# Figure 1.1 - Various standards
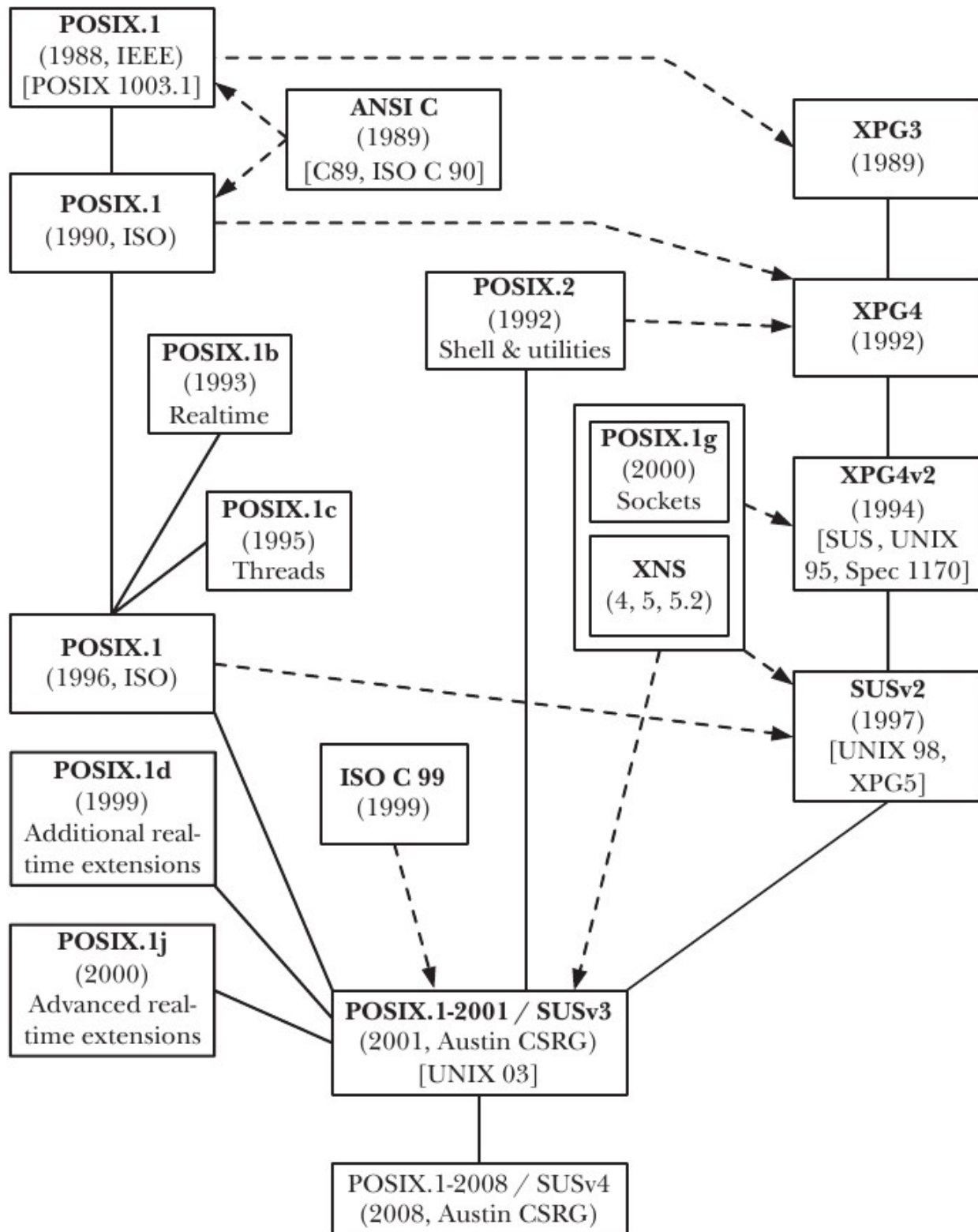
POSIX.1
(1988, IEEE)
[POSIX 1003.1]

ANSI C
(1989)
[C89, ISO C 90]

XPG3
(1989)

POSIX.1
(1990, ISO)

POSIX.2
(1992)
Shell & utilities

XPG4
(1992)

POSIX.1b
(1993)
Realtime

POSIX.1c
(1995)
Threads

POSIX.1g
(2000)
Sockets

XPG4v2
(1994)
[SUS, UNIX
95, Spec 1170]

XNS
(4, 5, 5.2)

POSIX.1
(1996, ISO)

SUSv2
(1997)
[UNIX 98,
XPG5]

POSIX.1d
(1999)
Additional real-
time extensions

ISO C 99
(1999)

POSIX.1j
(2000)
Advanced real-
time extensions

POSIX.1-2001 / SUSv3
(2001, Austin CSRG)
[UNIX 03]

POSIX.1-2008 / SUSv4
(2008, Austin CSRG)

# Figure 1.2 - system call flow



Figure 1.2 - system call flow

# Figure 1.3 - process memory layout

Virtual memory address
(hexadecimal)

| | |
|---|---|
| Kernel (mapped into process virtual memory, but not accessible to program) | ← /proc/kallsyms provides addresses of kernel symbols in this region ( /proc/ksyms in kernel 2.4 and earlier) |

0xC0000000

*argv, environ*

Stack
(grows downwards)

Top of stack →

↓

(unallocated memory)

↑

Program break →

Heap
(grows upwards)

Uninitialized data (bss)                    ← *&end*

Initialized data                            ← *&edata*

Text (program code)                         ← *&etext*

0x08048000

0x00000000

increasing virtual addesses ↑

**Figure 1.4 - process address space to page-frames**

# Figure 1.5 - changing process credentials

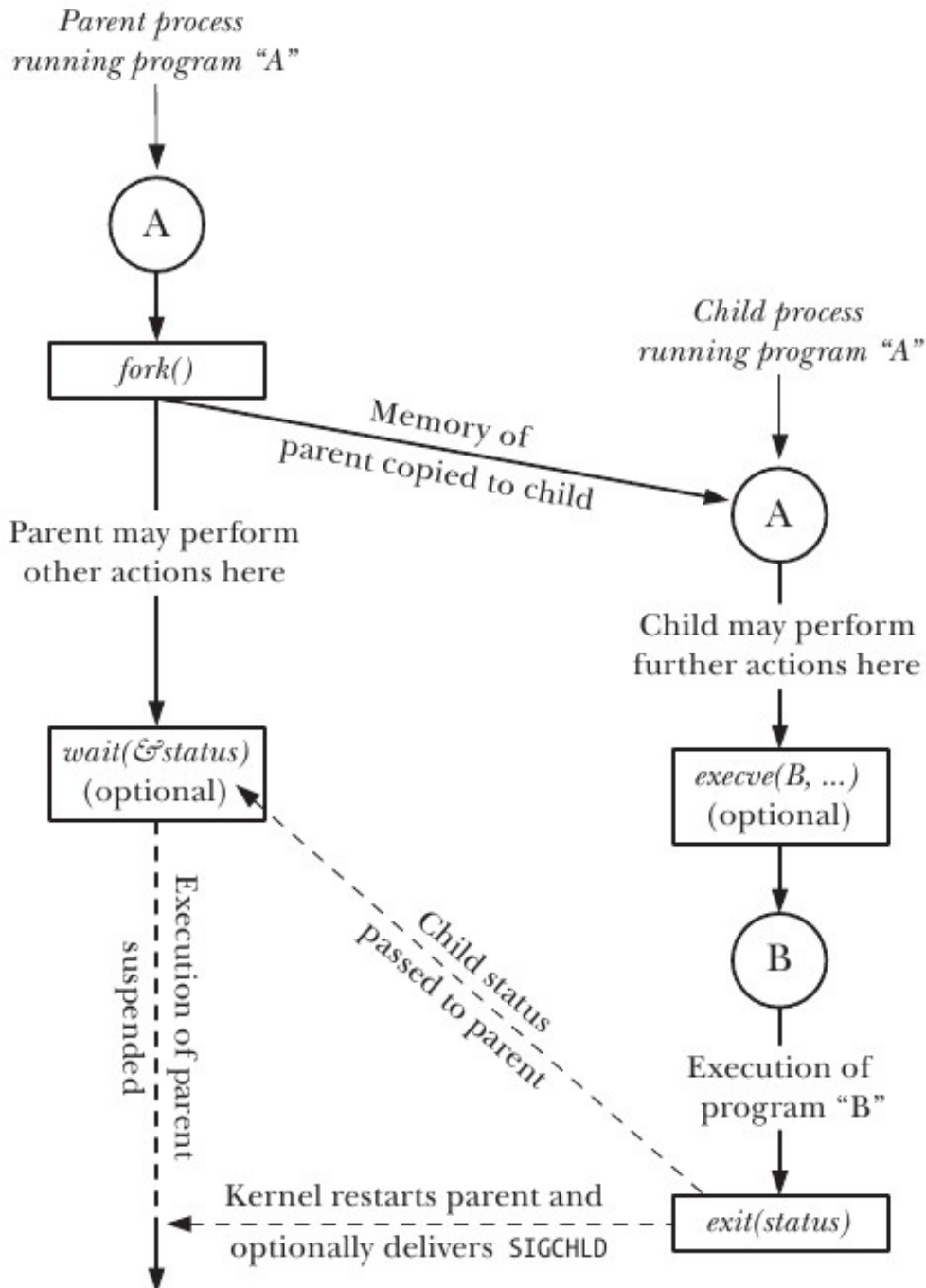| Interface | Purpose and effect within: | | Portability |
| --- | --- | --- | --- |
| | unprivileged process | privileged process | |
| *setuid(u)*<br>*setgid(g)* | Change effective ID to the same value as current real or saved set ID | Change real, effective, and saved set IDs to any (single) value | Specified in SUSv3; BSD derivatives have different semantics |
| *seteuid(e)*<br>*setegid(e)* | Change effective ID to the same value as current real or saved set ID | Change effective ID to any value | Specified in SUSv3 |
| *setreuid(r, e)*<br>*setregid(r, e)* | (Independently) change real ID to same value as current real or effective ID, and effective ID to same value as current real, effective, or saved set ID | (Independently) change real and effective IDs to any values | Specified in SUSv3, but operation varies across implementations |
| *setresuid(r, e, s)*<br>*setresgid(r, e, s)* | (Independently) change real, effective, and saved set IDs to same value as current real, effective, or saved set ID | (Independently) change real, effective, and saved set IDs to any values | Not in SUSv3 and present on few other UNIX implementations |
| *setfsuid(u)*<br>*setfsgid(u)* | Change file-system ID to same value as current real, effective, file system, or saved set ID | Change file-system ID to any value | Linux-specific |
| *setgroups(n, l)* | Can't be called from an unprivileged process | Set supplementary group IDs to any values | Not in SUSv3, but available on all UNIX implementations |

**Figure 1.6 - process creation and related**
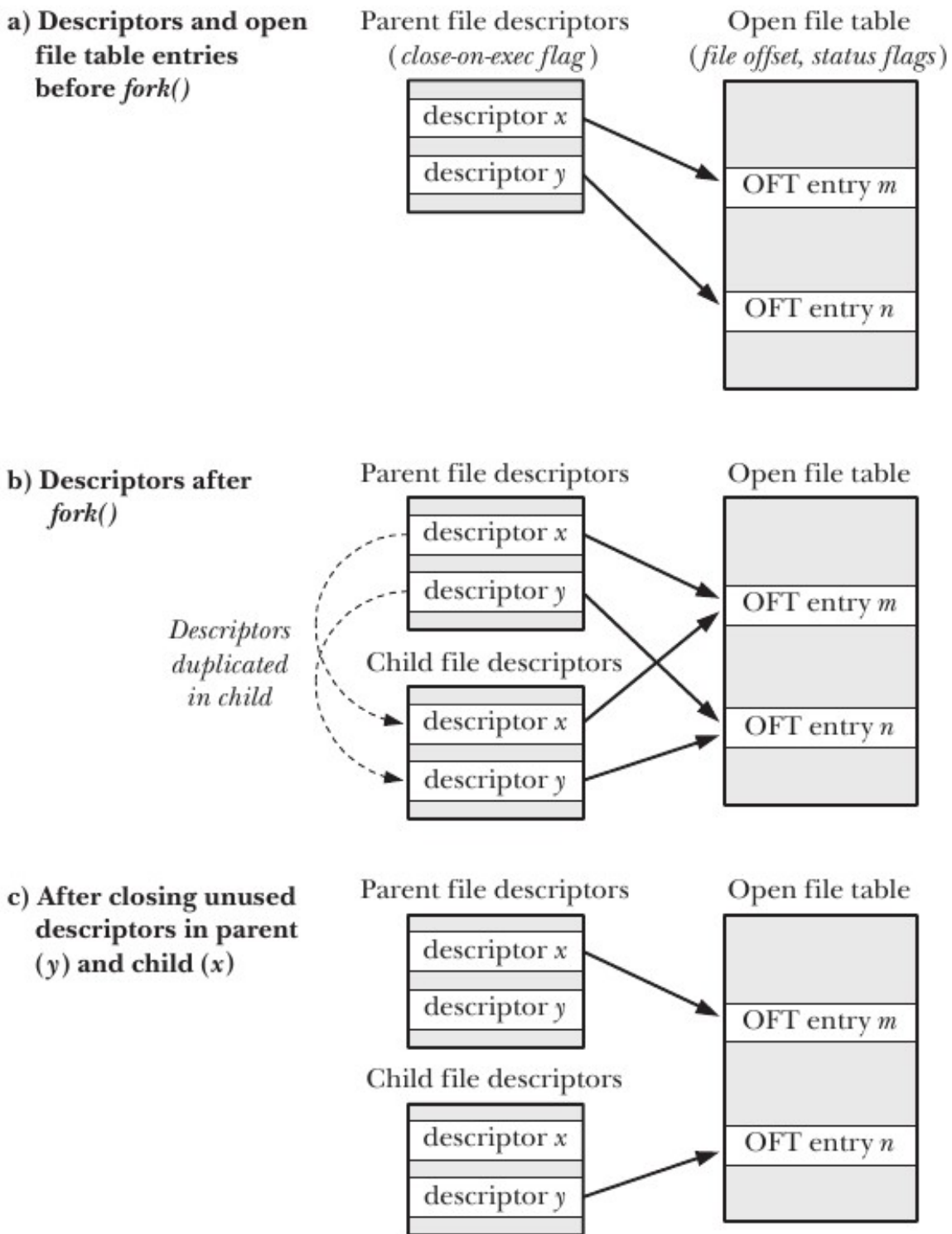
# Figure 1.7 - process creation and file descriptors

**a) Descriptors and open file table entries before *fork()***

Parent file descriptors
( *close-on-exec flag* )

Open file table
( *file offset, status flags* )

| descriptor x |
| descriptor y |

OFT entry m

OFT entry n

**b) Descriptors after *fork()***

Parent file descriptors

Open file table

| descriptor x |
| descriptor y |

*Descriptors duplicated in child*

Child file descriptors

| descriptor x |
| descriptor y |

OFT entry m

OFT entry n

**c) After closing unused descriptors in parent (*y*) and child (*x*)**

Parent file descriptors

Open file table

| descriptor x |
| descriptor y |

Child file descriptors

| descriptor x |
| descriptor y |

OFT entry m

OFT entry n

# Figure 1.8 - copy – on – write

**Before modification**

Parent page table     Physical page frames

PT entry 211

Child page table

PT entry 211

Frame 1998

*Unused page frames*

**After modification**

Parent page table     Physical page frames

PT entry 211

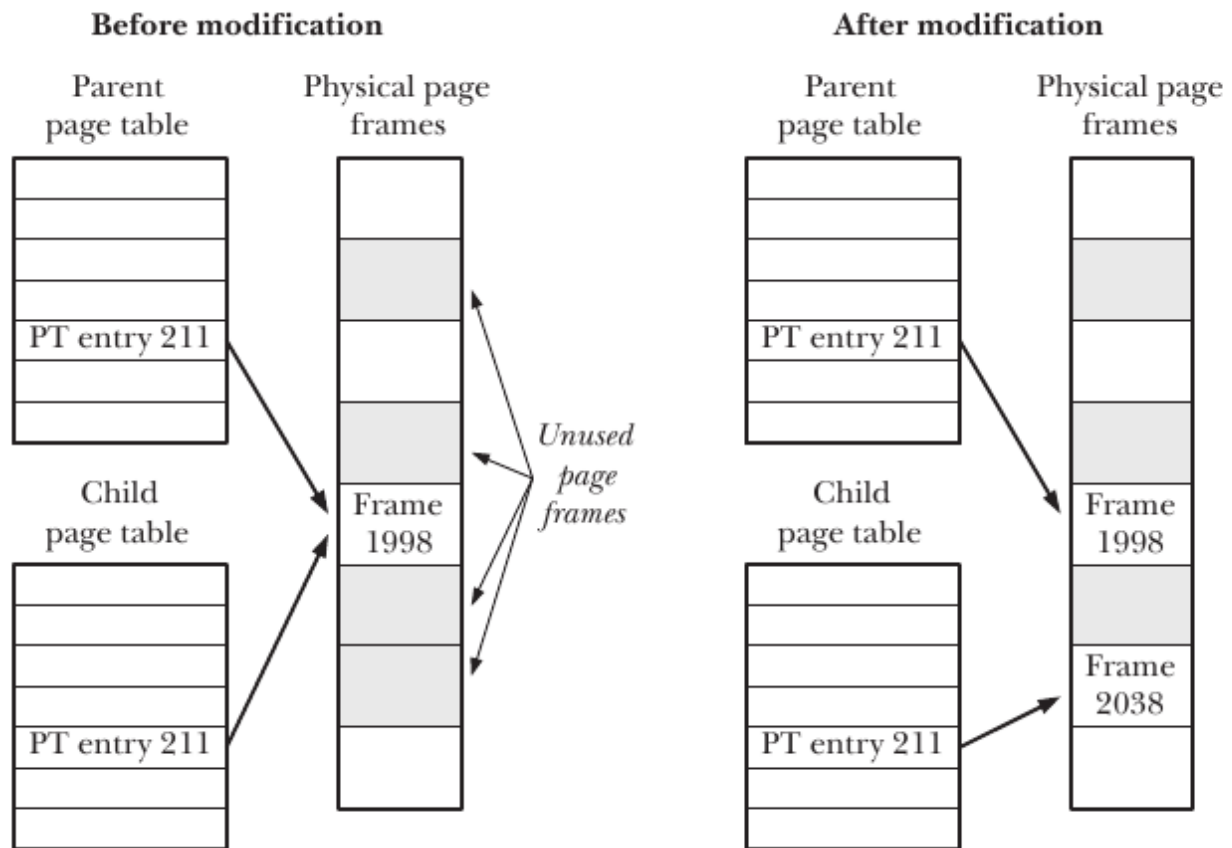Child page table

PT entry 211

Frame 1998

Frame 2038

**Figure 1.9 - status variable of wait() /waitpid()**

**Figure 1.10 - execve() and related library APIs**

```
#include <unistd.h>

int execve(const char *pathname, char *const argv[], char *const envp[]);
```
Never returns on success; returns −1 on error

```
#include <unistd.h>

int execle(const char *pathname, const char *arg, ...
              /* , (char *) NULL, char *const envp[] */ );
int execlp(const char *filename, const char *arg, ...
              /* , (char *) NULL */);
int execvp(const char *filename, char *const argv[]);
int execv(const char *pathname, char *const argv[]);
int execl(const char *pathname, const char *arg, ...
              /* , (char *) NULL */);
```
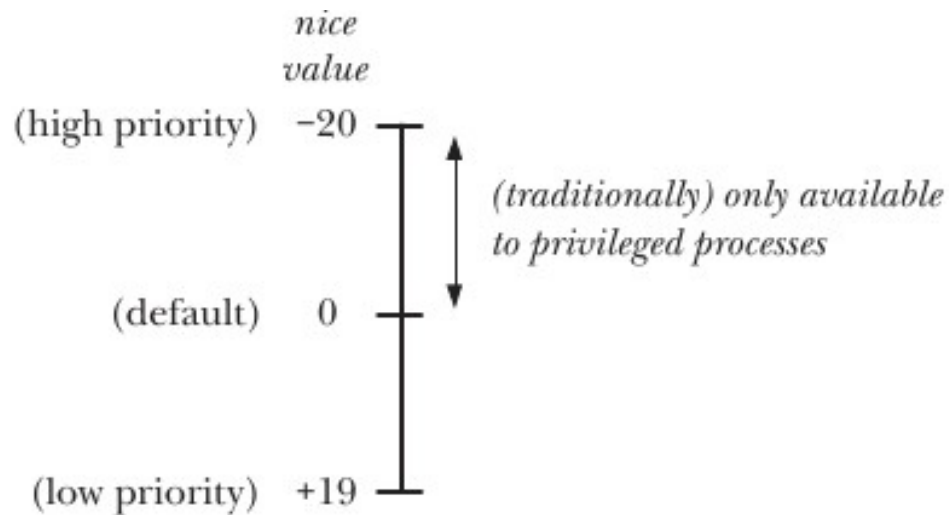None of the above returns on success; all return −1 on error

**Figure 1.11 - execve() and library APIs continued...**

| Function | Specification of program file (−, p) | Specification of arguments (v, l) | Source of environment (e, −) |
|---|---|---|---|
| execve() | pathname | array | envp argument |
| execle() | pathname | list | envp argument |
| execlp() | filename + PATH | list | caller's environ |
| execvp() | filename + PATH | array | caller's environ |
| execv() | pathname | array | caller's environ |
| execl() | pathname | list | caller's environ |

**Figure 1.12  -  nice priorities and scheduling policies**



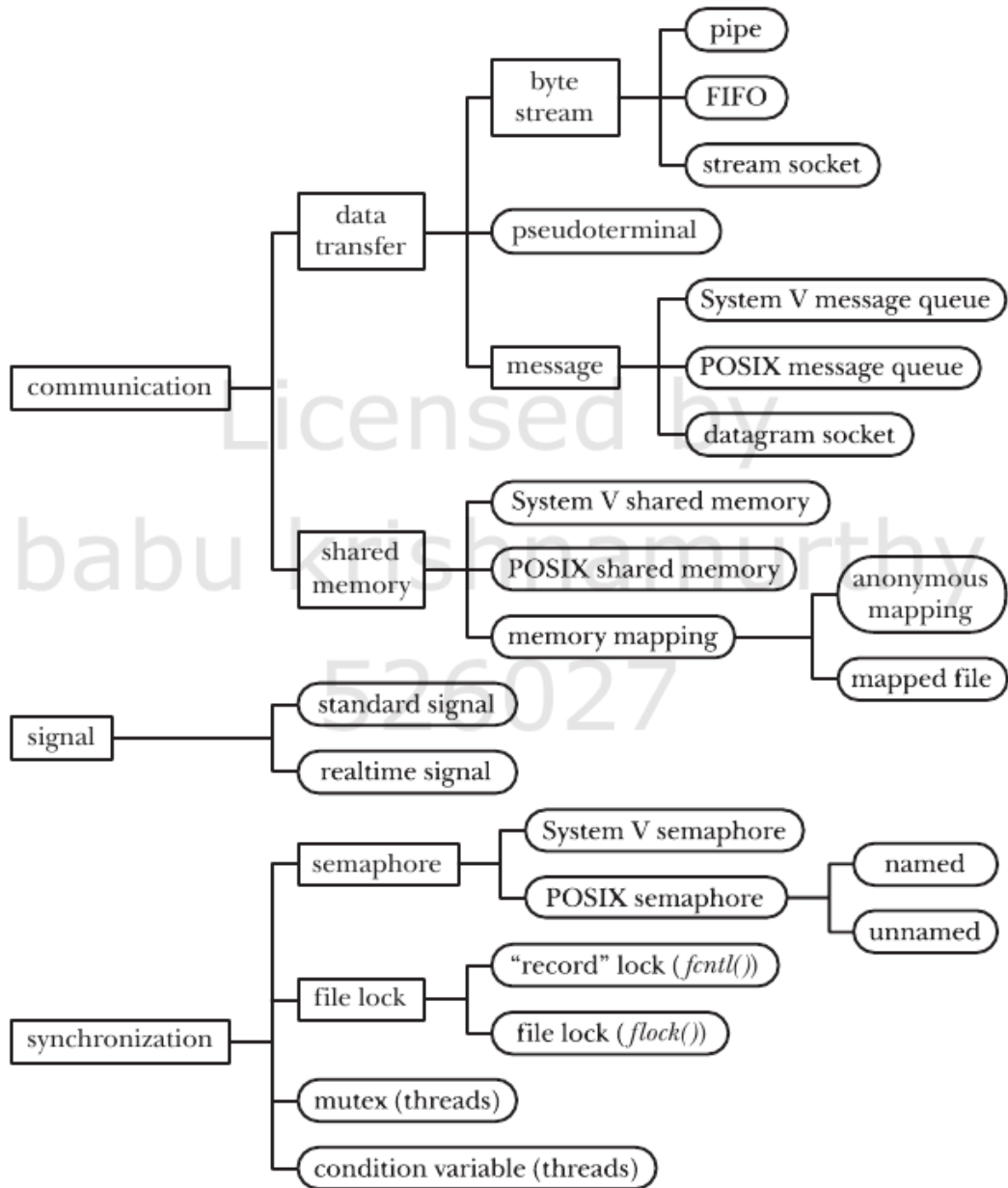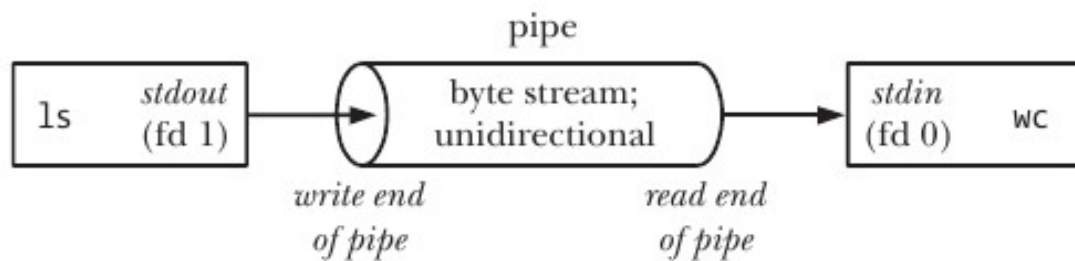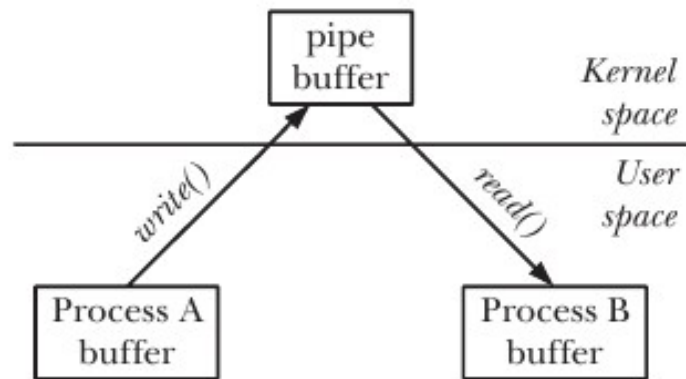| Policy | Description | SUSv3 |
|---|---|---|
| SCHED_FIFO | Realtime first-in first-out | • |
| SCHED_RR | Realtime round-robin | • |
| SCHED_OTHER | Standard round-robin time-sharing | • |
| SCHED_BATCH | Similar to SCHED_OTHER, but intended for batch execution (since Linux 2.6.16) | |
| SCHED_IDLE | Similar to SCHED_OTHER, but with priority even lower than nice value +19 (since Linux 2.6.23) | |

# Figure 1.13 - various IPC mechanisms



communication
- data transfer
  - byte stream
    - pipe
    - FIFO
    - stream socket
  - pseudoterminal
  - message
    - System V message queue
    - POSIX message queue
    - datagram socket
- shared memory
  - System V shared memory
  - POSIX shared memory
  - memory mapping
    - anonymous mapping
    - mapped file

signal
- standard signal
- realtime signal

synchronization
- semaphore
  - System V semaphore
  - POSIX semaphore
    - named
    - unnamed
- file lock
  - "record" lock (*fcntl()*)
  - file lock (*flock()*)
- mutex (threads)
- condition variable (threads)

**Figure 1.14 - pipe communication**

pipe
buffer

*Kernel
space*

*User
space*

*write()*

*read()*

Process A
buffer

Process B
buffer

pipe

ls    *stdout*
*(fd 1)*

*byte stream;*
*unidirectional*

*stdin*
*(fd 0)*    WC

*write end*
*of pipe*

*read end*
*of pipe*

# Figure 1.15 - summary of IPC mechanisms

| Facility type | Name used to identify object | Handle used to refer to object in programs |
|---|---|---|
| Pipe | no name | file descriptor |
| FIFO | pathname | file descriptor |
| UNIX domain socket | pathname | file descriptor |
| Internet domain socket | IP address + port number | file descriptor |
| System V message queue | System V IPC key | System V IPC identifier |
| System V semaphore | System V IPC key | System V IPC identifier |
| System V shared memory | System V IPC key | System V IPC identifier |
| POSIX message queue | POSIX IPC pathname | *mqd_t* (message queue descriptor) |
| POSIX named semaphore | POSIX IPC pathname | *sem_t* * (semaphore pointer) |
| POSIX unnamed semaphore | no name | *sem_t* * (semaphore pointer) |
| POSIX shared memory | POSIX IPC pathname | file descriptor |
| Anonymous mapping | no name | none |
| Memory-mapped file | pathname | file descriptor |
| *flock()* lock | pathname | file descriptor |
| *fcntl()* lock | pathname | file descriptor |

**Figure 1.16 - summary of IPC mechanisms**

| Facility type | Accessibility | Persistence |
|---|---|---|
| Pipe | only by related processes | process |
| FIFO | permissions mask | process |
| UNIX domain socket | permissions mask | process |
| Internet domain socket | by any process | process |
| System V message queue | permissions mask | kernel |
| System V semaphore | permissions mask | kernel |
| System V shared memory | permissions mask | kernel |
| POSIX message queue | permissions mask | kernel |
| POSIX named semaphore | permissions mask | kernel |
| POSIX unnamed semaphore | permissions of underlying memory | depends |
| POSIX shared memory | permissions mask | kernel |
| Anonymous mapping | only by related processes | process |
| Memory-mapped file | permissions mask | file system |
| *flock()* file lock | *open()* of file | process |
| *fcntl()* file lock | *open()* of file | process |

# Figure 1.17 - POSIX SYS V IPC mechanisms

| Interface | Message queues | Semaphores | Shared memory |
|---|---|---|---|
| Header file | <sys/msg.h> | <sys/sem.h> | <sys/shm.h> |
| Associated data structure | *msqid_ds* | *semid_ds* | *shmid_ds* |
| Create/open object | *msgget()* | *semget()* | *shmget() + shmat()* |
| Close object | (none) | (none) | *shmdt()* |
| Control operations | *msgctl()* | *semctl()* | *shmctl()* |
| Performing IPC | *msgsnd()*—write message *msgrcv()*—read message | *semop()*—test/adjust semaphore | access memory in shared region |

## Figure 1.18 - struct ipc_perm { }, struct ipc_ids{ } and struct semid_ds { }

```
struct ipc_perm {
    key_t           __key;          /* Key, as supplied to 'get' call */
    uid_t           uid;            /* Owner's user ID */
    gid_t           gid;            /* Owner's group ID */
    uid_t           cuid;           /* Creator's user ID */
    gid_t           cgid;           /* Creator's group ID */
    unsigned short mode;            /* Permissions */
    unsigned short __seq;           /* Sequence number */
};
```
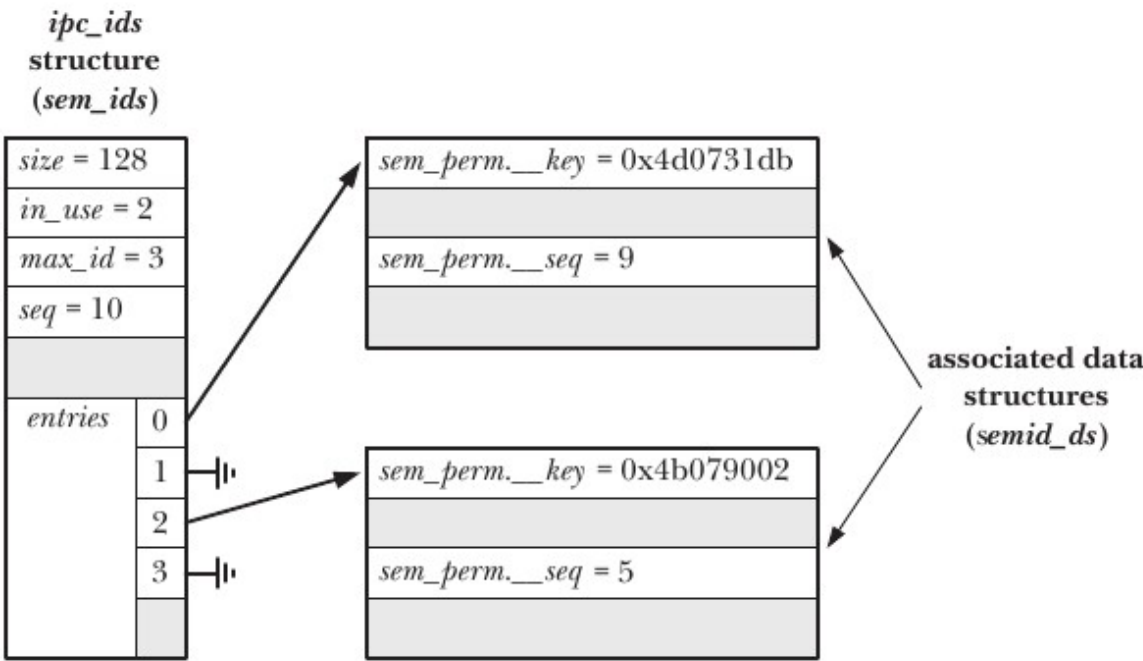
*ipc_ids*
**structure**
(*sem_ids*)

| | |
|---|---|
| *size* = 128 | |
| *in_use* = 2 | |
| *max_id* = 3 | |
| *seq* = 10 | |
| | |
| *entries* | 0 |
| | 1 |
| | 2 |
| | 3 |

*sem_perm.__key* = 0x4d0731db

*sem_perm.__seq* = 9

*sem_perm.__key* = 0x4b079002

*sem_perm.__seq* = 5

**associated data structures**
(*semid_ds*)

**Figure 1.19 - POSIX SYS V message queue related objects**

```
struct msqid_ds {
    struct ipc_perm msg_perm;        /* Ownership and permissions */
    time_t          msg_stime;       /* Time of last msgsnd() */
    time_t          msg_rtime;       /* Time of last msgrcv() */
    time_t          msg_ctime;       /* Time of last change */
    unsigned long   __msg_cbytes;    /* Number of bytes in queue */
    msgqnum_t       msg_qnum;        /* Number of messages in queue */
    msglen_t        msg_qbytes;      /* Maximum bytes in queue */
    pid_t           msg_lspid;       /* PID of last msgsnd() */
    pid_t           msg_lrpid;       /* PID of last msgrcv() */
};
```

# Figure 1.20 - a POSIX SYS V message queue usage scenario

Server

Server sends
response (*mtype*
= PID of client)  (3)

message queue

(2) Server reads
request (select
*msgtyp* = 1)

Client sends request
(*mtype* = 1, *mtext*
includes client PID)  (1)

(4) Client reads
response (select
*msgtyp* = own PID)

Client

---

Server creates child
to handle request  (4)

Server

Server child

*fork()*

(3) Server reads
request

(5) Server child sends
response(s)

Server MQ

Client MQ

*msgget(
IPC_PRIVATE, ...)*

(1) Client creates
private queue

Client

(2) Client sends request to
Server MQ (*mtext* includes
ID of client queue)

(6) Client reads
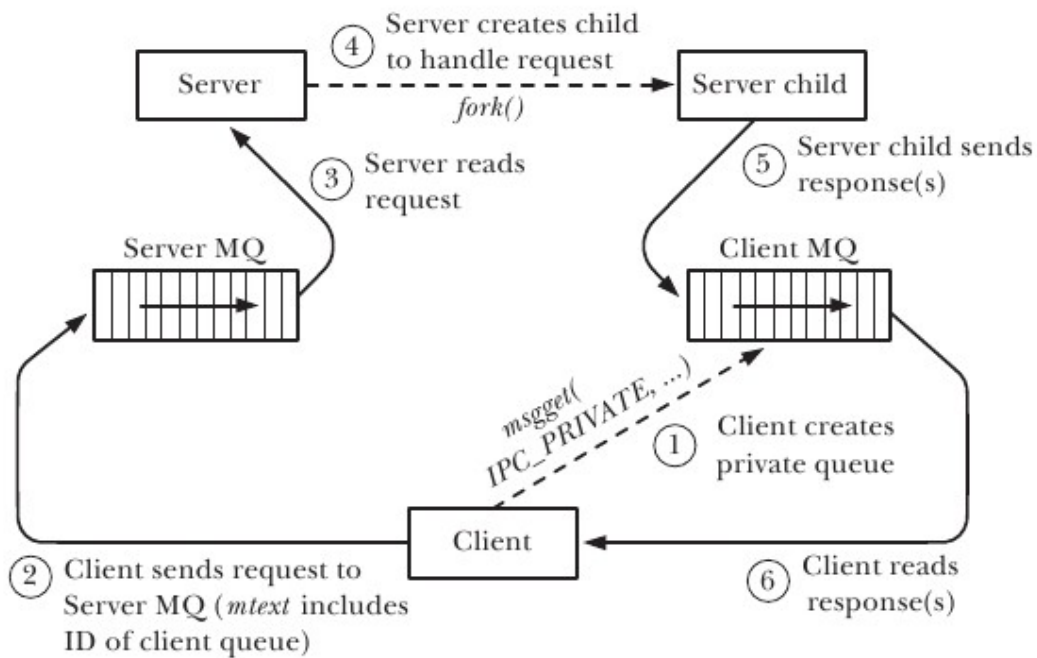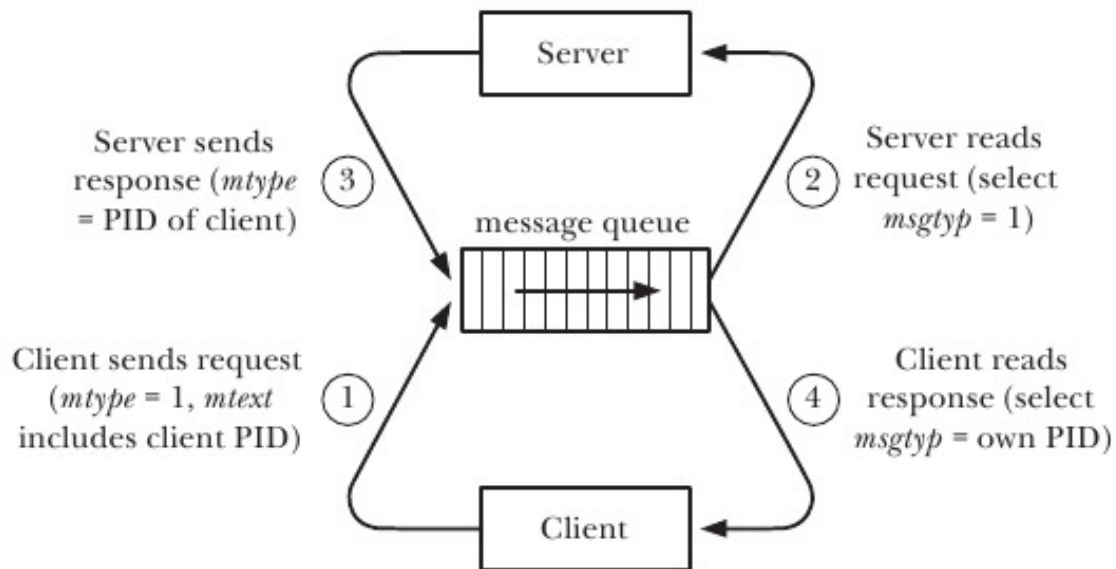response(s)

# Figure 1.21 - POSIX SYS V semaphore IPC OBJECT

```
struct semid_ds {
    struct ipc_perm sem_perm;       /* Ownership and permissions */
    time_t          sem_otime;      /* Time of last semop() */
    time_t          sem_ctime;      /* Time of last change */
    unsigned long   sem_nsems;      /* Number of semaphores in set */
};
```

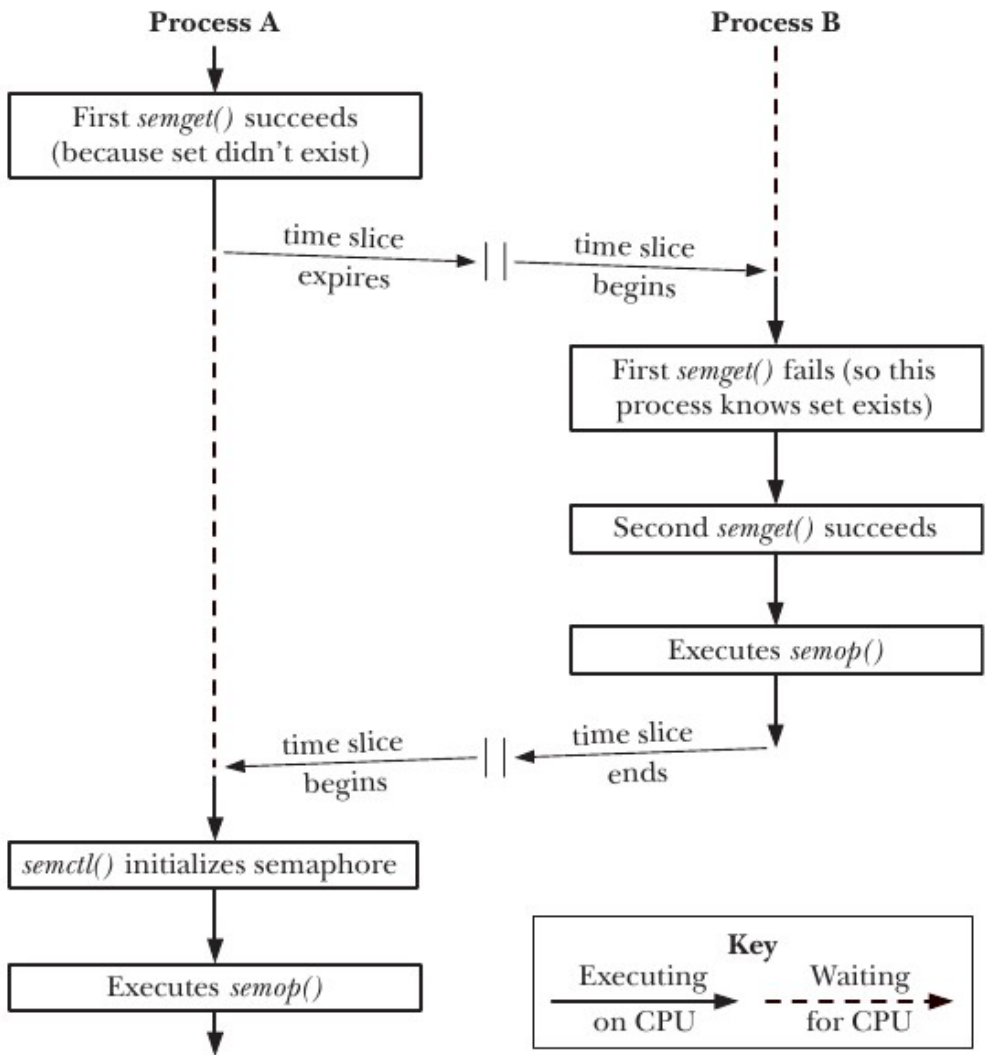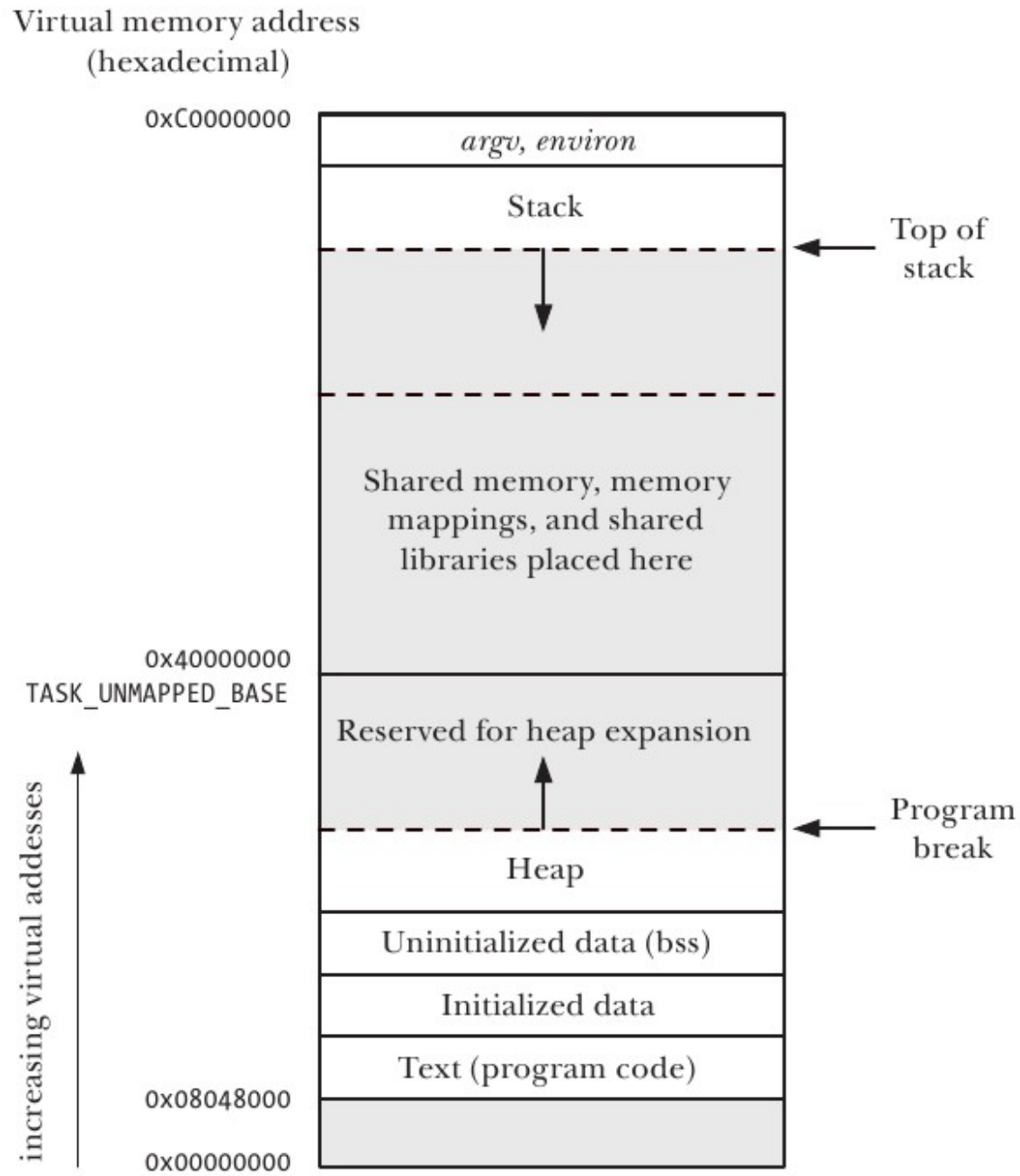# Figure 1.22 - subtle SYS V semaphore initialization problem

**Process A**                                    **Process B**

First *semget()* succeeds
(because set didn't exist)

time slice                  time slice
expires                     begins

First *semget()* fails (so this
process knows set exists)

Second *semget()* succeeds

Executes *semop()*

time slice                  time slice
begins                      ends

*semctl()* initializes semaphore

Executes *semop()*

**Key**

Executing          Waiting
on CPU            for CPU

# Figure 1.23 - shared memory mappings

Virtual memory address
(hexadecimal)

| | |
|---|---|
| 0xC0000000 | *argv, environ* |
| | Stack |
| | Shared memory, memory mappings, and shared libraries placed here |
| 0x40000000 TASK_UNMAPPED_BASE | Reserved for heap expansion |
| | Heap |
| | Uninitialized data (bss) |
| | Initialized data |
| | Text (program code) |
| 0x08048000 | |
| 0x00000000 | |

Top of stack

Program break

increasing virtual addesses

**Figure 1.24  -  SYSV shared memory IPC object**

```
struct shmid_ds {
    struct ipc_perm shm_perm;    /* Ownership and permissions */
    size_t   shm_segsz;          /* Size of segment in bytes */
    time_t   shm_atime;          /* Time of last shmat() */
    time_t   shm_dtime;          /* Time of last shmdt() */
    time_t   shm_ctime;          /* Time of last change */
    pid_t    shm_cpid;           /* PID of creator */
    pid_t    shm_lpid;           /* PID of last shmat() / shmdt() */
    shmatt_t shm_nattch;         /* Number of currently attached processes */
};
```

**Figure 1.25  -   memory mappings**

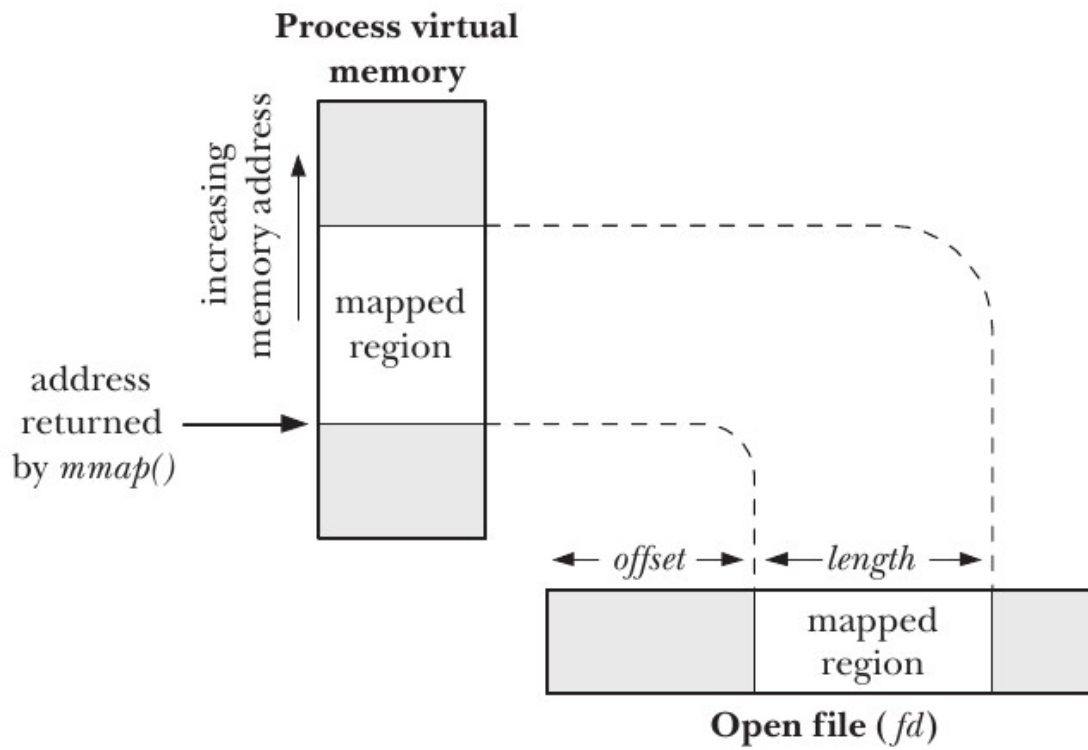| Visibility of modifications | Mapping type | |
|---|---|---|
| | **File** | **Anonymous** |
| **Private** | Initializing memory from contents of file | Memory allocation |
| **Shared** | Memory-mapped I/O; sharing memory between processes (IPC) | Sharing memory between processes (IPC) |

**Figure 1.26 - file memory – mapping**
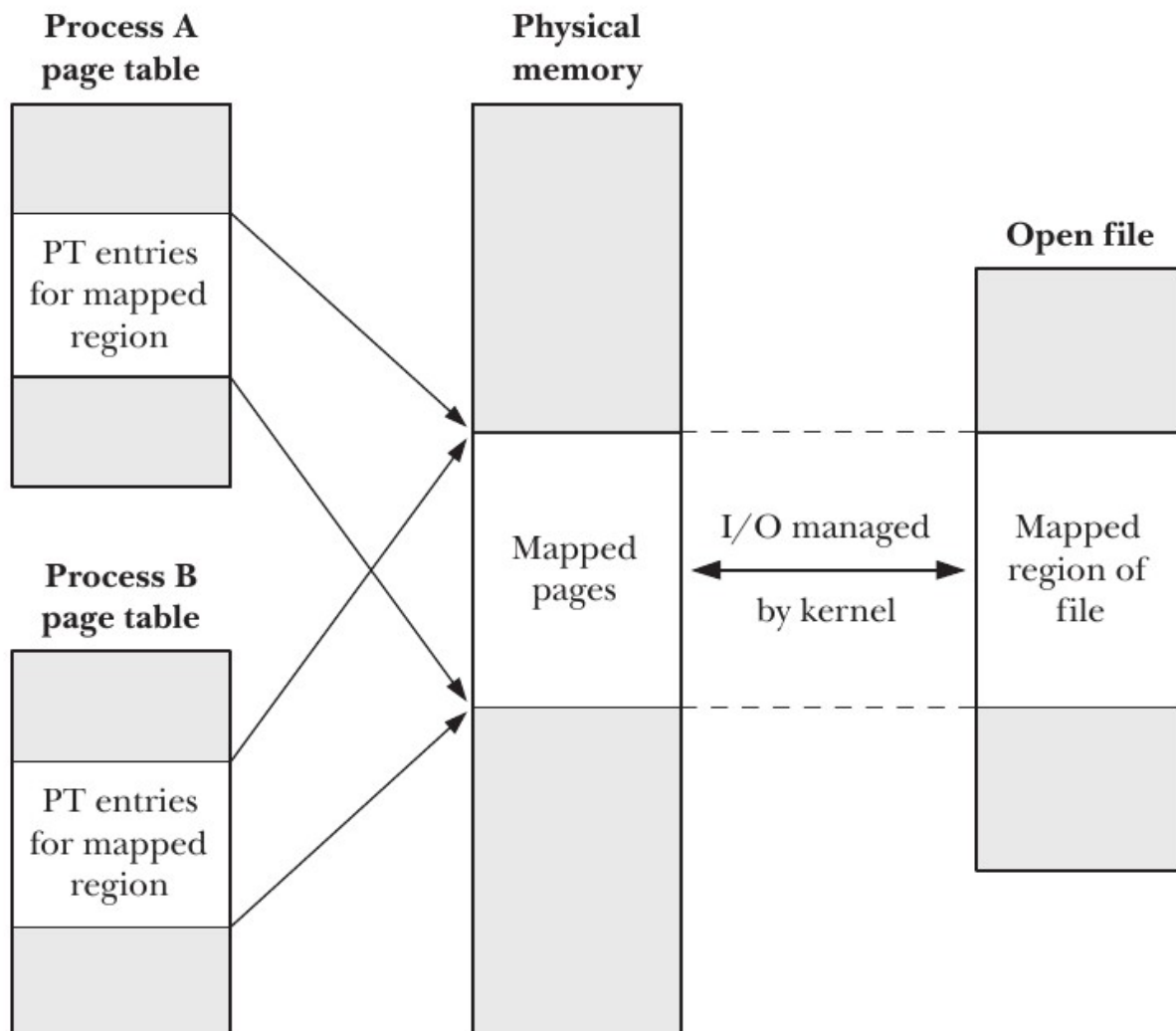
# Figure 1.27 - shared file memory mapping

**Process A**
**page table**

PT entries
for mapped
region

**Process B**
**page table**

PT entries
for mapped
region

**Physical**
**memory**

Mapped
pages

I/O managed

by kernel

**Open file**

Mapped
region of
file

# Figure 1.28 - file mapping access errors

*mmap(0, 6000, prot, MAP_SHARED, fd, 0);*

byte offset:   0                                    5999 6000      8191 8192

**Memory region**

| requested size of mapping | remainder of page |
|---|---|

← accessible, mapped to file → | ← references yield SIGSEGV →

**Mapped file (9500 bytes)**

| actual mapped region of file | unmapped |
|---|---|

file offset:   0                                              8191 8192        9499

---

*mmap(0, 8192, prot, MAP_SHARED, fd, 0);*

byte offset:   0        2199 2200      4095 4096                 8191 8192

**Memory region**

| | remainder of page (0s) | |
|---|---|---|

← accessible, mapped to file → | ← accessible, not mapped to file → | ← references yield SIGBUS → | ← references yield SIGSEGV →

**Mapped file (2200 bytes)**

file offset:   0           2199

**Figure 1.29 - mmap and virtual memory issues**

| overcommit_memory value | MAP_NORESERVE specified in *mmap()* call? | |
| --- | --- | --- |
| | **No** | **Yes** |
| 0 | Deny obvious overcommits | Allow overcommits |
| 1 | Allow overcommits | Allow overcommits |
| 2 (since Linux 2.6) | Strict overcommitting | |

# Figure 1.30 - different types of signals

| Name | Signal number | Description | SUSv3 | Default |
|---|---|---|---|---|
| SIGABRT | 6 | Abort process | • | core |
| SIGALRM | 14 | Real-time timer expired | • | term |
| SIGBUS | 7 (SAMP=10) | Memory access error | • | core |
| SIGCHLD | 17 (SA=20, MP=18) | Child terminated or stopped | • | ignore |
| SIGCONT | 18 (SA=19, M=25, P=26) | Continue if stopped | • | cont |
| SIGEMT | undef (SAMP=7) | Hardware fault | | term |
| SIGFPE | 8 | Arithmetic exception | • | core |
| SIGHUP | 1 | Hangup | • | term |
| SIGILL | 4 | Illegal instruction | • | core |
| SIGINT | 2 | Terminal interrupt | • | term |
| SIGIO / SIGPOLL | 29 (SA=23, MP=22) | I/O possible | • | term |
| SIGKILL | 9 | Sure kill | • | term |
| SIGPIPE | 13 | Broken pipe | • | term |
| SIGPROF | 27 (M=29, P=21) | Profiling timer expired | • | term |
| SIGPWR | 30 (SA=29, MP=19) | Power about to fail | | term |
| SIGQUIT | 3 | Terminal quit | • | core |
| SIGSEGV | 11 | Invalid memory reference | • | core |
| SIGSTKFLT | 16 (SAM=undef, P=36) | Stack fault on coprocessor | | term |
| SIGSTOP | 19 (SA=17, M=23, P=24) | Sure stop | • | stop |
| SIGSYS | 31 (SAMP=12) | Invalid system call | • | core |
| SIGTERM | 15 | Terminate process | • | term |
| SIGTRAP | 5 | Trace/breakpoint trap | • | core |
| SIGTSTP | 20 (SA=18, M=24, P=25) | Terminal stop | • | stop |
| SIGTTIN | 21 (M=26, P=27) | Terminal read from BG | • | stop |
| SIGTTOU | 22 (M=27, P=28) | Terminal write from BG | • | stop |
| SIGURG | 23 (SA=16, M=21, P=29) | Urgent data on socket | • | ignore |
| SIGUSR1 | 10 (SA=30, MP=16) | User-defined signal 1 | • | term |
| SIGUSR2 | 12 (SA=31, MP=17) | User-defined signal 2 | • | term |
| SIGVTALRM | 26 (M=28, P=20) | Virtual timer expired | • | term |
| SIGWINCH | 28 (M=20, P=23) | Terminal window size change | | ignore |
| SIGXCPU | 24 (M=30, P=33) | CPU time limit exceeded | • | core |
| SIGXFSZ | 25 (M=31, P=34) | File size limit exceeded | • | core |

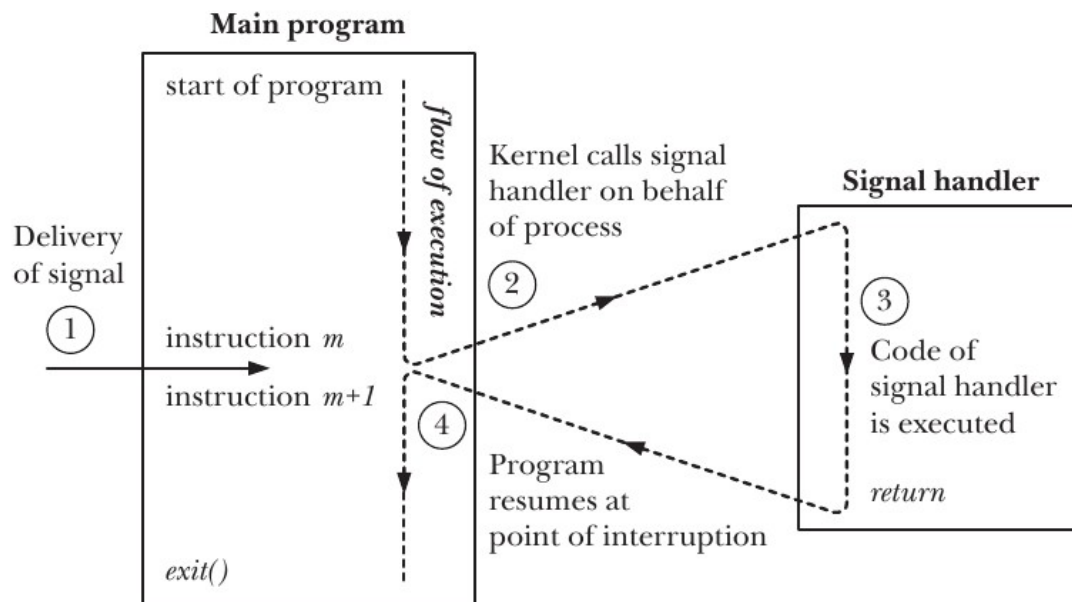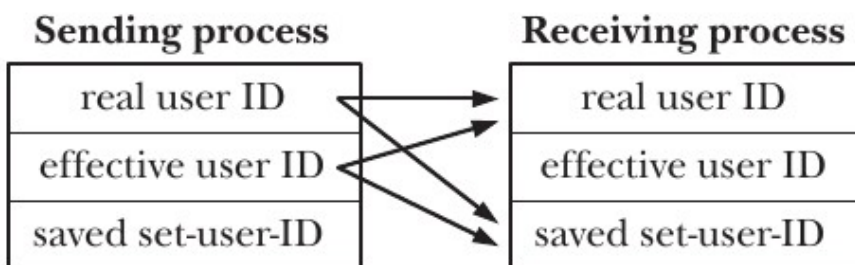# Figure 1.31 - signal handling and process/program flow

**Main program**

start of program

*flow of execution*

Kernel calls signal
handler on behalf
of process

**Signal handler**

Delivery
of signal

② 3

① instruction *m* Code of
signal handler
is executed

instruction *m+1*

④ *return*

Program
resumes at
point of interruption

*exit()*

**Figure 1.32 - permission to send a signal**



Sending process                  Receiving process

| real user ID | | real user ID |
|---|---|
| effective user ID | effective user ID |
| saved set-user-ID | saved set-user-ID |

indicates that if IDs match,
then sender has permission
to send a signal to receiver

# Figure 1.33  -  async signal safe library functions

| | | |
|---|---|---|
| _Exit() (v3) | getpid() | sigdelset() |
| _exit() | getppid() | sigemptyset() |
| abort() (v3) | getsockname() (v3) | sigfillset() |
| accept() (v3) | getsockopt() (v3) | sigismember() |
| access() | getuid() | signal() (v2) |
| aio_error() (v2) | kill() | sigpause() (v2) |
| aio_return() (v2) | link() | sigpending() |
| aio_suspend() (v2) | listen() (v3) | sigprocmask() |
| alarm() | lseek() | sigqueue() (v2) |
| bind() (v3) | lstat() (v3) | sigset() (v2) |
| cfgetispeed() | mkdir() | sigsuspend() |
| cfgetospeed() | mkfifo() | sleep() |
| cfsetispeed() | open() | socket() (v3) |
| cfsetospeed() | pathconf() | sockatmark() (v3) |
| chdir() | pause() | socketpair() (v3) |
| chmod() | pipe() | stat() |
| chown() | poll() (v3) | symlink() (v3) |
| clock_gettime() (v2) | posix_trace_event() (v3) | sysconf() |
| close() | pselect() (v3) | tcdrain() |
| connect() (v3) | raise() (v2) | tcflow() |
| creat() | read() | tcflush() |
| dup() | readlink() (v3) | tcgetattr() |
| dup2() | recv() (v3) | tcgetpgrp() |
| execle() | recvfrom() (v3) | tcsendbreak() |
| execve() | recvmsg() (v3) | tcsetattr() |
| fchmod() (v3) | rename() | tcsetpgrp() |
| fchown() (v3) | rmdir() | time() |
| fcntl() | select() (v3) | timer_getoverrun() (v2) |
| fdatasync() (v2) | sem_post() (v2) | timer_gettime() (v2) |
| fork() | send() (v3) | timer_settime() (v2) |
| fpathconf() (v2) | sendmsg() (v3) | times() |
| fstat() | sendto() (v3) | umask() |
| fsync() (v2) | setgid() | uname() |
| ftruncate() (v3) | setpgid() | unlink() |
| getegid() | setsid() | utime() |
| geteuid() | setsockopt() (v3) | wait() |
| getgid() | setuid() | waitpid() |
| getgroups() | shutdown() (v3) | write() |
| getpeername() (v3) | sigaction() | |
| getpgrp() | sigaddset() | |

# Figure 1.34 - struct sigaction { } and supported handlers

```
struct sigaction {
    union {
        void (*sa_handler)(int);
        void (*sa_sigaction)(int, siginfo_t *, void *);
    } __sigaction_handler;
    sigset_t   sa_mask;
    int        sa_flags;
    void       (*sa_restorer)(void);
};

/* Following defines make the union fields look like simple fields
   in the parent structure */

#define sa_handler __sigaction_handler.sa_handler
#define sa_sigaction __sigaction_handler.sa_sigaction
```

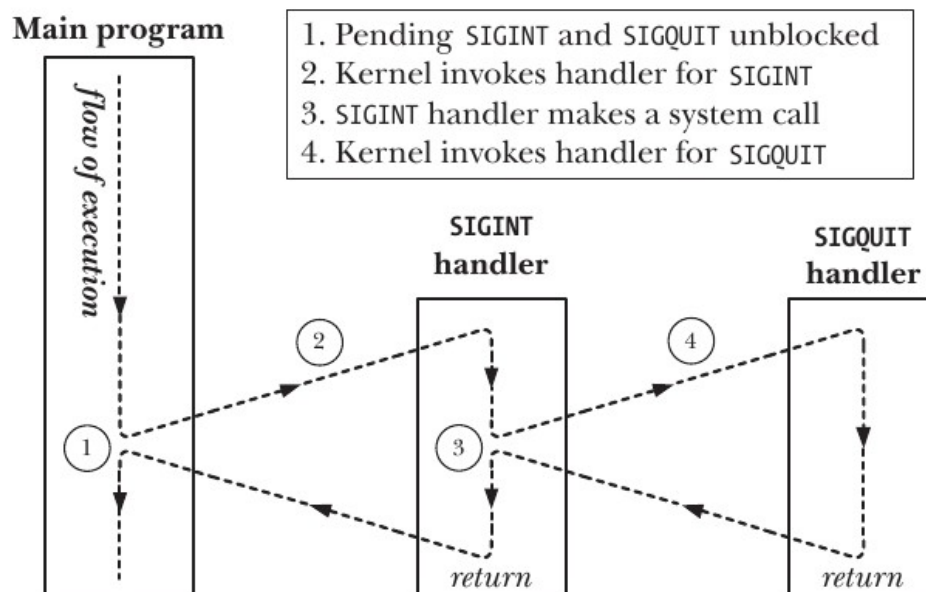**Figure 1.35 - a special case of nested signal handling**

**Main program**

*flow of execution*

1. Pending SIGINT and SIGQUIT unblocked
2. Kernel invokes handler for SIGINT
3. SIGINT handler makes a system call
4. Kernel invokes handler for SIGQUIT

**SIGINT handler**

**SIGQUIT handler**

(1)  (2)  (3)  (4)

*return*  *return*

# Figure 1.36  -  threads in process address space

Virtual memory address
(hexadecimal)

0xC0000000

| |
|---|
| *argv, environ* |
| Stack for main thread |
| |
| Stack for thread 3 |
| Stack for thread 2 |
| Stack for thread 1 |
| Shared libraries, shared memory |
| |
| Heap |
| Uninitialized data (bss) |
| Initialized data |
| Text (program code) |
| |

0x40000000
TASK_UNMAPPED_BASE

increasing virtual addesses

← thread 3 executing here
← main thread executing here
← thread 1 executing here
← thread 2 executing here

0x08048000

0x00000000

# Figure 1.37  -   non thread-safe library functions

| | | | |
|---|---|---|---|
| asctime() | fcvt() | getpwnam() | nl_langinfo() |
| basename() | ftw() | getpwuid() | ptsname() |
| catgets() | gcvt() | getservbyname() | putc_unlocked() |
| crypt() | getc_unlocked() | getservbyport() | putchar_unlocked() |
| ctime() | getchar_unlocked() | getservent() | putenv() |
| dbm_clearerr() | getdate() | getutxent() | pututxline() |
| dbm_close() | getenv() | getutxid() | rand() |
| dbm_delete() | getgrent() | getutxline() | readdir() |
| dbm_error() | getgrgid() | gmtime() | setenv() |
| dbm_fetch() | getgrnam() | hcreate() | setgrent() |
| dbm_firstkey() | gethostbyaddr() | hdestroy() | setkey() |
| dbm_nextkey() | gethostbyname() | hsearch() | setpwent() |
| dbm_open() | gethostent() | inet_ntoa() | setutxent() |
| dbm_store() | getlogin() | l64a() | strerror() |
| dirname() | getnetbyaddr() | lgamma() | strtok() |
| dlerror() | getnetbyname() | lgammaf() | ttyname() |
| drand48() | getnetent() | lgammal() | unsetenv() |
| ecvt() | getopt() | localeconv() | wcstombs() |
| encrypt() | getprotobyname() | localtime() | wctomb() |
| endgrent() | getprotobynumber() | lrand48() | |
| endpwent() | getprotoent() | mrand48() | |
| endutxent() | getpwent() | nftw() | |

# Figure 1.38 - Threads in Linux and Ids



Process with PID 2001

| Thread A | Thread B | Thread C | Thread D |
|----------|----------|----------|----------|
| PPID=1900 | PPID=1900 | PPID=1900 | PPID=1900 |
| TGID=2001 | TGID=2001 | TGID=2001 | TGID=2001 |
| TID=2001 | TID=2002 | TID=2003 | TID=2004 |

Thread group leader (TID matches TGID)