# System Programming

## Week 14: Interprocess Commnunication

---

Seongjin Lee

Updated: 2016-12-07
11_ipc


insight@hanyang.ac.kr
http://esos.hanyang.ac.kr
Esos Lab. Hanyang University
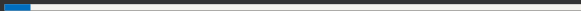
## Table of contents

## Introduction

This chapter covers other techniques for processes to communicate with one another: interprocess communication (IPC)

- ○ pipes
- ○ FIFOs
- ○ Message Queues
- ○ Shared Memory
- ○ Semaphores

# Pipes

## Pipes

Pipes are the oldest form of UNIX system IPC

Pipes have two limitations

1. They have been half duplex (most commonly used), now provide full-duplex pipes (but never assume it is full-duplex)
2. Pipes can be used only between processes that have a common ancestor

## Pipes

```
#include <unistd.h>
int pipe(int fd[2]);
// Returns: 0 if OK, 1 on error
```

Two file descriptors are returned through the fd argument:

○ fd[0] is open for reading,

○ fd[1] is open for writing.

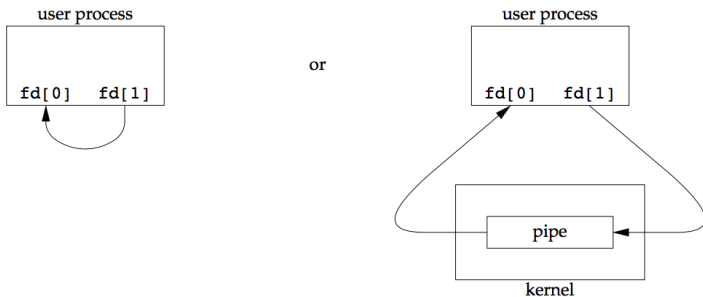○ The output of fd[1] is the input for fd[0].

# Pipes



**Figure:** Two ways to veiw a half-duplex pipe

# Pipes

Normally, the process that calls pipe then calls fork, creating an IPC channel from the parent to the child, or vice versa.
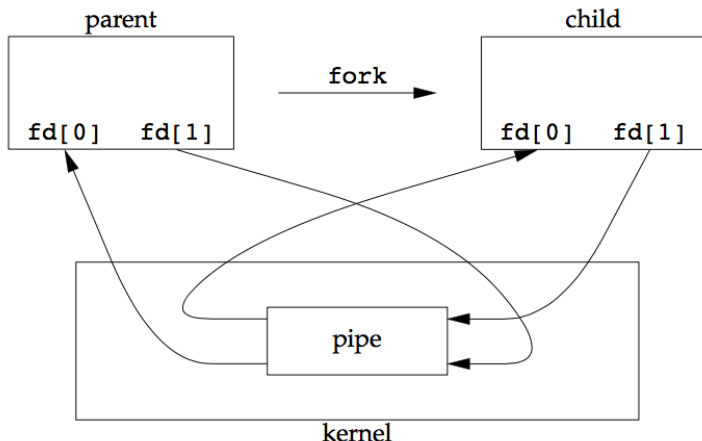


**Figure:** Half-duplex pipe after a fork

# Pipes

Normally, the process that calls pipe then calls fork, creating an IPC channel from the parent to the child, or vice versa.



**Figure:** Pipe from parent to child
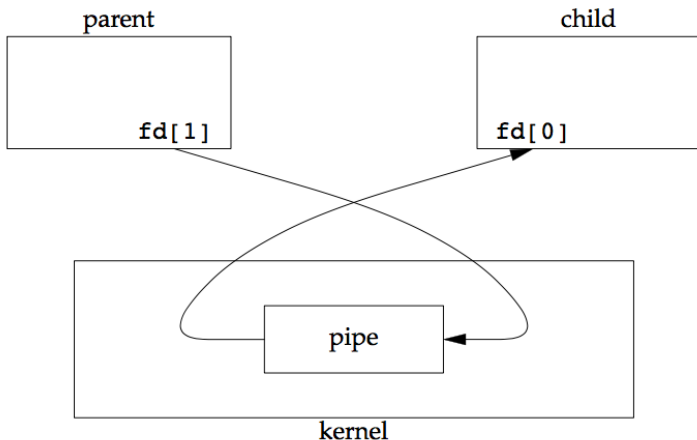
## Pipes

When one end of a pipe is closed, two rules apply

1. If we read from a pipe whose write end has been closed, read returns 0 to indicate an end of file after all the data has been read.
2. If we write to a pipe whose read end has been closed, the signal SIGPIPE is generated.

When we're writing to a pipe (or FIFO), the constant PIPE_BUF specifies the kernel's pipe buffer size.

In this example, we call read and write directly on the pipe descriptors.

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <unistd.h>
4
5
6   #define MAXLINE 4096    /* max line length */
7
8   int
9   main(void)
10  {
11     int  n;
12     int  fd[2];
13     pid_t pid;
14     char line[MAXLINE];
15
16     if (pipe(fd) < 0){
17        fprintf(stderr, "pipe error");
18        exit(1);
19     }
20     if ((pid = fork()) < 0) {
```

```
21        fprintf(stderr, "fork error");
22        exit(1);
23    } else if (pid > 0) {  /* parent */
24        close(fd[0]);
25        write(fd[1], "hello world\n", 12);
26    } else {      /* child */
27        close(fd[1]);
28        n = read(fd[0], line, MAXLINE);
29        write(STDOUT_FILENO, line, n);
30    }
31    exit(0);
32 }
```

○ Before calling fork, we create a pipe.

○ After the fork, the parent closes its read end, and the child closes its write end.

○ The child then calls dup2 to have its standard input be the read end of the pipe.

○ When the pager program is executed, its standard input will be the read end of the pipe.

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <unistd.h>
4   #include <string.h>
5   #include <sys/wait.h>
6
7   #define MAXLINE 4096    /* max line length */
8   #define DEF_PAGER "/usr/bin/more"  /* default pager program */
9
10  int
11  main(int argc, char *argv[])
12  {
13     int  n;
14     int  fd[2];
15     pid_t pid;
16     char *pager, *argv0;
17     char line[MAXLINE];
18     FILE *fp;
19
20     if (argc != 2){
21        fprintf(stderr, "usage: a.out <pathname>");
22        exit(1);
23     }
```

```
24
25     if ((fp = fopen(argv[1], "r")) == NULL){
26        fprintf(stderr, "can't open %s", argv[1]);
27        exit(1);
28     }
29     if (pipe(fd) < 0){
30        fprintf(stderr, "pipe error");
31        exit(1);
32     }
33
34     if ((pid = fork()) < 0) {
35        fprintf(stderr, "fork error");
36        exit(1);
37     } else if (pid > 0) {        /* parent */
38        close(fd[0]);  /* close read end */
39
40        /* parent copies argv[1] to pipe */
41        while (fgets(line, MAXLINE, fp) != NULL) {
42           n = strlen(line);
43           if (write(fd[1], line, n) != n)
44              fprintf(stderr, "write error to pipe");
45        }
46        if (ferror(fp))
```

```
47          fprintf(stderr, "fgets error");
48
49       close(fd[1]); /* close write end of pipe for reader */
50
51       if (waitpid(pid, NULL, 0) < 0){
52          fprintf(stderr, "waitpid error");
53          exit(1);
54       }
55       exit(0);
56    } else {          /* child */
57       close(fd[1]); /* close write end */
58       if (fd[0] != STDIN_FILENO) {
59          if (dup2(fd[0], STDIN_FILENO) != STDIN_FILENO){
60             fprintf(stderr, "dup2 error to stdin");
61             exit(1);
62          }
63          close(fd[0]); /* don't need this after dup2 */
64       }
65
66       /* get arguments for execl() */
67       if ((pager = getenv("PAGER")) == NULL)
68          pager = DEF_PAGER;
69       if ((argv0 = strrchr(pager, '/')) != NULL)
```

```
70          argv0++;  /* step past rightmost slash */
71        else
72          argv0 = pager; /* no slash in pager */
73
74        if (execl(pager, argv0, (char *)0) < 0){
75          fprintf(stderr, "execl error for %s", pager);
76          exit(1);
77        }
78      }
79    exit(0);
80  }
```

### Defensive programming measure

Whenever we call dup2 and close to duplicate one descriptor onto another, we'll always compare the descriptors first.

# PIPES : popen AND pclose FUNCTIONS

## popen and pclose Functions

Standard I/O library to create a pipe to another process to either read or write

```
#include <stdio.h>
FILE *popen(const char *cmdstring, const char *type);

// Returns: file pointer if OK, NULL on error

int pclose(FILE *fp);
// Returns: termination status of cmdstring, or 1 on error
```

The function popen does a fork and exec to execute the *cmdstring* and returns a standard I/O file pointer.

If *type* is "r", the file pointer is connected to the standard output of *cmdstring*



**Figure:** Result of fp = popen(cmdstring, ''r'')

If type is "w", the file pointer is connected to the standard input of *cmdstring*

## popen and pclose Functions

The pclose function closes the standard I/O stream, waits for the command to terminate, and returns the termination status of the shell.

## popen and pclose Example I

redo of the program using popen: codes/popen2.c

The shell command ${PAGER:-more} says to use the value of the shell variable PAGER if it is defined and non-null; otherwise, use the string more.

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <unistd.h>
4   #include <sys/wait.h>
5
6   #define MAXLINE 4096    /* max line length */
7   #define PAGER "${PAGER:-more}" /* environment variable, or default */
8
9   int
10  main(int argc, char *argv[])
11  {
12      char line[MAXLINE];
13      FILE *fpin, *fpout;
14
15      if (argc != 2){
```

```
16        fprintf(stderr, "usage: a.out <pathname>");
17        exit(1);
18    }
19    if ((fpin = fopen(argv[1], "r")) == NULL){
20        fprintf(stderr, "can't open %s", argv[1]);
21        exit(1);
22    }
23
24    if ((fpout = popen(PAGER, "w")) == NULL){
25        fprintf(stderr, "popen error");
26        exit(1);
27    }
28
29    /* copy argv[1] to pager */
30    while (fgets(line, MAXLINE, fpin) != NULL) {
31        if (fputs(line, fpout) == EOF){
32            fprintf(stderr, "fputs error to pipe");
33            exit(1);
34        }
35    }
36    if (ferror(fpin)){
37        fprintf(stderr, "fgets error");
38        exit(1);
```

```
39        }
40        if (pclose(fpout) == -1){
41            fprintf(stderr, "pclose error");
42            exit(1);
43        }
44
45        exit(0);
46    }
```

## popen and pclose Example

codes/myuclc.c : a filter program that changes upper cases to lower cases.

```c
1  #include <stdio.h>
2  #include <string.h>
3  #include <stdlib.h>
4  #include <ctype.h>
5
6  int
7  main(void)
8  {
9     int  c;
10
11    while ((c = getchar()) != EOF) {
12       if (isupper(c))
13          c = tolower(c);
14       else if (islower(c))
15          c = toupper(c);
16       if (putchar(c) == EOF){
17          fprintf(stderr, "output error");
18          exit(1);
19       }
20       if (c == '\n')
21          fflush(stdout);
22    }
23    exit(0);
24
```

codes/popen1.c

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <unistd.h>
4   #include <string.h>
5   #include <sys/wait.h>
6
7   #define MAXLINE 4096    /* max line length */
8
9   int
10  main(void)
11  {
12     char line[MAXLINE];
13     FILE *fpin;
14
15     if ((fpin = popen("./myuclc", "r")) == NULL){
16        fprintf(stderr, "popen error");
17        exit(1);
18     }
19
20      /* type 111 to exit the program */
```
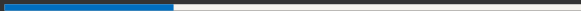
```
21      while (strcmp(line, "111\n") != 0){
22         fputs("prompt> ", stdout);
23         fflush(stdout);
24         if (fgets(line, MAXLINE, fpin) == NULL) /* read from pipe */
25            break;
26         if (fputs(line, stdout) == EOF)
27            fprintf(stderr, "fputs error to pipe");
28      }
29      if (pclose(fpin) == -1){
30         fprintf(stderr, "pclose error");
31         exit(1);
32      }
33      putchar('\n');
34      exit(0);
35   }
```

# COPROCESSES

# Coprocesses

A Unix system filter is a program that reads from standard input and writes to standard output.

Filters are normally connected linearly in shell pipelines.

A filter becomes a *coprocess* when the same program generates the filter's input and reads the filter output
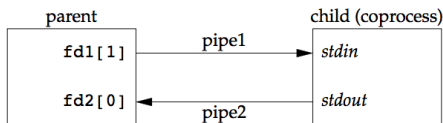


**Figure:** Driving a coprocess by writing its standard input and reading its standard output

## Coprocess Example I

A simple coprocess that reads two numbers from its standard input, computes their sum, and writes the sum to its standard output.

`codes/add2.c; make add2`

```c
1   #include <stdio.h>
2   #include <unistd.h>
3   #include <stdlib.h>
4   #include <string.h>
5
6   #define MAXLINE 4096   /* max line length */
7
8   int
9   main(void)
10  {
11      int  n, int1, int2;
12      char line[MAXLINE];
13
14      while ((n = read(STDIN_FILENO, line, MAXLINE)) > 0) {
15          line[n] = 0;  /* null terminate */
16          if (sscanf(line, "%d%d", &int1, &int2) == 2) {
17              sprintf(line, "%d\n", int1 + int2);
18              n = strlen(line);
19              if (write(STDOUT_FILENO, line, n) != n){
```

```
20            fprintf(stderr, "write error");
21            exit(1);
22        }
23    } else {
24        if (write(STDOUT_FILENO, "invalid args\n", 13) != 13){
25            fprintf(stderr, "write error");
26            exit(1);
27        }
28    }
29    }
30    exit(0);
31 }
```

Here, we create two pipes, with the parent and the child closing the ends they don't need.

## Coprocess Example III

```
codes/pipe4.c; make pipe4
```

```c
1   #include <stdio.h>
2   #include <unistd.h>
3   #include <stdlib.h>
4   #include <string.h>
5
6   #define MAXLINE 4096   /* max line length */
7
8   static void sig_pipe(int);  /* our signal handler */
9
10  int
11  main(void)
12  {
13    int  n, fd1[2], fd2[2];
14    pid_t pid;
15    char line[MAXLINE];
16
17    if (signal(SIGPIPE, sig_pipe) == SIG_ERR){
18      fprintf(stderr, "signal error");
19      exit(1);
20    }
21
22    if (pipe(fd1) < 0 || pipe(fd2) < 0){
23      fprintf(stderr, "pipe error");
24      exit(1);
```

```
25      }
26
27      if ((pid = fork()) < 0) {
28         fprintf(stderr, "fork error");
29         exit(1);
30      } else if (pid > 0) {        /* parent */
31         close(fd1[0]);
32         close(fd2[1]);
33
34         while (fgets(line, MAXLINE, stdin) != NULL) {
35            if(strcmp(line, "quit\n") == 0 )
36               exit(0);
37            n = strlen(line);
38            if (write(fd1[1], line, n) != n){
39               fprintf(stderr, "write error to pipe");
40               exit(1);
41            }
42            if ((n = read(fd2[0], line, MAXLINE)) < 0){
43               fprintf(stderr, "read error from pipe");
44               exit(1);
45            }
46            if (n == 0) {
47               fprintf(stderr, "child closed pipe");
48               break;
49            }
50            line[n] = 0; /* null terminate */
```

```
51          if (fputs(line, stdout) == EOF){
52              fprintf(stderr, "fputs error");
53              exit(1);
54          }
55      }
56
57      if (ferror(stdin)){
58          fprintf(stderr, "fgets error on stdin");
59          exit(1);
60      }
61      exit(0);
62  } else {          /* child */
63      close(fd1[1]);
64      close(fd2[0]);
65      if (fd1[0] != STDIN_FILENO) {
66          if (dup2(fd1[0], STDIN_FILENO) != STDIN_FILENO){
67              fprintf(stderr, "dup2 error to stdin");
68              exit(1);
69          }
70          close(fd1[0]);
71      }
72
73      if (fd2[1] != STDOUT_FILENO) {
74          if (dup2(fd2[1], STDOUT_FILENO) != STDOUT_FILENO){
75              fprintf(stderr, "dup2 error to stdout");
76              exit(1);
```

```
77            }
78            close(fd2[1]);
79        }
80        if (execl("./add2", "add2", (char *)0) < 0){
81            fprintf(stderr, "execl error");
82            exit(1);
83        }
84    }
85    exit(0);
86 }
87
88 static void
89 sig_pipe(int signo)
90 {
91    printf("SIGPIPE caught\n");
92    exit(1);
93 }
```

If we kill the add2 coprocess while the program is waiting for our input and then enter two numbers, the signal handler is invoked when the program writes to the pipe that has no reader.

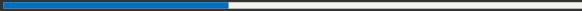From another terminal type the following

```
$ ps | grep add2
40467 ttys000 0:00.00 add2
$ kill 40467
```

then input two numbers to pipe4

# FIFOs

## FIFOs

FIFOs are sometimes called named pipes.

- ○ FIFOs are sometimes called named pipes.
- ○ Unnamed pipes can be used only between related processes when a common ancestor has created the pipe.
- ○ With FIFOs, unrelated processes can exchange data.

## FIFOs

Creating a FIFO is similar to creating a file. Indeed, the pathname for a FIFO exists in the file system

```
 #include <sys/stat.h>
int mkfifo(const char *path, mode_t mode);
int mkfifoat(int fd, const char *path, mode_t mode);
// Both return: 0 if OK, 1 on error
```

# FIFOs

- *mode* argument is the same as for the open function (Section 3.3)
- The rules for the user and group ownership of the new FIFO are the same as we described in Section 4.6
- `mkfifoat` creates a FIFO in a location relative to the directory represented by the `fd` file descriptor argument.
  - *path* == absolute pathname: `mkfifoat` function behaves like the `mkfifo` function.
  - *path* == relative pathname && `fd` == valid directory: the pathname is evaluated relative to this directory.
  - *path* == relative pathname && `fd` == `AT_FDCWD`: starts incurrent working directory, and mkfifoat behaves like mkfifo.

Once we have used `mkfifo` or `mkfifoat` to create a FIFO, we open it using open. Indeed, the normal file I/O functions (e.g., close, read, write, unlink) all work with FIFOs.

## FIFOs

When we open a FIFO, the nonblocking flag (O_NONBLOCK) affects what happens.

- ○ without O_NONBLOCK:
  - ○ open for read-only blocks until some other process opens the FIFO for writing
  - ○ open for write-only blocks until some other process opens the FIFO for reading.
- ○ with O_NONBLOCK:
  - ○ open for read-only returns immediately.
  - ○ open for write-only returns -1 with errno set to ENXIO if no process has the FIFO open for reading.

## FIFO Use Cases

There are two uses for FIFO

1. FIFOs are used by shell commands to pass data from one shell pipeline to another *without creating intermediate temporary files*.
2. FIFOs are used as rendezvous points in client–server applications to pass data between the clients and the servers.

# Example - Using FIFOs to Duplicate Output Streams

FIFOs can be used to duplicate an output stream in a series of shell commands.

A FIFO has a name, so it can be used for nonlinear connections.



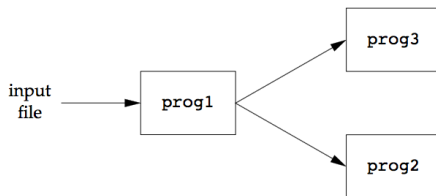**Figure:** Procedure that processes a filtered input stream twice

```
mkfifo fifo1
prog3 < fifo1 &
prog1 < infile | tee fifo1 | prog2
```



**Figure:** Using a FIFO and tee to send a stream to two different processes

**Figure:** Client-server communication using FIFOs

codes/create_fifo.c

make `create_fifo` to create a FIFO for server and client

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <sys/stat.h>
4
5  int main()
6  {
7    int server, client;
8
9    /* create a FIFO for server */
10   server = mkfifo("FIFO-server", 0666);
11   if(server < 0) {
12     fprintf(stderr, "Failed to create a FIFO for server");
13     exit(-1);
14   }
15
16   /* create a FIFO for client */
17   client = mkfifo("FIFO-client", 0666);
```

```
18    if(client < 0) {
19        fprintf(stderr, "Failed to create a FIFO for client");
20        exit(-1);
21    }
22
23    printf("FIFO for server and child created successfuly");
24
25    return 0;
26 }
```

codes/server.c

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <unistd.h>
4   #include <fcntl.h>
5   #include <sys/stat.h>
6
7   int main()
8   {
9       int FIFO_server,FIFO_client;
10      int choice;
11      char *buf;
12
13      FIFO_server = open("FIFO-server", O_RDWR);
14      if(FIFO_server < 0) {
15          fprintf(stderr, "Error opening server side FIFO file");
16          exit(-1);
17      }
18
19      read(FIFO_server, &choice, sizeof(int));
20
21      sleep(1);
22
23      FIFO_client = open("FIFO-client", O_RDWR);
24      if(FIFO_client < 0) {
```

```
25        fprintf(stderr, "Error opening client side FIFO file");
26        exit(-1);
27     }
28
29     switch(choice) {
30        case 1:
31           buf="Linux";
32           break;
33        case 2:
34           buf="Debian";
35           break;
36        case 3:
37           buf="4.0";
38     }
39
40     write(FIFO_client,buf,10*sizeof(char));
41     printf("\n### Data sent to client ###\n");
42
43     close(FIFO_server);
44     close(FIFO_client);
45     return 0;
46  }
```

codes/client.c

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <unistd.h>
4   #include <fcntl.h>
5
6   int main()
7   {
8       char *buf;
9       int choice = 1;
10      int fifo_server, fifo_client;
11
12      printf("Chose one of following options to send to the server\n");
13      printf("\t\t 1 for the name of the OS \n \
14               2 for the name of distribution \n \
15               3 for the version of Kernel \n");
16      printf("Your choice: ");
17      scanf("%d",&choice);
18
19      fifo_server=open("FIFO-server",O_RDWR);
20      if(fifo_server < 0) {
21          fprintf(stderr, "Error in opening file");
22          exit(-1);
23       }
24
```

```
25      write(fifo_server, &choice, sizeof(int));
26
27      fifo_client=open("FIFO-client", O_RDWR);
28
29      if(fifo_client < 0) {
30          fprintf(stderr, "Error in opening file");
31          exit(-1);
32       }
33
34      buf=malloc(10*sizeof(char));
35      read (fifo_client, buf, 10*sizeof(char));
36      printf("\n ### Server replied with %s ###\n", buf);
37
38      close(fifo_server);
39      close(fifo_client);
40      return 0;
41  }
```

# XSI IPC

# XSI IPC : General Background

# XSI IPC

There are three types of IPC that we call XSI IPC

1. message queues
2. semaphores
3. shared memory

# Identifiers and Keys

Each IPC *structure* in the kernel is referred to by a non-negative integer *identifier*

- ○ to send a message to or fetch a message from a message queue, for example, all we need to know is the identifier for the queue
- ○ when a given IPC structure is created and then removed
  - ○ identifier associated with that structure conitunally increases
- ○ identifier is an internal name for an IPC object

## Identifiers and Keys

Whenever an IPC structure is being created, a key must be specified

```
#include <sys/ipc.h>
key_t ftok(const char *path, int id);
// Returns: key if OK, (key_t)1 on error
```

- ○ *path* must refer to an existing file
- ○ only the lower 8 bits of *id* are used when generating the key
- ○ key created by ftok is foremd by taking parts of st_dev and st_ino fields in the stat structure

## Identifiers and Keys

The three get functions (msgget, semget, and shmget) all have two
similar arguments: a *key* and an integer *flag*

If we want to create a new unique IPC structure, we must specify a *flag*
with both the IPC_CREAT and IPC_EXCL bits set

# Permission Structure

XSI IPC associates an `ipc_perm` structure with each IPC structure

All the fields are initialized when the IPC structure is created

`{msg | sem | shm}ctl` can modify the uid, gid, and mode fields

```
struct ipc_perm {
    uid_t uid; /* owner's effective user ID */
    gid_t gid; /* owner's effective group ID */
    uid_t cuid; /* creator's effective user ID */
    gid_t cgid; /* creator's effective group ID */
    mode_t mode; /* access modes */
    .
    .
    .
};
```

| Permission | Bit |
|---|---|
| user-read | 0400 |
| user-write (alter) | 0200 |
| group-read | 0040 |
| group-write (alter) | 0020 |
| other-read | 0004 |
| other-write (alter) | 0002 |

| File | $\longleftrightarrow$ | IPC |
|------|:---:|------|
| File name | $\longleftrightarrow$ | Key |
| file | $\longleftrightarrow$ | Object |
| fd | $\longleftrightarrow$ | Identifier |

Difference between object identifiers and fds

○ IPC identifiers are kernel-persistent; fds are process persistent

○ IPC identifiers are globally visible; fds are restricted to related processes

## Advantages and Disadvantages

Fundamental problem with XSI IPC:

1. IPC structures are systemwide and do not have a reference count
   - After creating a message queue, place some messages on the queue, and the terminate, the message queue and its contents ar not deleted
   - They remain until specifically read or deleted by `msgrcv` or `msgctl`, by executing the `ipcrm(1)` command, or by system being rebooted
2. IPC structures are not kwown by names in the file system
   - new system calls required (`msgget`, `semop`, `shmat`, and so on)
   - can't see the IPC objects with `ls` command or remove them with `rm`; instead two new commands—`ipcs(1)` and `ipcrm(1)`—were added

| IPC type | Connectionless? | Reliable? | Flow control? | Records? | Message types or priorities? |
|---|---|---|---|---|---|
| message queues | no | yes | yes | yes | yes |
| STREAMS | no | yes | yes | yes | yes |
| UNIX domain stream socket | no | yes | yes | no | no |
| UNIX domain datagram socket | yes | yes | no | yes | no |
| FIFOs (non-STREAMS) | no | yes | yes | no | no |

**Figure:** Comparison of features of various forms of IPC

By "connectionless," we mean the ability to send a message without having to call some form of an open function first

"Flow control" means that the sender is put to sleep if there is a shortage of system resources (buffers) or if the receiver can't accept any more messages.

XSI IPC : XSI MESSAGE QUEUES

## Message Queues

A message queue is a linked list of messages stored within the kernel and identified by a message queue identifier.

○ A new queue is created or an existing queue opened by `msgget`.

○ New messages are added to the end of a queue by `msgsnd`.

○ Messages are fetched from a queue by `msgrcv`.

○ We don't have to fetch the messages in a first-in, first-out order. Instead, we can fetch messages based on their type field.

## Message Qeues Datastructure

```
struct msqid_ds {
  struct ipc_perm msg_perm; /* see Section 15.6.2 */
  msgqnum_t msg_qnum; /* # of messages on queue */
  msglen_t msg_qbytes; /* max # of bytes on queue */
  pid_t msg_lspid; /* pid of last msgsnd() */
  pid_t msg_lrpid; /* pid of last msgrcv() */
  time_t msg_stime; /* last-msgsnd() time */
  time_t msg_rtime; /* last-msgrcv() time */
  time_t msg_ctime; /* last-change time */
. };
```

## Message Queues: msgget

msgget to either open an existing queue or create a new queue.

```
#include <sys/msg.h>
int msgget(key_t key, int flag);
// Returns: message queue ID if OK, 1 on error
```

- ○ The ipc_perm structure is initialized as described in Section 15.6.2. The mode member of this structure is set to the corresponding permission bits of flag.
- ○ msg_qnum, msg_lspid, msg_lrpid, msg_stime, and msg_rtime are all set to 0.
- ○ msg_ctime is set to the current time.
- ○ msg_qbytes is set to the system limit.

On success, msgget returns the non-negative queue ID. This value is then used with the other three message queue functions.

The `msgctl` function performs various operations on a queue

```
#include <sys/msg.h>
int msgctl(int msqid, int cmd, struct msqid_ds *buf );
// Returns: 0 if OK, 1 on error
```

cmd specifies the command to be performed on the queue specified by *msqid*

- ○ IPC_STAT: Fetch the msqid_ds structure for this queue, storing it in the structure pointed to by buf.
- ○ IPC_SET: Copy the msg_perm.uid, msg_perm.gid, msg_perm.mode, and msg_qbytes fields from the structure pointed to by buf to the msqid_ds structure associated with this queue.
- ○ IPC_RMID: Remove the message queue from the system and any data still on the queue. This removal is immediate. Any other process still using the message queue will get an error of EIDRM on its next attempted operation on the queue.

## Message Queues: msgsnd

Data is placed onto a message queue by calling msgsnd.

```
#include <sys/msg.h>
int msgsnd(int msqid, const void *ptr, size_t nbytes, int flag);
// Returns: 0 if OK, 1 on error
```

Messages are always placed at the end of the queue.

*ptr* argument is a pointer to a mymesg structure.

○ The message type can be used by the receiver to fetch messages in an order other than first in, first out.

```
struct mymesg {
    long mtype; /* positive message type */
    char mtext[512]; /* message data, of length nbytes */
};
```

Data is placed onto a message queue by calling msgsnd.

```
#include <sys/msg.h>
int msgsnd(int msqid, const void *ptr, size_t nbytes, int flag);
// Returns: 0 if OK, 1 on error
```

*flag* value of IPC_NOWAIT can be specified (similar to nonblocking I/O flag for file I/O)

○ If IPC_NOWAIT is set and message queue is full: msgsnd to return immediately with an error of EAGAIN

○ If IPC_NOWAIT is not set, we are blocked until there is room for the message, the queue is removed from the system, or a signal is caught and the signal handler returns. EINTR is returned

Since reference count is not maintained, removal of a queue simply generates error on the next queue operation

Data is placed onto a message queue by calling msgsnd.

```
#include <sys/msg.h>
int msgsnd(int msqid, const void *ptr, size_t nbytes, int flag);
// Returns: 0 if OK, 1 on error
```

When msgsnd returns successfully, the msqid_ds structure associated
with the message queue is updated

○ to indicate the process ID that made the call (msg_lspid),

○ the time that the call was made (msg_stime),

○ and that one more message is on the queue (msg_qnum).

Messages are retrieved from a queue by msgrcv.

```
#include <sys/msg.h>
ssize_t msgrcv(int msqid, void *ptr, size_t nbytes, long type, int flag);
// Returns: size of data portion of message if OK, 1 on error
```

As with msgsnd, the *ptr* argument points to a long integer (where the message type of the returned message is stored) followed by a data buffer for the actual message data.

○ if returned message is larger than *nbytes*, MSG_NOERROR is set: the message is truncated

○ if returned message is larger than *nbytes*, MSG_NOERROR is not set: error E2BIG is returned

Messages are retrieved from a queue by msgrcv.

```
#include <sys/msg.h>
ssize_t msgrcv(int msqid, void *ptr, size_t nbytes, long type, int flag);
// Returns: size of data portion of message if OK, 1 on error
```

*type* lets us specify which message we want

- ○ type == 0 : The first message on the queue is returned.
- ○ type > 0 : The first message on the queue whose message type equals *type* is returned.
- ○ type < 0 : The first message on the queue whose message type is the lowest value less than or equal to the absolute value of *type* is returned.

A nonzero type is used to read the messages in an order other than first in, first out.

Messages are retrieved from a queue by msgrcv.

```
#include <sys/msg.h>
ssize_t msgrcv(int msqid, void *ptr, size_t nbytes, long type, int flag);
// Returns: size of data portion of message if OK, 1 on error
```

We can specify a flag value of IPC_NOWAIT to make the operation nonblocking,

When msgrcv succeeds, the kernel updates the msqid_ds structure associated with the message queue to indicate

- the caller's process ID (msg_lrpid),
- the time of the call (msg_rtime),

codes/msg_send.c

```
1   #include <stdio.h>
2   #include <string.h>
3   #include <stdlib.h>
4   #include <sys/ipc.h>
5   #include <sys/msg.h>
6   #include <sys/types.h>
7
8   #define MSGSZ 128
9
10  /* the message structure. */
11  typedef struct mymesg {
12          long mtype;
13          char mtext[MSGSZ];
14          } message_buf;
15
16  int main()
17  {
18      int msqid;
19      int msgflg = IPC_CREAT | 0666;
20      key_t key;
21      message_buf sbuf;
```

```
22      size_t buf_length;
23
24      if((key = ftok("Makefile", 'R')) == -1){
25          perror("ftok");
26          exit(1);
27      }
28
29      fprintf(stderr, "\nmsgget: Calling msgget(%#x, %#o)\n", key, msgflg);
30
31      if ((msqid = msgget(key, msgflg )) < 0) {
32          perror("msgget");
33          exit(1);
34      }
35      fprintf(stderr,"msgget: msgget succeeded: msqid = %d\n", msqid);
36
37      /* We'll send message type 1 */
38      sbuf.mtype = 1;
39      printf("Enter a message to add to the message queue: ");
40      scanf("%[^\n]", sbuf.mtext);
41      getchar();
42
43      buf_length = strlen(sbuf.mtext) + 1 ;
44
45      /* Send a message.*/
46      if (msgsnd(msqid, &sbuf, buf_length, IPC_NOWAIT) < 0) {
47          printf ("%d, %ld, %s, %ld\n", msqid, sbuf.mtype, sbuf.mtext, buf_length);
48          perror("msgsnd");
```

```
49          exit(1);
50      } else
51        printf("Message: \"%s\" Sent\n", sbuf.mtext);
52
53      exit(0);
54  }
```

codes/msg_receive.c

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <sys/ipc.h>
4   #include <sys/msg.h>
5   #include <sys/types.h>
6
7   #define MSGSZ 128
8
9   /* Declare the message structure */
10
11  typedef struct msgbuf {
12      long mtype;
13      char mtext[MSGSZ];
14  } message_buf;
15
16
17  int main()
18  {
19      int msqid;
20      key_t key;
21      message_buf rbuf;
22
23      if((key = ftok("Makefile", 'R')) == -1){
24          perror("ftok");
```

```
25          exit(1);
26      }
27
28      if ((msqid = msgget(key, 0666)) < 0) {
29          perror("msgget");
30          exit(1);
31      }
32
33      /* Receive an answer of message type 1. */
34      if (msgrcv(msqid, &rbuf, MSGSZ, 1, 0) < 0) {
35          perror("msgrcv");
36          exit(1);
37      }
38
39      /* Print the answer. */
40      printf("%s\n", rbuf.mtext);
41      exit(0);
42  }
```

# XSI IPC : XSI Semaphore

## Semaphores

Semaphore is a counter a counter used to provide access to a shared data object for multiple processes.

To obtain a shared resource, a process needs to do the following:

1. Test the semaphore that controls the resource.
2. If the value of the semaphore **is positive**, the process can use the resource. In this case, the process **decrements the semaphore value by 1**, indicating that it has used one unit of the resource.
3. Otherwise, if the value of the **semaphore is 0**, the process goes to **sleep** until the semaphore value **is greater than 0**. When the process wakes up, it returns to step 1.

A common form of semaphore is called a *binary semaphore*

○ controls a signle resource, and it is initialized to 1

## XSI Semahores

XSI semaphores are more complicated than this. Three features contribute to this unnecessary complication.

1. A semaphore is defined as a set of one or more semaphore values. When we create a semaphore, we specify the number of values in the set.

2. The creation of a semaphore (semget) is independent of its initialization (semctl). We cannot atomically create a new semaphore set and initialize all the values in the set.

3. Since all forms of XSI IPC remain in existence even when no process is using them, we have to worry about a program that terminates without releasing the semaphores it has been allocated.

## XSI Sempahores Data Structures

The kernel maintains a semid_ds structure for each semaphore set:

```
struct semid_ds {
    struct ipc_perm sem_perm; /* see Section 15.6.2 */
    unsigned short sem_nsems; /* # of semaphores in set */
    time_t sem_otime; /* last-semop() time */
    time_t sem_ctime; /* last-change time */
    .
    .
    .
};
```

## XSI Sempahores Data Structures

Each semaphore is represented by an anonymous structure containing at least the following members:

```
struct {
    unsigned short semval; /* semaphore value, always >= 0 */
    pid_t sempid; /* pid for last operation */
    unsigned short semncnt; /* # processes awaiting semval>curval */
    unsigned short semzcnt; /* # processes awaiting semval==0 */
    .
    .
    .
};
```

When we want to use XSI semaphores, we first need to obtain a semaphore ID by calling the semget function.

```
#include <sys/sem.h>
int semget(key_t key, int nsems, int flag);
// Returns: semaphore ID if OK, 1 on error
```

- ○ The ipc_perm structure is initialized as described in Section 15.6.2. The mode member of this structure is set to the corresponding permission bits of flag.
- ○ sem_otime is set to 0.
- ○ sem_ctime is set to the current time.
- ○ sem_nsems is set to *nsems*.
  - The number of semaphores in the set is nsems.
  - If a new set is being created (typically by a server), we must specify nsems.
  - If we are referencing an existing set (a client), we can specify nsems as 0.

The semctl function is the catchall for various semaphore operations.

```
#include <sys/sem.h>
int semctl(int semid, int semnum, int cmd, ... /* union semun arg */ );
// Returns: (see following)
```

The fourth argument is optional, depending on the command requested, and if present, is of type semun, a union of various command-specific arguments:

```
union semun {
    int val; /* for SETVAL */
    struct semid_ds *buf; /* for IPC_STAT and IPC_SET */
    unsigned short *array; /* for GETALL and SETALL */ };
```

# XSI Semaphores: `semctl`

The *cmd* argument specifies one of the following ten commands to be performed on the set specified by semid.

- ○ `IPC_STAT`: Fetch the `semid_ds` structure for this set, storing it in the structure pointed to by *arg.buf*.
- ○ `IPC_SET`: Set the `sem_perm.uid`, `sem_perm.gid`, and `sem_perm.mode` fields from the structure pointed to by `arg.buf` in the `semid_ds` structure associated with this set.
- ○ `IPC_RMID`: Remove the semaphore set from the system. This removal is immediate. Any other process still using the semaphore will get an error of `EIDRM` on its next attempted operation on the semaphore.
- ○ `GETVAL`: Return the value of `semval` for the member *semnum*.
- ○ `SETVAL`: Set the value of `semval` for the member *semnum*. The value is specified by *arg.val*.
- ○ `GETPID`: Return the value of `sempid` for the member *semnum*.
- ○ `GETNCNT`: Return the value of `semncnt` for the member *semnum*.
- ○ `GETZCNT`: Return the value of `semzcnt` for the member *semnum*.
- ○ `GETALL`: Fetch all the semaphore values in the set. These values are stored in the array pointed to by *arg.array*.
- ○ `SETALL`: Set all the semaphore values in the set to the values pointed to by *arg.array*.

# XSI IPC : XSI Shared Memory

## XSI Shared Memory

Shared memory allows two or more processes to share a given region of memory. This is the fastest form of IPC, because the data does not need to be copied between the client and the server.

The kernel maintains a structure with at least the following members for each shared memory segment:

```
struct shmid_ds {
  struct ipc_perm shm_perm; /* see Section 15.6.2 */
  size_t shm_segsz; /* size of segment in bytes */
  pid_t shm_lpid; /* pid of last shmop() */
  pid_t shm_cpid; /* pid of creator */
  shmatt_t shm_nattach; /* number of current attaches */
  time_t shm_atime; /* last-attach time */
  time_t shm_dtime; /* last-detach time */
  time_t shm_ctime; /* last-change time */
  .
};
```

## XSI Shared Memory: shmget

The first function called is usually shmget, to obtain a shared memory identifier.

```
#include <sys/shm.h>
int shmget(key_t key, size_t size, int flag);
// Returns: shared memory ID if OK, 1 on error
```

- ○ The ipc_perm structure is initialized as described in Section 15.6.2. The mode member of this structure is set to the corresponding permission bits of flag.
- ○ shm_lpid, shm_nattch, shm_atime, and shm_dtime are all set to 0.
- ○ shm_ctime is set to the current time.
- ○ shm_segsz is set to the *size* requested
- ○ The *size* parameter is the size of the shared memory segment in *bytes*

The shmctl function is the catchall for various operations.

```
#include <sys/shm.h>
int shmctl(int shmid, int cmd, struct shmid_ds *buf );
// Returns: 0 if OK, 1 on error
```

The *cmd* argument specifies one of the following five commands to be performed, on the segment specified by *shmid*.

○ IPC_STAT: Fetch the shmid_ds structure for this segment, storing it in the structure pointed to by *buf*.

○ IPC_SET: Set the following three fields from the structure pointed to by *buf* in the shmid_ds structure associated with this shared memory segment: shm_perm.uid, shm_perm.gid, and shm_perm.mode.

○ IPC_RMID: Remove the shared memory segment set from the system. Since an attachment count is maintained for shared memory segments (the shm_nattch field in the shmid_ds structure), the segment is not removed until the last process using the segment terminates or detaches it. Regardless of whether the segment is still in use, the segment's identifier is immediately removed so that shmat can no longer attach the segment.

## XSI Shared Memory: attaching

Once a shared memory segment has been created, a process attaches it to its address space by calling shmat.

```
#include <sys/shm.h>
void *shmat(int shmid, const void *addr, int flag);
// Returns: pointer to shared memory segment if OK, 1 on error
```

The address in the calling process at which the segment is attached depends on the *addr* argument and whether the SHM_RND bit is specified in *flag*.

○ If *addr* is 0, the segment is attached at the first available address selected by the kernel. **This is the recommended technique.**

○ If *addr* is nonzero and SHM_RND is not specified, the segment is attached at the address given by *addr*.

○ If *addr* is nonzero and SHM_RND is specified, the segment is attached at the address given by (*addr*-(*addr* modulus SHMLBA)).

## XSI Shared Memory: detaching

When we're done with a shared memory segment, we call shmdt to detach it.

```
#include <sys/shm.h>
int shmdt(const void *addr);
// Returns: 0 if OK, 1 on error
```

The *addr* argument is the value that was returned by a previous call to shmat. If successful, shmdt will decrement the shm_nattch counter in the associated shmid_ds structure.

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <sys/shm.h>
4
5   #define ARRAY_SIZE 40000
6   #define MALLOC_SIZE 100000
7   #define SHM_SIZE 100000
8   #define SHM_MODE 0600 /* user read/write */
9
10  char array[ARRAY_SIZE]; /* uninitialized data = bss */
11
12  int
13  main(void)
14  {
15     int  shmid;
16     char *ptr, *shmptr;
17
18     printf("array[] from %p to %p\n", (void *)&array[0],
19       (void *)&array[ARRAY_SIZE]);
20     printf("stack around %p\n", (void *)&shmid);
21
22     if ((ptr = malloc(MALLOC_SIZE)) == NULL){
23        perror("malloc error");
24        exit(1);
25     }
26     printf("malloced from %p to %p\n", (void *)ptr,
```

```
27      (void *)ptr+MALLOC_SIZE);
28
29      if ((shmid = shmget(IPC_PRIVATE, SHM_SIZE, SHM_MODE)) < 0){
30        perror("shmget error");
31        exit(1);
32      }
33      if ((shmptr = shmat(shmid, 0, 0)) == (void *)-1){
34        perror("shmat error");
35        exit(1);
36      }
37      printf("shared memory attached from %p to %p\n", (void *)shmptr,
38        (void *)shmptr+SHM_SIZE);
39
40      if (shmctl(shmid, IPC_RMID, 0) < 0){
41        perror("shmctl error");
42        exit(1);
43      }
44
45      exit(0);
46    }
```

## Shared Memory Example: codes/tshm.c III

```
$ ./tshm
array[] from 0x1085d3050 to 0x1085dcc90
stack around 0x7fff5762cf88
malloced from 0x7f8212802000 to 0x7f821281a6a0
shared memory attached from 0x108611000 to 0x1086296a0
```
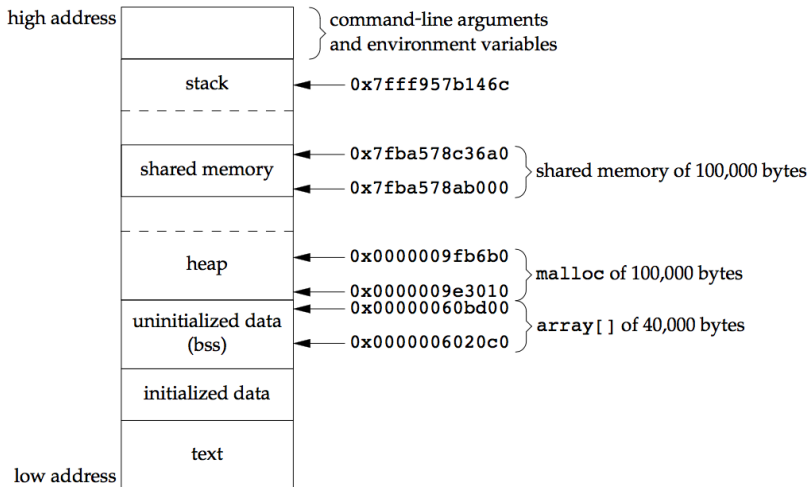
**Figure:** Memory layout on an Intel-based Linux System

## Shared Memory Example: Server - Client

codes/shm_send.c

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <string.h>
4   #include <unistd.h>
5   #include <sys/types.h>
6   #include <sys/ipc.h>
7   #include <sys/shm.h>
8
9   #define SHSIZE 100
10
11  int main(int argc, char *argv[])
12  {
13      int shmid;
14      key_t key;
15      char *shm, *s;
16      char *message = "The server sends a line of text0";
17
18      // generate key
19      if ( (key = ftok("Makefile", 'R')) == -1 )
20      {
21          perror("ftok");
22          exit(1);
23      }
24
25      // connect to the segment
26      if ( (shmid = shmget( key, SHSIZE, IPC_CREAT | 0666 )) < 0 )
27      {
28          perror( "shmget" );
29          exit(1);
```

# LAST WORDS

## Last Words

○ prepare for the exam