

# Advanced Programming in the UNIX Environment — *Advanced I/O*

Hop Lee  
hoplee@bupt.edu.cn

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Nonblocking I/O</b>	<b>1</b>
<b>3</b>	<b>Record Locking</b>	<b>2</b>
<b>4</b>	<b>STREAMS</b>	<b>2</b>
4.1	STREAMS Messages . . . . .	3
4.2	putmsg and putpmsg Functions . . . . .	4
4.3	STREAMS ioctl Operations . . . . .	4
4.4	getmsg and getpmsg Functions . . . . .	4
<b>5</b>	<b>I/O Multiplexing</b>	<b>5</b>
5.1	select and pselect Functions . . . . .	6
<b>6</b>	<b>readv and writev Functions</b>	<b>7</b>
<b>7</b>	<b>Memory-Mapped I/O</b>	<b>8</b>

## 1 Introduction

- This chapter covers numerous topics and functions that we lump under the term **advanced I/O**:
  - nonblocking I/O,
  - record locking,
  - System V STREAMS,
  - I/O multiplexing (the **select** and **poll** functions),
  - the **readv** and **writev** functions,
  - and memory-mapped I/O (**mmap**).

## 2 Nonblocking I/O

- In previous chapter, we said that the system calls are divided into two categories: the “slow” ones and all the others. The slow system calls are those that can block forever. They include
  - Reads that can block the caller forever if data isn’t present with certain file types (pipes, terminal devices, and network devices)

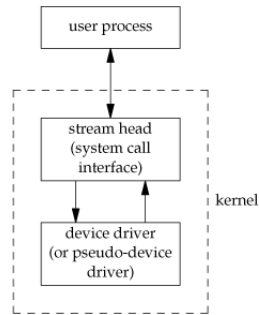


Figure 1: A simple stream

- Writes that can block the caller forever if the data can't be accepted immediately by these same file types (no room in the pipe, network flow control, etc.)
- Opens that block until some condition occurs on certain file types
- Reads and writes of files that have mandatory record locking enabled
- Certain `ioctl` operations
- Some of the IPC functions
- We also said that system calls related to disk I/O are not considered slow, even though the read or write of a disk file can block the caller temporarily.
- Nonblocking I/O lets us issue an I/O operation and not have it block forever.
- There are two ways to specify nonblocking I/O for a given descriptor:
  1. Call `open` with the `O_NONBLOCK` flag (Section 3.3).
  2. For a descriptor that is already open, we call `fcntl` to turn on the `O_NONBLOCK` file status flag (Section 3.14).
- Example 12.1

### 3 Record Locking

- Omitted.

## 4 STREAMS

- The **STREAMS** mechanism is provided by System V as a general way to interface communication drivers into the kernel.
- We need to discuss STREAMS to understand the terminal interface in System V, the use of the `poll` function for I/O multiplexing, and the implementation of STREAMS-based pipes and named pipes.
- A stream provides a full-duplex path between a user process and a device driver. There is no need for a stream to talk to a hardware device.
- Beneath the stream head, we can push processing modules onto the stream. This is done using an `ioctl` command.
- Any number of processing modules can be *pushed* onto a stream.

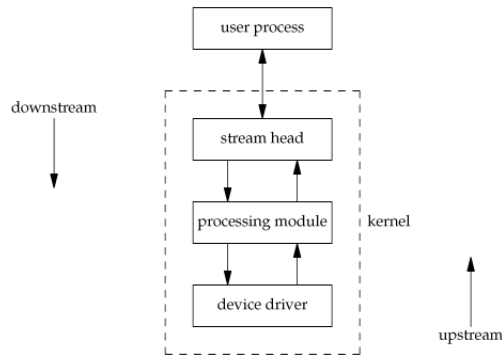


Figure 2: A simple stream

- STREAMS modules are similar to device drivers in that they execute as part of the kernel, and they are normally link edited into the kernel when the kernel is built.
- We access a stream with the functions from Chapter 3: `open`, `close`, `read`, `write`, and `ioctl`. Additionally, three new functions were added to the SVR3 kernel to support STREAMS (`getmsg`, `putmsg`, and `poll`), and another two (`getpmsg` and `putpmsg`) were added with SVR4 to handle messages with different priority bands within a stream.

## 4.1 STREAMS Messages

- All input and output under STREAMS is based on messages.
- The stream head and the user process exchange messages using `read`, `write`, `ioctl`, `getmsg`, `getpmsg`, `putmsg`, and `putpmsg`.
- Messages are also passed up and down a stream between the stream head, the processing modules, and the device driver.
- Between the user process and the stream head, a message consists of a message type, optional control information, and optional data.
- The control information and data are specified by `strbuf` structures:

```

1 struct strbuf{
2     int maxlen; /* size of buffer */
3     int len;    /* number of bytes currently in buffer */
4     char *buf;  /* pointer to buffer */
5 };

```

- There are about 25 different types of messages, but only a few of these are used between the user process and the stream head. The rest are passed up and down a stream within the kernel.
- We'll encounter only three of these message types with the functions we use:
  1. `M_DATA` (user data for I/O)
  2. `M_PROTO` (protocol control information)
  3. `M_PCPROTO` (high-priority protocol control information)
- Every message on a stream has a queuing priority:

- High-priority messages (highest priority)
  - Priority band messages
  - Ordinary messages (lowest priority)
- Each STREAMS module has two input queues. One receives messages from the module above, and one receives messages from the module below. The messages on an input queue are arranged by priority.

## 4.2 putmsg and putpmsg Functions

- A STREAMS message (control information or data, or both) is written to a stream using either `putmsg` or `putpmsg`. The difference between these two functions is that the latter allows us to specify a priority band for the message.

```

1 #include <stropts.h>
2 int putmsg(int fildes, const struct strbuf *ctlptr,
3           const struct strbuf *dataptr, int flag);
4 int putpmsg(int fildes, const struct strbuf *ctlptr,
5            const struct strbuf *dataptr, int band,
6            int flag);

```

- We can also `write` to a stream, which is equivalent to a `putmsg` without any control information and with a *flag* of 0.
- These two functions can generate the three different priorities of messages: ordinary, priority band, and high priority.

## 4.3 STREAMS ioctl Operations

- Between Linux and Solaris, there are almost 40 different operations that can be performed on a stream using `ioctl`. Most of these operations are documented in the `streamio(7)` manual page.
- Example: `isastream` function.
- Example 12.10: List the names of the modules on a stream

## 4.4 getmsg and getpmsg Functions

- STREAMS messages are read from a stream head using `read`, `getmsg`, or `getpmsg`.

```

1 #include <stropts.h>
2 int getmsg(int fildes, struct strbuf *restrict ctlptr,
3           struct strbuf *restrict dataptr,
4           int *restrict flagptr);
5 int getpmsg(int fildes, struct strbuf *restrict ctlptr,
6            struct strbuf *restrict dataptr,
7            int *restrict bandptr,
8            int *restrict flagptr);

```

- Example 14.19. Copy standard input to standard output using `getmsg`.

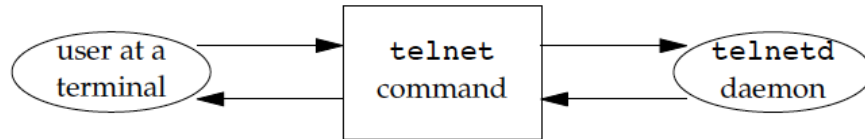


Figure 3: Overview of telnet program

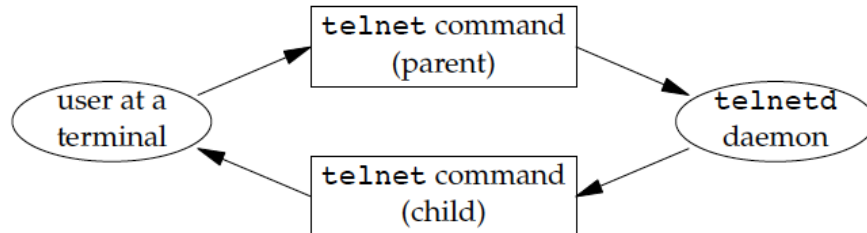


Figure 4: The telnet program using two processes

## 5 I/O Multiplexing

- When we read from one descriptor and write to another, we can use blocking I/O in a loop.
- What if we have to read from two descriptors? In this case, we can't do a blocking **read** on either descriptor, as data may appear on one descriptor while we're blocked in a **read** on the other. A different technique is required to handle this case.
- Let's look at the structure of the **telnet(1)** command.
- The **telnet** process has two inputs and two outputs. We can't do a blocking **read** on either of the inputs, as we never know which input will have data for us.
- One way to handle this particular problem is to divide the process in two pieces (using **fork**), with each half handling one direction of data.
- But this leads to a problem when the operation terminates: it does complicate the program somewhat.
- Instead of two processes, we could use two threads in a single process. This avoids the termination complexity, but requires that we deal with synchronization between the threads, which could add more complexity than it saves.
- We could use nonblocking I/O in a single process by setting both descriptors nonblocking and adopting poll mechanism. The problem is that it wastes CPU time.
- Although it works on any system that supports nonblocking I/O, polling should be avoided on a multitasking system.
- Another technique is called **asynchronous I/O**. To do this, we tell the kernel to notify us with a signal when a descriptor is ready for I/O. There are two problems with this. First, not all systems support this feature. Second, there is only one of these signals per process.
- A better technique is to use **I/O multiplexing**. To do this, we build a list of the descriptors that we are interested in and call a function that doesn't return until one of the descriptors is ready for I/O. On return from the function, we are told which descriptors are ready for I/O.
- Three functions **poll**, **pselect**, and **select** allow us to perform I/O multiplexing.

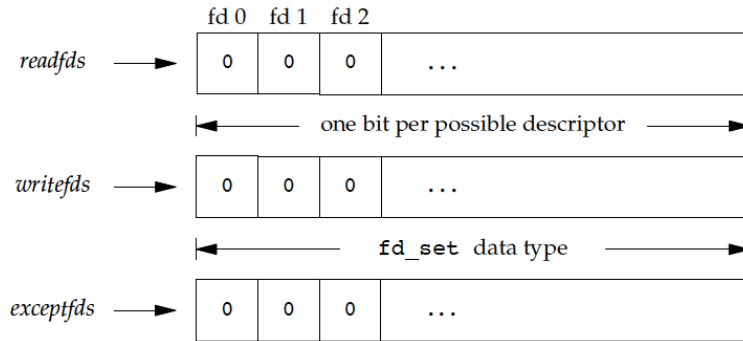


Figure 5: Specifying the read, write, and exception descriptors for `select`

## 5.1 `select` and `pselect` Functions

- The `select` function lets us do I/O multiplexing under all POSIX-compatible platforms. The arguments we pass to `select` tell the kernel
  - Which descriptors we’re interested in.
  - What conditions we’re interested in for each descriptor.
  - How long we want to wait.
- On the return from `select`, the kernel tells us
  - The total count of the number of descriptors that are ready
  - Which descriptors are ready for each of the three conditions
- With this return information, we can call the appropriate I/O function and know that the function won’t block.

```

1 #include <sys/select.h>
2 int select(int maxfdp1, fd_set *restrict readfds,
3           fd_set *restrict writefds,
4           fd_set *restrict exceptfds,
5           struct timeval *restrict tvptr);
6 struct timeval{
7     long tv_sec;    /* seconds */
8     long tv_usec;   /* and microseconds */
9 };

```

- First parameter `maxfdp1` is *maximum fd plus 1*, indicates the number of fd we want check.
- The middle three arguments `readfds`, `writefds`, and `exceptfds` are pointers to **descriptor sets**.
- The only thing we can do with the `fd_set` data type is allocate a variable of this type, assign a variable of this type to another variable of the same type, or use one of the following four functions on a variable of this type.

```

1 #include <sys/select.h>
2 int FD_ISSET(int fd, fd_set *fdset);
3 void FD_CLR(int fd, fd_set *fdset);
4 void FD_SET(int fd, fd_set *fdset);
5 void FD_ZERO(fd_set *fdset);

```

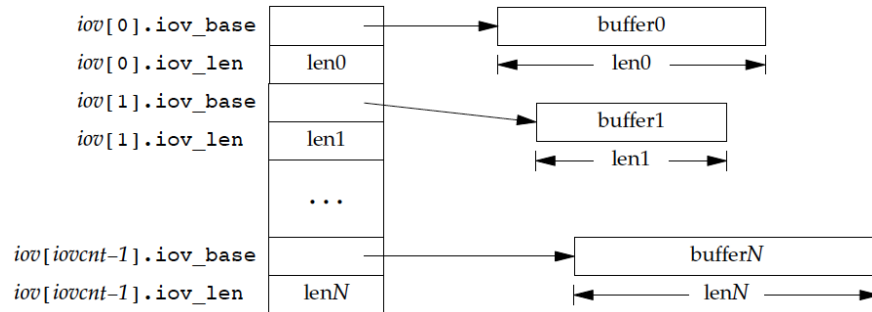


Figure 6: The `iovec` structure for `readv` and `writev`

- These interfaces can be implemented as either macros or functions.
- There are three possible return values from `select`.
  1. A return value of `-1` means that an error occurred. In this case, none of the descriptor sets will be modified.
  2. A return value of `0` means that no descriptors are ready. This happens if the time limit expires before any of the descriptors are ready. When this happens, all the descriptor sets will be zeroed out.
  3. A positive return value specifies the number of descriptors that are ready. This value is the sum of the descriptors ready in all three sets, so if the same descriptor is ready to be read and written, it will be counted twice in the return value. The only bits left on in the three descriptor sets are the bits corresponding to the descriptors that are ready.

## 6 `readv` and `writev` Functions

- The `readv` and `writev` functions let us read into and write from multiple noncontiguous buffers in a single function call. These operations are called **scatter read** and **gather write**.

```

1 #include <sys/uio.h>
2 ssize_t readv(int filedes, const struct iovec *iov,
3               int iovcnt);
4 ssize_t writev(int filedes, const struct iovec *iov,
5                int iovcnt);

```

- The second argument to both functions is a pointer to an array of `iovec` structures:

```

1 struct iovec{
2     void *iov_base; /* starting address of buffer */
3     size_t iov_len; /* size of buffer */
4 };

```

- The number of elements in the `iov` array is specified by `iovcnt`. It is limited to `IOV_MAX`.
- The `writev` function gathers the output data from the buffers in order: `iov[0]`, `iov[1]`, through `iov[iovcnt-1]`; `writev` returns the total number of bytes output, which should normally equal the sum of all the buffer lengths.
- The `readv` function scatters the data into the buffers in order, always filling one buffer before proceeding to the next. `readv` returns the total number of bytes that were read. A count of `0` is returned if there is no more data and the end of file is encountered.

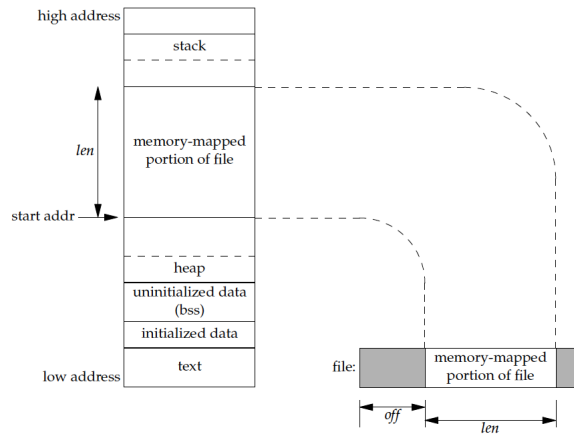


Figure 7: Example of a memory-mapped file

- Example. As we expect, the system time increases when we call `write` twice, compared to calling either `write` or `writew` once.

## 7 Memory-Mapped I/O

- **Memory-mapped I/O** lets us map a file on disk into a buffer in memory so that, when we fetch bytes from the buffer, the corresponding bytes of the file are read. Similarly, when we store data in the buffer, the corresponding bytes are automatically written to the file. This lets us perform I/O without using `read` or `write`.
- To use this feature, we have to tell the kernel to map a given file to a region in memory. This is done by the `mmap` function.

```
1 #include <sys/mman.h>
2 void *mmap(void *addr, size_t len, int prot,
3           int flag, int filed, off_t off);
```

- The return value of this function is the starting address of the mapped area.
- The `addr` argument lets us specify the address of where we want the mapped region to start. We normally set this to 0 to allow the system to choose the starting address.
- The `filedes` argument is the file descriptor specifying the file that is to be mapped. We have to open this file before we can map it into the address space.
- The `len` argument is the number of bytes to map, and `off` is the starting offset in the file of the bytes to map.
- The `prot` argument specifies the protection of the mapped region.
  - We can specify the protection as either `PROT_NONE` or the bitwise OR of any combination of `PROT_READ`, `PROT_WRITE`, and `PROT_EXEC`.
  - The protection specified for a region can't allow more access than the open mode of the file.
- The `flag` argument affects various attributes of the mapped region.



**MAP\_FIXED** The return value must equal *addr*. Use of this flag is discouraged, as it hinders portability. If this flag is not specified and if *addr* is nonzero, then the kernel uses *addr* as a hint of where to place the mapped region, but there is no guarantee that the requested address will be used. Maximum portability is obtained by specifying *addr* as 0.

**MAP\_SHARED** This flag describes the disposition of store operations into the mapped region by this process. This flag specifies that store operations modify the mapped file that is, a store operation is equivalent to a *write* to the file. Either this flag or the next (**MAP\_PRIVATE**), but not both, must be specified.

**MAP\_PRIVATE** This flag says that store operations into the mapped region cause a private copy of the mapped file to be created. All successive references to the mapped region then reference the copy.

- The value of *off* and the value of *addr* (if **MAP\_FIXED** is specified) are required to be multiples of the system's virtual memory page size.
- We cannot append to a file with *mmap*. We must first grow the file.
- A memory-mapped region is inherited by a child across a *fork*, but is not inherited by the new program across an *exec*.
- Example.
- We can change the permissions on an existing mapping by calling *mprotect*.

```
1 #include <sys/mman.h>
2 int mprotect(void *addr, size_t len, int prot);
```

- If the pages in a shared mapping have been modified, we can call *msync* to flush the changes to the file that backs the mapping. The *msync* function is similar to *fsync*, but works on memory-mapped regions.

```
1 #include <sys/mman.h>
2 int msync(void *addr, size_t len, int flags);
```

- A memory-mapped region is automatically unmapped when the process terminates or by calling *munmap* directly. Closing the file descriptor *filedes* does **NOT** unmap the region.

```
1 #include <sys/mman.h>
2 int munmap(caddr_t addr, size_t len);
```

- The call to *munmap* does not cause the contents of the mapped region to be written to the disk file.
- Memory-mapped I/O is faster when copying one regular file to another. There are limitations. We can't use it to copy between certain devices (such as a network device or a terminal device), and we have to be careful if the size of the underlying file could change after we map it.

## The End of Chapter 14.