sfence.vma zero, zero
109       csrw satp, a0
110       sfence.vma zero, zero
111
112       li a0, TRAPFRAME
113
114       # restore all but a0 from TRAPFRAME
115       ld ra, 40(a0)
116       ld sp, 48(a0)
117       ld gp, 56(a0)
118       ld tp, 64(a0)
119       ld t0, 72(a0)
120       ld t1, 80(a0)
121       ld t2, 88(a0)
122       ld s0, 96(a0)
```

trapframe->a0 contains the return value

```asm
146    # restore user a0
147    ld a0, 112(a0)
148
149    # return to user mode and user pc.
150    # usertrapret() set up sstatus and sepc.
151    sret
152
132       ld s3, 184(a0)
133       ld s4, 192(a0)
```

# A series of things happen [RECAP & NOTE]

user/gibuid.c ->     user/Usys.s ->     kernel/trampoline.S/uservec()  -> kernel/trap.c/usertrap()

```
int uid = getuid();
```

```
108    .global getuid
109    getuid:
110     li a7, SYS_getuid
111     ecall
112     ret
```

```
33   // handle an interrupt, exception, or system call from user space.
34   // called from trampoline.S
35   //
36   void
37   usertrap(void)
38   {
39     int which_dev = 0;
40
41     if((r_sstatus() & SSTATUS_SPP) != 0)
42       panic("usertrap: not from user mode");
43
44     // send interrupts and exceptions to kerneltrap(),
45     // since we're now in the kernel.
46     w_stvec((uint64)kernelvec);
47
48     struct proc *p = myproc();
49
50     // save user program counter.
51     p->trapframe->epc = r_sepc();
52
53     if(r_scause() == 8){
54       // system call
55
56       if(killed(p))
57         exit(-1);
58
59       // sepc points to the ecall instruction,
60       // but we want to return to the next instruction.
61       p->trapframe->epc += 4;
62
63       // an interrupt will change sepc, scause, and sstatus,
64       // so enable only now that we're done with those registers.
65       intr_on();
66
67       syscall();
```

```
19   trampoline:
20   .align 4
21   .globl uservec
22   uservec:
23         #
24         # trap.c sets stvec to point here, so
25         # traps from user space start here,
26         # in supervisor mode, but with a
27         # user page table.
28         #
29
30         # save user a0 in sscratch so
31         # a0 can be used to get at TRAPFRAME.
32         csrw sscratch, a0
33
34         # each process has a separate p->trapframe memory area,
35         # but it's mapped to the same virtual address
36         # (TRAPFRAME) in every process's user page table.
37         li a0, TRAPFRAME
38
39         # save the user registers in TRAPFRAME
40         sd ra, 40(a0)
41         sd sp, 48(a0)
         sd gp, 56(a0)
97         # jump to usertrap(), which does not return
98         jr t0
99
47         sd s0, 96(a0)
48         sd s1, 104(a0)
49         sd a1, 120(a0)
50         sd a2, 128(a0)
51         sd a3, 136(a0)
52         sd a4, 144(a0)
53         sd a5, 152(a0)
54         sd a6, 160(a0)
55         sd a7, 168(a0)
56         sd s2, 176(a0)
57         sd s3, 184(a0)
```

# A series of things happenIng

user/gibuid.c->user/Usys.s ->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()->kernel/trampoline.S/userret()

```
int uid = getuid();
```

```
89    void
90    usertrapret(void)
91    {
92        struct proc *p = myproc();
93
94        // we're about to switch the destination of traps from
95        // kerneltrap() to usertrap(), so turn off interrupts until
96        // we're back in user space, where usertrap() is correct.
97        intr_off();
98
99        // send syscalls, interrupts, and exceptions to uservec in trampo
100       uint64 trampoline_uservec = TRAMPOLINE + (uservec - trampoline);
101       w_stvec(trampoline_uservec);
102
103       // set up trapframe values that uservec will need when
104       // the process next traps into the kernel.
105       p->trapframe->kernel_satp = r_satp();        // kernel page tab
106       p->trapframe->kernel_sp = p->kstack + PGSIZE; // process's kerne
107       p->trapframe->kernel_trap = (uint64)usertrap;
108       p->trapframe->kernel_hartid = r_tp();        // hartid for cpui
109
110       // set up the registers that trampoline.S's sret will use
111       // to get to user space.
112
113       // set S Previous Privilege mode to User.
114       unsigned long x = r_sstatus();
115       x &= ~SSTATUS_SPP; // clear SPP to 0 for user mode
116       x |= SSTATUS_SPIE; // enable interrupts in user mode
117       w_sstatus(x);
118
119       // set S Exception Program Counter to the saved user pc.
120       w_sepc(p->trapframe->epc);
121
122       // tell trampoline.S the user page table to switch to.
123       uint64 satp = MAKE_SATP(p->pagetable);
124
125       // jump to userret in trampoline.S at the top of memory, which
126       // switches to the user page table, restores user registers,
127       // and switches to user mode with sret.
128       uint64 trampoline_userret = TRAMPOLINE + (userret - trampoline);
129       ((void (*)(uint64))trampoline_userret)(satp);
130   }
131
```

kernel/trampoline.S

```
100       .globl userret
101   userret:
102           # userret(pagetable)
10            # called by usertrapret() in trap.c to
04            # switch from kernel to user.
105           # a0: user page table, for satp.
106
107           # switch to the user page table.
108           sfence.vma zero, zero
109           csrw satp, a0
110           sfence.vma zero, zero
111
112           li a0, TRAPFRAME
113
114           # restore all but a0 from TRAPFRAME
115           ld ra, 40(a0)
116           ld sp, 48(a0)
117           ld gp, 56(a0)
118           ld tp, 64(a0)
119           ld t0, 72(a0)
120           ld t1, 80(a0)
121           ld t2, 88(a0)
122           ld s0, 96(a0)
```

trapframe->a0 contains the return value

```
146       # restore user a0
147       ld a0, 112(a0)
148
149       # return to user mode and user pc.
150       # usertrapret() set up sstatus and sepc.
151       sret
152
132           ld s3, 184(a0)
133           ld s4, 192(a0)
```

# A series of things happened

user/gibuid.c->user/Usys.s

```
int uid = getuid();        108    .global getuid
                           109    getuid:
                           110     li a7, SYS_getuid
                           111     ecall
                           112     ret  ⬅
```

# A series of things happened (We are back!)

user/gibuid.c

```c
#include "kernel/types.h"
#include "kernel/stat.h"
#include "user/user.h"

int main()
{
    int uid = getuid();

    printf("%d\n", uid);

    return 0;
}
```

# Okay, but how to pass arguments?

Say, we want to set the uid to 456! (As we know it's simply a variable)

The call will be setuid(456)

```
if(!setuid(456)) {
    printf("setuid success\n");
} else {
    printf("setuid failed\n");
}
```

xv6 stores the arguments passed in registers a0, a1, a2, a3, a4, a5

# Yet another system call

user/gibuid.c->user/Usys.s->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()*[stores registers in trapframe]*->
kernel/syscall.c/syscall()

kernel/syscall.c

```
135  void
136  syscall(void)
137  {
138    int num;
139    struct proc *p = myproc();
140
141    num = p->trapframe->a7;
142    if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143      // Use num to lookup the system call function for num, call it,
144      // and store its return value in p->trapframe->a0
145      p->trapframe->a0 = syscalls[num]();
146    } else {
147      printf("%d %s: unknown sys call %d\n",
148              p->pid, p->name, num);
149      p->trapframe->a0 = -1;
150    }
151  }
```

# Yet another system call

user/gibuid.c->user/Usys.s->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()*[stores registers in trapframe]*->
kernel/syscall.c/syscall()

kernel/syscall.c

```
static uint64
argraw(int n)
{
  struct proc *p = myproc();
  switch (n) {
  case 0:
    return p->trapframe->a0;
  case 1:
    return p->trapframe->a1;
  case 2:
    return p->trapframe->a2;
  case 3:
    return p->trapframe->a3;
  case 4:
    return p->trapframe->a4;
  case 5:
    return p->trapframe->a5;
  }
  panic("argraw");
  return -1;
}
```

```
void
argint(int n, int *ip)
{
  *ip = argraw(n);
}

// Retrieve an argument as a pointer.
// Doesn't check for legality, since
// copyin/copyout will do that.
void
argaddr(int n, uint64 *ip)
{
  *ip = argraw(n);
}

// Fetch the nth word-sized system ca
// Copies into buf, at most max.
// Returns string length if OK (inclu
int
argstr(int n, char *buf, int max)
{
  uint64 addr;
  argaddr(n, &addr);
  return fetchstr(addr, buf, max);
}
```

kernel/syscall.c

```
135  void
136  syscall(void)
137  {
138    int num;
139    struct proc *p = myproc();
140
141    num = p->trapframe->a7;
142    if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
143      // Use num to lookup the system call function for num, call it,
144      // and store its return value in p->trapframe->a0
145      p->trapframe->a0 = syscalls[num]();
146    } else {
147      printf("%d %s: unknown sys call %d\n",
148              p->pid, p->name, num);
149      p->trapframe->a0 = -1;
150    }
151  }
```

# Yet another system call

user/gibuid.c->user/Usys.s->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()*[stores registers in trapframe]*->
kernel/syscall.c/syscall() *[say, it called sys_setuid() in sysproc.c]*

kernel/syscall.c

```
static uint64
argraw(int n)
{
  struct proc *p = myproc();
  switch (n) {
  case 0:
    return p->trapframe->a0;
  case 1:
    return p->trapframe->a1;
  case 2:
    return p->trapframe->a2;
  case 3:
    return p->trapframe->a3;
  case 4:
    return p->trapframe->a4;
  case 5:
    return p->trapframe->a5;
  }
  panic("argraw");
  return -1;
}
```

```
void
argint(int n, int *ip)
{
  *ip = argraw(n);
}

// Retrieve an argument as a pointer.
// Doesn't check for legality, since
// copyin/copyout will do that.
void
argaddr(int n, uint64 *ip)
{
  *ip = argraw(n);
}

// Fetch the nth word-sized system cal
// Copies into buf, at most max.
// Returns string length if OK (inclu
int
argstr(int n, char *buf, int max)
{
  uint64 addr;
  argaddr(n, &addr);
  return fetchstr(addr, buf, max);
}
```

kernel/sysproc.c

```
uint64
sys_setuid(void)
{
  int uid;
  argint(0, &uid);

  if (uid < 0 || uid > 65535) {
    return -1;
  }

  setuid(uid);

  return 0;
}
```

# Yet another system call

user/gibuid.c->user/Usys.s->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()*[stores registers in trapframe]*->
kernel/syscall.c/syscall() *[say, it called sys_setuid() in sysproc.c]*

kernel/syscall.c

```
static uint64
argraw(int n)
{
  struct proc *p = myproc();
  switch (n) {
  case 0:
    return p->trapframe->a0;
  case 1:
    return p->trapframe->a1;
  case 2:
    return p->trapframe->a2;
  case 3:
    return p->trapframe->a3;
  case 4:
    return p->trapframe->a4;
  case 5:
    return p->trapframe->a5;
  }
  panic("argraw");
  return -1;
}
```

```
void
argint(int n, int *ip)
{
  *ip = argraw(n);
}

// Retrieve an argument as a pointer.
// Doesn't check for legality, since
// copyin/copyout will do that.
void
argaddr(int n, uint64 *ip)
{
  *ip = argraw(n);
}

// Fetch the nth word-sized system ca
// Copies into buf, at most max.
// Returns string length if OK (inclu
int
argstr(int n, char *buf, int max)
{
  uint64 addr;
  argaddr(n, &addr);
  return fetchstr(addr, buf, max);
}
```

kernel/sysproc.c

```
uint64
sys_setuid(void)
{
  int uid;
  argint(0, &uid);

  if (uid < 0 || uid > 65535) {
    return -1;
  }

  setuid(uid);

  return 0;
}
```

user/gibuid.c

```
if(!setuid(456)) {
    printf("setuid success\n");
} else {
    printf("setuid failed\n");
}
```

**Only one argument. For others, n would be 1,2,3,..**

# Yet another system call

user/gibuid.c->user/Usys.s->kernel/trampoline.S/uservec()->kernel/trap.c/usertrap()*[stores registers in trapframe]*->
kernel/syscall.c/syscall() *[say, it called sys_setuid() in sysproc.c]*

kernel/syscall.c

```
static uint64
argraw(int n)
{
  struct proc *p = myproc();
  switch (n) {
  case 0:
    return p->trapframe->a0;
  case 1:
    return p->trapframe->a1;
  case 2:
    return p->trapframe->a2;
  case 3:
    return p->trapframe->a3;
  case 4:
    return p->trapframe->a4;
  case 5:
    return p->trapframe->a5;
  }
  panic("argraw");
  return -1;
}
```

```
void
argint(int n, int *ip)
{
  *ip = argraw(n);
}

// Retrieve an argument as a pointer.
// Doesn't check for legality, since
// copyin/copyout will do that.
void
argaddr(int n, uint64 *ip)
{
  *ip = argraw(n);
}

// Fetch the nth word-sized system cal
// Copies into buf, at most max.
// Returns string length if OK (inclu
int
argstr(int n, char *buf, int max)
{
  uint64 addr;
  argaddr(n, &addr);
  return fetchstr(addr, buf, max);
}
```

kernel/sysproc.c

```
uint64
sys_setuid(void)
{
  int uid;
  argint(0, &uid);

  if (uid < 0 || uid > 65535) {
    return -1;
  }

  setuid(uid);

  return 0;
}
```

kernel/proc.c

```
int setuid(int newuid) {
  uid = newuid;
  return 0;
}
```

# Thank You



**Catalin Pit**
@catalinmpit

I'm so glad I did a CS degree learning Data Structures, Algorithms, Operating Systems, Databases, Linear Algebra, Software Engineering, Information Analysis, Networking and many more.

How else would I have been able to center a div with CSS or change the background color?

ProgrammerHumor.io