

# SYSTEM PROGRAMMING

## WEEK 13: ADVANCED I/O

---

Seongjin Lee

Updated: 2018/07/06

10\_adv-io

[insight@gnu.ac.kr](mailto:insight@gnu.ac.kr)

<http://open.gnu.ac.kr>

Systems Research Lab.

Gyeongsang National University



# Table of contents

1. Nonblocking I/O
2. Record Locking
3. Asynchronous I/O
4. `readv` and `writv` Functions
5. `readn` and `writen` Functions
6. Memory-Mapped I/O



# Introduction

This chapter covers

- Nonblocking I/O
- Record Locking
- Asynchronous I/O
- memory-mapped I/O (mmap)



# NONBLOCKING I/O

---

# Nonblocking I/O

## Disk I/O are not considered slow

Nonblocking I/O lets us issue an I/O operation, such as an open, read, or write, and not have it block forever.

- if the operation cannot be completed, the call returns immediately with an error noting that the operation would have blocked

Two ways to specify nonblocking I/O for a given descriptor

- open with `O_NONBLOCK` flag
- For already opened descriptors, use `fcntl` to turn on the `O_NONBLOCK` file status flag



# Nonblocking I/O Example: codes/nonblocking.c |

```
1  #include <stdio.h>
2  #include <unistd.h>
3  #include <stdlib.h>
4  #include <errno.h>
5  #include <fcntl.h>
6
7  char buf[500000];
8
9
10
11 /* flags are the file status flags to turn off */
12 void
13 clr_fl(int fd, int flags)
14 {
15     int val;
16
17     if ((val = fcntl(fd, F_GETFL, 0)) < 0){
18         fprintf(stderr, "fcntl F_GETFL error");
19         exit(1);
20     }
21
22     val &= ~flags; /* turn flags off */
```



# Nonblocking I/O Example: codes/nonblocking.c II

```
23
24     if (fcntl(fd, F_SETFL, val) < 0){
25         fprintf(stderr, "fcntl F_SETFL error");
26         exit(1);
27     }
28 }
29
30 /* flags are file status flags to turn on */
31 void
32 set_fl(int fd, int flags)
33 {
34     int val;
35
36     if ((val = fcntl(fd, F_GETFL, 0)) < 0){
37         fprintf(stderr, "fcntl F_GETFL error");
38         exit(1);
39     }
40
41     val |= flags; /* turn on flags */
42
43     if (fcntl(fd, F_SETFL, val) < 0){
44         fprintf(stderr, "fcntl F_SETFL error");
45         exit(1);
```



# Nonblocking I/O Example: codes/nonblocking.c III

```
46     }
47 }
48
49 int
50 main(void)
51 {
52     int ntowrite, nwrite;
53     char *ptr;
54     ntowrite = read(STDIN_FILENO, buf, sizeof(buf));
55     fprintf(stderr, "read %d bytes\n", ntowrite);
56
57     set_fl(STDOUT_FILENO, O_NONBLOCK); /* set nonblocking */
58
59     ptr = buf;
60
61     while (ntowrite > 0) {
62         errno = 0;
63         nwrite = write(STDOUT_FILENO, ptr, ntowrite);
64         fprintf(stderr, "nwrite = %d, errno = %d\n", nwrite, errno);
65         if (nwrite > 0) {
66             ptr += nwrite;
67             ntowrite -= nwrite;
68         }
69     }
```





# Nonblocking I/O Example: codes/nonblocking.c IV

```
69     }  
70  
71     clr_fl(STDOUT_FILENO, O_NONBLOCK); /* clear nonblocking */  
72  
73     exit(0);  
74 }
```

## running the example

```
James@maker:codes$ ls -lh /etc/services  
-rw-r--r--  1 root  wheel   662K Jul 31 04:32 /etc/services  
James@maker:codes$ ./nonblocking < /etc/services > temp  
read 500000 bytes  
nwrite = 500000, errno = 0  
James@maker:codes$ ls -lh temp  
-rw-r--r--  1 James  staff   488K Nov 28 15:29 temp  
James@maker:codes$ ./nonblocking < /etc/services 2>stderr  
#  
# Network services, Internet style  
# ...
```



# Nonblocking I/O Example: codes/nonblocking.c V

understanding the result

errno = 35 is EAGAIN

```
James@maker:codes$ head stderr
```

```
read 500000 bytes
```

```
nwrite = 999, errno = 0
```

```
nwrite = -1, errno = 35
```

```
nwrite = -1, errno = 35
```

```
nwrite = -1, errno = 35
```

```
nwrite = -1, errno = 35
```

```
nwrite = -1, errno = 35
```

```
nwrite = -1, errno = 35
```

```
nwrite = 1001, errno = 0
```

```
nwrite = 1002, errno = 0
```



# Nonblocking I/O Example: codes/nonblocking.c VI

Finding the total sum of writes

```
grep "nwrite" stderr |  
  grep -v "\-1" |  
  awk -F "," '{print $1}' |  
  awk -F "=" '{sum=sum+$2};END{print sum}'
```

Finding the total number of errors

```
grep "nwrite" stderr |  
  grep "\-1" | wc -l
```



# RECORD LOCKING



# Record Locking

What happens when two people edit the same file at the same time?

- the final state of the file corresponds to the last process that wrote the file

In some applications, a process needs to be certain that it alone is writing to a file

- To provide this capability, commercial UNIX systems provide record locking



## Record Locking

It is the term normally used to describe the ability of a process to prevent other processes from modifying a region of a file while the first process is reading or modifying that portion of the file

**A better term is byte-range locking**



# Record Locking: History

Early Berkeley releases supported only the `flock` function

- This function locks only entire files, not regions of a file.

Record locking was added to System V Release 3 through the `fcntl` function.

- The `lockf` function was built on top of this, providing a simplified interface.
- These functions allowed callers to lock arbitrary byte ranges in a file, ranging from the entire file down to a single byte within the file.



# Record Locking: History cnt'd

POSIX.1 chose to standardize on the `fcntl` approach

System	Advisory	Mandatory	<code>fcntl</code>	<code>lockf</code>	<code>flock</code>
SUS	•		•	XSI	
FreeBSD 8.0	•		•	•	•
Linux 3.2.0	•	•	•	•	•
Mac OS X 10.6.8	•		•	•	•
Solaris 10	•	•	•	•	•

**Figure:** Forms of record locking supported by various UNIX systems





# fcntl Record Locking

## Prototype for the fcntl function

```
#include <fcntl.h>
int fcntl(int fd, int cmd, ... /* struct flock *flockptr */ );
// Returns: depends on cmd if OK (see following), 1 on error
```

- *cmd*: F\_GETLK, F\_SETLK, F\_SETLKW
- *flockptr*: is a pointer to an flock structure

```
struct flock {
    short  l_type;    /* F_RDLCK, F_WRLCK, or F_UNLCK */
    short  l_whence;  /* SEEK_SET, SEEK_CUR, or SEEK_END */
    off_t   l_start;   /* offset in bytes, relative to l_whence */
    off_t   l_len;     /* length, in bytes; 0 means lock to EOF */
    pid_t   l_pid;     /* returned with F_GETLK */
};
```



# fcntl Record Locking cnt'd

Rules in specification of the region to be locked or unlocked

- The two elements that specify the starting offset of the region are similar to the last two arguments of the `lseek` function.
  - the `l_whence` member is specified as `SEEK_SET`, `SEEK_CUR`, or `SEEK_END`.
- Locks can start and extend beyond the current end of file, but cannot start or extend before the beginning of the file.
- If `l_len` is 0, it means that the lock extends to the largest possible offset of the file.
- To lock the entire file, we set `l_start` and `l_whence` to point to the beginning of the file and specify a length (`l_len`) of 0. (There are several ways to specify the beginning of the file, but most applications specify `l_start` as 0 and `l_whence` as `SEEK_SET`.)



# fcntl Record Locking cnt'd

The basic rule

- any number of processes can have a shared read lock on a given byte
- only one process can have an exclusive write lock on a given byte.

Furthermore, if there are one or more read locks on a byte, there can't be any write locks on that byte;

if there is an exclusive write lock on a byte, there can't be any read locks on that byte.



		Request for	
		read lock	write lock
Region currently has	no locks	OK	OK
	one or more read locks	OK	denied
	one write lock	denied	denied

**Figure:** Compatibility between different lock types

If a process has an existing lock on a range of a file, a subsequent attempt to place a lock on the same range by the same process **will replace the existing lock with the new one.**



# fcntl Record Locking cnt'd

## Three commands for the fcntl function

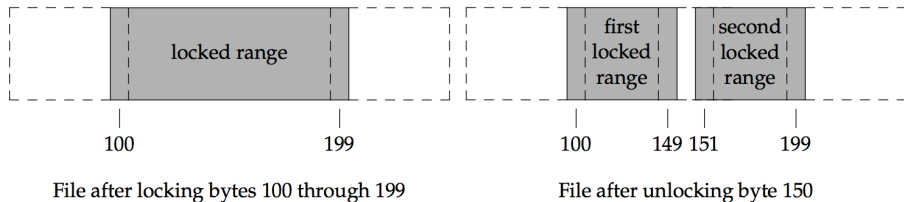
- F\_GETLK: Determine whether the lock described by *flockptr* is blocked by some other lock.
  - If a lock exists, the information on that existing lock overwrites the information pointed to by *flockptr*.
  - If no lock exists, the structure pointed to by *flockptr* is left unchanged except for the *l\_type* member, which is set to F\_UNLCK.
- F\_SETLK: Set the lock described by *flockptr*.
  - If the compatibility rule prevents the system from giving us the lock, fcntl returns immediately with *errno* set to either EACCES or EAGAIN.
- F\_SETLKW: a blocking version of F\_SETLK.
  - If the requested lock cannot be granted the calling process is put to sleep. The process wakes up either when the lock becomes available or when interrupted by a signal.



# fcntl Record Locking cnt'd

Be aware that testing for a lock with `F_GETLK` and then trying to obtain that lock with `F_SETLK` or `F_SETLKW` is not an atomic operation.

When setting or releasing a lock on a file, the system combines or splits adjacent areas as required.



**Figure:** File byte-range lock diagram



# Example: Requesting and Releasing a Lock

## Definition of lock\_reg

```
#include <fcntl.h>
int
lock_reg(int fd, int cmd, int type, off_t offset, int whence, off_t len)
{
    struct flock lock;
    lock.l_type = type; /* F_RDLCK, F_WRLCK, F_UNLCK */
    lock.l_start = offset; /* byte offset, relative to l_whence */
    lock.l_whence = whence; /* SEEK_SET, SEEK_CUR, SEEK_END */
    lock.l_len = len; /* #bytes (0 means to EOF) */
    return(fcntl(fd, cmd, &lock));
}
```



# Example: Requesting and Releasing a Lock

Use following five macros

```
#define read_lock(fd, offset, whence, len) \  
    lock_reg((fd), F_SETLK, F_RDLCK, (offset), (whence), (len))  
#define readw_lock(fd, offset, whence, len) \  
    lock_reg((fd), F_SETLKW, F_RDLCK, (offset), (whence), (len))  
#define write_lock(fd, offset, whence, len) \  
    lock_reg((fd), F_SETLK, F_WRLCK, (offset), (whence), (len))  
#define writew_lock(fd, offset, whence, len) \  
    lock_reg((fd), F_SETLKW, F_WRLCK, (offset), (whence), (len))  
#define un_lock(fd, offset, whence, len) \  
    lock_reg((fd), F_SETLK, F_UNLCK, (offset), (whence), (len))
```





# Example: Requesting and Releasing a Lock I

```
#include <fcntl.h>
#include <stdlib.h>

pid_t lock_test(int fd, int type, off_t offset, int whence, off_t len)
{
    struct flock lock;
    lock.l_type = type; /* F_RDLCK or F_WRLCK */
    lock.l_start = offset; /* byte offset, relative to l_whence */
    lock.l_whence = whence; /* SEEK_SET, SEEK_CUR, SEEK_END */
    lock.l_len = len; /* #bytes (0 means to EOF) */

    if (fcntl(fd, F_GETLK, &lock) < 0){
        fprintf(stderr, "fcntl error");
        exit(1);
    }

    if (lock.l_type == F_UNLCK)
        return(0); /* false, region isn't locked by another proc */
    return(lock.l_pid); /* true, return pid of lock owner */
}
```



# Example: Requesting and Releasing a Lock

- If a lock exists, this function returns the process ID of the process holding the lock.
- Otherwise, the function returns 0 (false).

## macros

```
#define is_read_lockable(fd, offset, whence, len) \  
    (lock_test((fd), F_RDLCK, (offset), (whence), (len)) == 0)  
#define is_write_lockable(fd, offset, whence, len) \  
    (lock_test((fd), F_WRLCK, (offset), (whence), (len)) == 0)
```



# Example: codes/deadlock.c |

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <fcntl.h>
5
6  // tellwait.c
7  void TELL_WAIT(void); /* parent/child from {Sec race_conditions} */
8  void TELL_PARENT(pid_t);
9  void TELL_CHILD(pid_t);
10 void WAIT_PARENT(void);
11 void WAIT_CHILD(void);
12
13 #define FILE_MODE (S_IRUSR | S_IWUSR | S_IRGRP | S_IROTH)
14
15 #define read_lock(fd, offset, whence, len) \
16     lock_reg((fd), F_SETLK, F_RDLCK, (offset), (whence), (len))
17 #define readw_lock(fd, offset, whence, len) \
18     lock_reg((fd), F_SETLKW, F_RDLCK, (offset), (whence), (len))
19 #define write_lock(fd, offset, whence, len) \
20     lock_reg((fd), F_SETLK, F_WRLCK, (offset), (whence), (len))
21 #define writew_lock(fd, offset, whence, len) \
22     lock_reg((fd), F_SETLKW, F_WRLCK, (offset), (whence), (len))
```



# Example: codes/deadlock.c II

```
23 #define un_lock(fd, offset, whence, len) \  
24     lock_reg((fd), F_SETLK, F_UNLCK, (offset), (whence), (len))  
25  
26 int  
27 lock_reg(int fd, int cmd, int type, off_t offset, int whence, off_t len)  
28 {  
29     struct flock lock;  
30     lock.l_type = type; /* F_RDLCK, F_WRLCK, F_UNLCK */  
31     lock.l_start = offset; /* byte offset, relative to l_whence */  
32     lock.l_whence = whence; /* SEEK_SET, SEEK_CUR, SEEK_END */  
33     lock.l_len = len; /* #bytes (0 means to EOF) */  
34     return(fcntl(fd, cmd, &lock));  
35 }  
36  
37 pid_t  
38 lock_test(int fd, int type, off_t offset, int whence, off_t len)  
39 {  
40     struct flock lock;  
41     lock.l_type = type; /* F_RDLCK or F_WRLCK */  
42     lock.l_start = offset; /* byte offset, relative to l_whence */  
43     lock.l_whence = whence; /* SEEK_SET, SEEK_CUR, SEEK_END */  
44     lock.l_len = len; /* #bytes (0 means to EOF) */  
45
```



## Example: codes/deadlock.c III

```
46     if (fcntl(fd, F_GETLK, &lock) < 0){
47         fprintf(stderr, "fcntl error");
48         exit(1);
49     }
50
51     if (lock.l_type == F_UNLCK)
52         return(0); /* false, region isn't locked by another proc */
53     return(lock.l_pid); /* true, return pid of lock owner */
54 }
55
56 static void
57 lockabyte(const char *name, int fd, off_t offset)
58 {
59     if (writew_lock(fd, offset, SEEK_SET, 1) < 0){
60         fprintf(stderr, "%s: writew_lock error", name);
61         exit(1);
62     }
63     printf("%s: got the lock, byte %lld\n", name, (long long)offset);
64 }
65
66 int
67 main(void)
68 {
```



## Example: codes/deadlock.c IV

```
69  int fd;
70  pid_t pid;
71
72  /*
73   * Create a file and write two bytes to it.
74   */
75  if ((fd = creat("templock", FILE_MODE)) < 0){
76      fprintf(stderr, "creat error");
77      exit(1);
78  }
79  if (write(fd, "ab", 2) != 2){
80      fprintf(stderr, "write error");
81      exit(1);
82  }
83
84  TELL_WAIT();
85  if ((pid = fork()) < 0) {
86      fprintf(stderr, "fork error");
87      exit(1);
88  } else if (pid == 0) {    /* child */
89      lockabyte("child", fd, 0);
90      TELL_PARENT(getppid());
91      WAIT_PARENT();
```



## Example: codes/deadlock.c V

```
92     lockabyte("child", fd, 1);
93 } else {      /* parent */
94     lockabyte("parent", fd, 1);
95     TELL_CHILD(pid);
96     WAIT_CHILD();
97     lockabyte("parent", fd, 0);
98 }
99     exit(0);
100 }
```

When a deadlock is detected, the kernel has to choose one process to receive the error return.

```
James@maker:codes$ ./deadlock
```

```
parent: got the lock, byte 1
```

```
child: got the lock, byte 0
```

```
child: writew_lock errorparent: got the lock, byte 0
```



# implied Inheritance and Release of Locks

Three rules govern the automatic inheritance and release of record locks

- Locks are associated with a process and a file.
  - when a process terminates, all its locks are released
  - whenever a descriptor is closed, any locks on the file referenced by that descriptor for that process are released.

## Example

After the `close(fd2)`, the lock on `fd1` is released

```
fd1 = open(pathname, ...);  
read_lock(fd1, ...);  
fd2 = dup(fd1);  
close(fd2);
```

The same would happen if we replace `dup` with `open`

```
fd1 = open(pathname, ...);  
read_lock(fd1, ...);  
// if pathname is same  
fd2 = open(pathname, ...);  
close(fd2);
```





# implied Inheritance and Release of Locks cnt'd

Three rules govern the automatic inheritance and release of record locks

- Locks are associated with a process and a file.
- Locks are never inherited by the child across a fork.
  - The child has to call `fcntl` to obtain its own locks on any descriptors that were inherited across the fork.
  - locks are meant to prevent multiple processes from writing to the same file at the same time.
- Locks are inherited by a new program across an exec.



# FreeBSD Implementation

To understand rule 1

```
fd1 = open(pathname, ...);
write_lock(fd1, 0, SEEK_SET, 1);
if ((pid = fork()) > 0) {
    fd2 = dup(fd1);
    fd3 = open(pathname, ...);
} else if (pid == 0) {
    /* parent write locks byte 0 */
    /* parent */
    read_lock(fd1, 1, SEEK_SET, 1); /* child read locks byte 1 */
}
pause();
```



# FreeBSD Implementation cnt'd

- lockf structures that are linked together from the i-node structure
- Each lockf structure describes one locked region for a given process
- In the parent, closing any one of fd1, fd2, or fd3 causes the parent's lock to be released. When any one of these three file descriptors is closed, the kernel goes through the linked list of locks for the corresponding i-node and releases the locks held by the calling process.

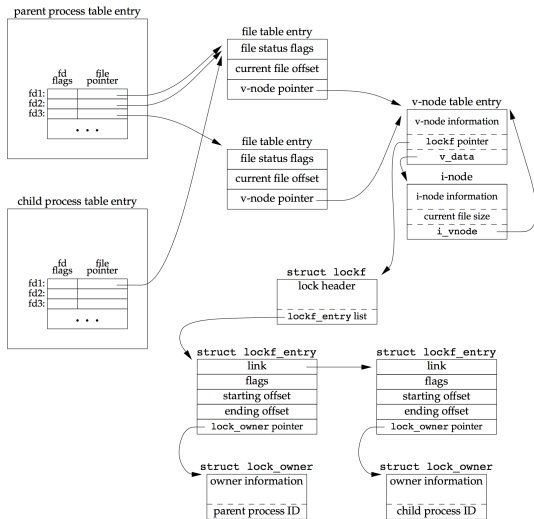


Figure: The FreeBSD data structures for record locking

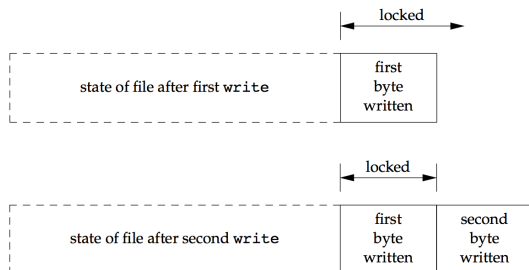


# Locks at End of File

We need to use caution when locking or unlocking byte ranges relative to the end of file.

- Most implementations convert an `l_whence` value of `SEEK_CUR` or `SEEK_END` into an absolute file offset, using `l_start` and the file's current position or current length.

```
writew_lock(fd, 0, SEEK_END, 0);  
write(fd, buf, 1);  
un_lock(fd, 0, SEEK_END);  
write(fd, buf, 1);
```



**Figure:** File range lock diagram



# Advisory versus Mandatory Locking

## Cooperating processes

If all the functions in the library handle record locking in a consistent way, then we say that any set of processes using these functions to access a database are *cooperating processes*.

It is feasible for these database access functions to use **advisory locking**, if they are the only ones being used to access the database.

**Mandatory Locking** or *enforcement-mode locking* causes the kernel to check open, read, write to verify that the calling process isn't violating a lock on the file being accessed



# ASYNCHRONOUS I/O

---

# Asynchronous I/O

## The cost of using Asynchronous I/O

- we complicate the design of our application by choosing to juggle multiple concurrent operations.
- We have to worry about three sources of errors for every asynchronous operation
  - one associated with the submission of the operation,
  - one associated with the result of the operation itself,
  - and one associated with the functions used to determine the status of the asynchronous operations.
- The interfaces themselves involve a lot of extra setup and processing rules compared to their conventional counterparts
- Recovering from errors can be difficult.



# POSIX Asynchronous I/O

The POSIX asynchronous I/O interfaces give us a consistent way to perform asynchronous I/O, regardless of the type of file. The

asynchronous I/O interfaces use AIO control blocks to describe I/O operations.

```
struct aiocb {  
    int          aio_fildes;      /* file descriptor */  
    off_t        aio_offset;      /* file offset fro I/O */  
    volatile void *aio_buf;       /* buffer for I/O */  
    size_t       aio_nbytes;      /* number of bytes to transfer */  
    int          aio_reqprio;     /* priority */  
    struct sigevent aio_sigevent; /* Signal information */  
    int          aio_lio_opcode;  /* operation for list I/O */  
};
```





# POSIX Asynchronous I/O cnt'd

The `aio_sigevent` field controls how the application is notified about the completion of the I/O event.

```
struct sigevent {  
    int     sigev_notify;           /* notify type */  
    int     sigev_signo;           /* signal number */  
    union sigval  sigev_value;     /* notify argument */  
    void (*sigev_notify_function)(union sigval); /* notify function */  
    pthread_attr_t *sigev_notify_attributes; /* notify attrs */  
};
```

`sigev_notify` field controls the type of notification

1. `SIGEV_NONE`: The process is not notified when the asynchronous I/O request completes.
2. `SIGEV_SIGNAL`: The signal specified by the `sigev_signo` field is generated when the asynchronous I/O request completes.
3. `SIGEV_THREAD`: The function specified by the `sigev_notify_function` field is called when the asynchronous I/O request completes. The function is executed in a separate thread



# POSIX Asynchronous I/O: Read and Write

To perform asynchronous I/O, we need to initialize an AIO control block

```
#include <aio.h>
int aio_read(struct aiocb *aiocb); int aio_write(struct aiocb *aiocb);
// Both return: 0 if OK, 1 on error
```

When these functions return success,

- the asynchronous I/O request has been queued for processing by the operating system.

The return value bears no relation to the result of the actual I/O operation



# POSIX Asynchronous I/O: Persistency

To force all pending asynchronous writes to persistent storage without waiting

```
#include <aio.h>
int aio_fsync(int op, struct aiocb *aiocb);
// Returns: 0 if OK, 1 on error
```

The `aio_fildes` field in the AIO control block indicates the file whose asynchronous writes are synched. If the `op` argument is set to

- `O_DSYNC`: the operation behaves like a call to `fdatasync`.
- `O_SYNC`: the operation behaves like a call to `fsync`.

Returns when the synch is scheduled. The data won't be persistent until the asynchronous synch completes.



# POSIX Asynchronous I/O: Status

To determine the completion status of an asynchronous read, write, or synch operation

```
#include <aio.h>
int aio_error(const struct aiocb *aiocb);
Returns: 0, -1, EINPROGRESS, others
```

## Four return types

- 0: The asynchronous operation completed successfully.
- -1: The call to `aio_error` failed
- EINPROGRESS: The asynchronous read, write, or synch is still pending
- *anything else*: error code corresponding to the failed asynchronous operation.



# POSIX Asynchronous I/O: Success

If the asynchronous operation succeeded, we can call the `aio_return` function to get the asynchronous operation's return value

```
#include <aio.h>
ssize_t aio_return(const struct aiocb *aiocb);
// Returns: ...
```

**Until the asynchronous operation completes, we need to be careful to avoid calling the `aio_return` function**

- Results undefined until the operation completes
- call `aio_return` only one time per asynchronous I/O operation



# POSIX Asynchronous I/O: suspend

When we have completed the processing and find that we still have asynchronous operations outstanding, we can call the `aio_suspend` function to block until an operation completes.

```
#include <aio.h>
int aio_suspend(const struct aiocb *const list[], int nent,
               const struct timespec *timeout);
// Returns: 0 if OK, 1 on error
```

Three things cause `aio_suspend` to return

- -1: Interrupted by a signal, `errno` is set to `EINTR`
- -1: *timeout* argument expires without any of the I/O operations completing, `errno` is set to `EAGAIN`
- 0: if any of the I/O operations complete

*list* is a pointer to an array of AIO control block

*nent* indicates the number of entries in the array



# POSIX Asynchronous I/O: Cancel

When we have pending asynchronous I/O operations that we no longer want to complete, we can attempt to cancel them with the `aio_cancel` function.

```
#include <aio.h>
int aio_cancel(int fd, /* file descriptor */
               struct aiocb *aiocb); /* NULL: all outstanding AIO */
// Returns: (see following)
```

## Four return values

- `AIO_ALLDONE`: All done before attempt to cancel
- `AIO_CANCELED`: All are canceled
- `AIO_NOTCANCELED`: At least one of the requested operation could not be canceled
- `-1`: Failed. Error code stored in `errno`



# POSIX Asynchronous I/O: set of I/O requests

The `lio_listio` function submits a set of I/O requests described by a list of AIO control blocks.

```
#include <aio.h>
int lio_listio(int mode, struct aiocb *restrict const list[restrict],
               int nent, struct sigevent *restrict sigev);
// Returns: 0 if OK, 1 on error
```

- `mode`
  - `LIO_WAIT`: won't return until all I/Os in the list are complete (`sigev` is ignored)
  - `LIO_NOWAIT`: returns as soon as the I/O requests are queued. Process is notified asynchronously when all of the I/O operations complete
- `sigev`: can be set to `NULL` if don't want to be notified
- `list`: list of AIO control blocks specifying the I/O operations to perform
- `nent`: number of elements in the array





# Example: codes/rot13a.c |

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <ctype.h>
5  #include <fcntl.h>
6
7  #define BSZ 4096
8  #define FILE_MODE (S_IRUSR | S_IWUSR | S_IRGRP | S_IROTH)
9
10 unsigned char buf[BSZ];
11
12 unsigned char
13 translate(unsigned char c)
14 {
15     if (isalpha(c)) {
16         if (c >= 'n')
17             c -= 13;
18         else if (c >= 'a')
19             c += 13;
20         else if (c >= 'N')
21             c -= 13;
22         else
```



## Example: codes/rot13a.c II

```
23         c += 13;
24     }
25     return(c);
26 }
27
28 int
29 main(int argc, char* argv[])
30 {
31     int ifd, ofd, i, n, nw;
32
33     if (argc != 3){
34         fprintf(stderr, "usage: rot13 infile outfile");
35         exit(1);
36     }
37     if ((ifd = open(argv[1], O_RDONLY)) < 0){
38         fprintf(stderr, "can't open %s", argv[1]);
39         exit(1);
40     }
41     if ((ofd = open(argv[2], O_RDWR|O_CREAT|O_TRUNC, FILE_MODE)) < 0){
42         fprintf(stderr, "can't create %s", argv[2]);
43         exit(1);
44     }
45
```



## Example: codes/rot13a.c III

```
46     while ((n = read(ifd, buf, BSZ)) > 0) {
47         for (i = 0; i < n; i++)
48             buf[i] = translate(buf[i]);
49         if ((nw = write(ofd, buf, n)) != n) {
50             if (nw < 0){
51                 fprintf(stderr, "write failed");
52                 exit(1);
53             }
54             else {
55                 fprintf(stderr, "short write (%d/%d)", nw, n);
56                 exit(1);
57             }
58         }
59     }
60
61     fsync(ofd);
62     exit(0);
63 }
```



## Example: codes/rot13c2.c |

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <ctype.h>
5  #include <fcntl.h>
6  #include <aio.h>
7  #include <errno.h>
8  #include <sys/stat.h>
9
10 #define BSZ 4096
11 #define NBUF 8
12 #define FILE_MODE (S_IRUSR | S_IWUSR | S_IRGRP | S_IROTH)
13
14 enum rwop {
15     UNUSED = 0,
16     READ_PENDING = 1,
17     WRITE_PENDING = 2
18 };
19
20 struct buf {
21     enum rwop op;
```



## Example: codes/rot13c2.c II

```
22     int last;
23     struct aiocb aiocb;
24     unsigned char data[BSZ];
25 };
26
27 struct buf bufs[NBUF];
28
29 unsigned char
30 translate(unsigned char c)
31 {
32     if (isalpha(c)) {
33         if (c >= 'n')
34             c -= 13;
35         else if (c >= 'a')
36             c += 13;
37         else if (c >= 'N')
38             c -= 13;
39         else
40             c += 13;
41     }
42     return(c);
43 }
44
```



## Example: codes/rot13c2.c III

```
45 int
46 main(int argc, char* argv[])
47 {
48     int ifd, ofd, i, j, n, err, numop;
49     struct stat sbuf;
50     const struct aiocb *aiolist[NBUF];
51     off_t off = 0;
52
53     if (argc != 3){
54         fprintf(stderr, "usage: rot13 infile outfile");
55         exit(1);
56     }
57     if ((ifd = open(argv[1], O_RDONLY)) < 0){
58         fprintf(stderr, "can't open %s", argv[1]);
59         exit(1);
60     }
61     if ((ofd = open(argv[2], O_RDWR | O_CREAT | O_TRUNC, FILE_MODE)) < 0){
62         fprintf(stderr, "can't create %s", argv[2]);
63         exit(1);
64     }
65     if (fstat(ifd, &sbuf) < 0){
```



## Example: codes/rot13c2.c IV

```
66     fprintf(stderr, "fstat failed");
67     exit(1);
68 }
69
70 /* initialize the buffers */
71 for (i = 0; i < NBUF; i++) {
72     bufs[i].op = UNUSED;
73     bufs[i].aiocb.aio_buf = bufs[i].data;
74     bufs[i].aiocb.aio_sigevent.sigev_notify = SIGEV_NONE;
75     aiolist[i] = NULL;
76 }
77
78 numop = 0;
79 for (;;) {
80     for (i = 0; i < NBUF; i++) {
81         switch (bufs[i].op) {
82             case UNUSED:
83                 /*
84                  * Read from the input file if more data
85                  * remains unread.
86                  */
87                 if (off < sbuf.st_size) {
```



## Example: codes/rot13c2.c V

```
88         bufs[i].op = READ_PENDING;
89         bufs[i].aiocb.aio_fildes = ifd;
90         bufs[i].aiocb.aio_offset = off;
91         off += BSZ;
92         if (off >= sbuf.st_size)
93             bufs[i].last = 1;
94         bufs[i].aiocb.aio_nbytes = BSZ;
95         if (aio_read(&bufs[i].aiocb) < 0){
96             fprintf(stderr, "aio_read failed");
97             exit(1);
98         }
99         aiolist[i] = &bufs[i].aiocb;
100        numop++;
101    }
102    break;
103
104    case READ_PENDING:
105        if ((err = aio_error(&bufs[i].aiocb)) == EINPROGRESS)
106            continue;
107        if (err != 0) {
108            if (err == -1){
109                fprintf(stderr, "aio_error failed");
```





## Example: codes/rot13c2.c VI

```
110         exit(1);
111     }
112     else{
113         fprintf(stderr, "read failed");
114         exit(1);
115     }
116 }
117
118 /*
119  * A read is complete; translate the buffer
120  * and write it.
121  */
122 if ((n = aio_return(&bufs[i].aiocb)) < 0){
123     fprintf(stderr, "aio_return failed");
124     exit(1);
125 }
126 if (n != BSZ && !bufs[i].last){
127     fprintf(stderr, "short read (%d/%d)", n, BSZ);
128     exit(1);
129 }
130 for (j = 0; j < n; j++)
```



## Example: codes/rot13c2.c VII

```
131         bufs[i].data[j] = translate(bufs[i].data[j]);
132     bufs[i].op = WRITE_PENDING;
133     bufs[i].aiocb.aio_fildes = ofd;
134     bufs[i].aiocb.aio_nbytes = n;
135     if (aio_write(&bufs[i].aiocb) < 0){
136         fprintf(stderr, "aio_write failed");
137         exit(1);
138     }
139     /* retain our spot in aiolist */
140     break;
141
142 case WRITE_PENDING:
143     if ((err = aio_error(&bufs[i].aiocb)) == EINPROGRESS)
144         continue;
145     if (err != 0) {
146         if (err == -1){
147             fprintf(stderr, "aio_error failed");
148             exit(1);
149         }
150         else{
151             fprintf(stderr, "write failed");
```



## Example: codes/rot13c2.c VIII

```
152         exit(1);
153     }
154 }
155
156 /*
157  * A write is complete; mark the buffer as unused.
158  */
159 if ((n = aio_return(&bufs[i].aiocb)) < 0){
160     fprintf(stderr, "aio_return failed");
161     exit(1);
162 }
163 if (n != bufs[i].aiocb.aio_nbytes){
164     fprintf(stderr, "short write (%d/%d)", n, BSZ);
165     exit(1);
166 }
167 aiolist[i] = NULL;
168 bufs[i].op = UNUSED;
169 numop--;
170 break;
171 }
172 }
```



## Example: codes/rot13c2.c IX

```
173     if (numop == 0) {
174         if (off >= sbuf.st_size)
175             break;
176     } else {
177         if (aio_suspend(aiolist, NBUF, NULL) < 0){
178             fprintf(stderr, "aio_suspend failed");
179             exit(1);
180         }
181     }
182 }
183
184 bufs[0].aiocb.aio_fildes = ofd;
185 if (aio_fsync(O_SYNC, &bufs[0].aiocb) < 0){
186     fprintf(stderr, "aio_fsync failed");
187     exit(1);
188 }
189 exit(0);
190 }
```



## readv AND writev FUNCTIONS

---

# 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*.

```
#include <sys/uio.h>
ssize_t readv(int fd, const struct iovec *iov, int iovcnt);
ssize_t writev(int fd, const struct iovec *iov, int iovcnt);
// Both return: number of bytes read or written, 1 on error
```

iov is a pointer to an array of iovec structures

```
struct iovec {
    void    *iov_base; /* starting address of buffer */
    size_t  iov_len;   /* size of buffer */
};
```

The number of elements in the iov array is specified iovcnt



# readv and writev Functions

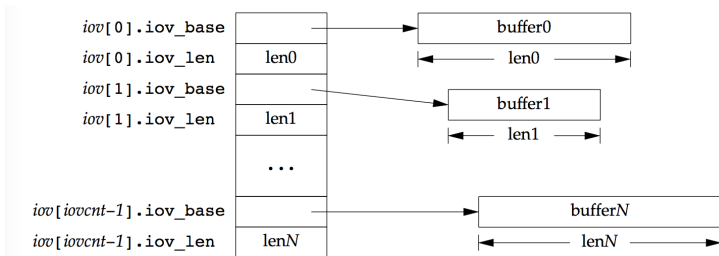


Figure: the iovec structure for readv and writev

- 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.



readn AND writen FUNCTIONS

---



# readn and writen Functions

Pipes, FIFOs, and some devices—notably terminals and networks—have the following two properties.

1. A read operation may return less than asked for, even though we have not encountered the end of file.
  - This is not an error, and we should simply continue reading from the device.
2. A write operation can return less than we specified. This may be caused by kernel output buffers becoming full.
  - it's not an error, and we should continue writing the remainder of the data.

Generally, when we read from or write to a pipe, network device, or terminal, we need to take these characteristics into consideration.



# readn and writen Functions

We can use the `readn` and `writen` functions to read and write `N` bytes of data, respectively, letting these functions handle a return value that's possibly less than requested.

- These two functions simply call `read` or `write` as many times as required to read or write the entire `N` bytes of data.

```
ssize_t readn(int fd, void *buf, size_t nbytes);  
ssize_t writen(int fd, void *buf, size_t nbytes);  
// Both return: number of bytes read or written, 1 on error
```



# readn Functions

```
1  ssize_t /* Read "n" bytes from a descriptor */
2  readn(int fd, void *ptr, size_t n)
3  {
4      size_t nleft;
5      ssize_t nread;
6
7      nleft = n;
8      while (nleft > 0) {
9          if ((nread = read(fd, ptr, nleft)) < 0) {
10             if (nleft == n)
11                 return(-1); /* error, return -1 */
12             else
13                 break; /* error, return amount read so far */
14         } else if (nread == 0) {
15             break; /* EOF */
16         }
17         nleft -= nread;
18         ptr += nread;
19     }
20     return(n - nleft); /* return >= 0 */
21 }
```



# written Functions

```
1  ssize_t /* Write "n" bytes to a descriptor */
2  written(int fd, const void *ptr, size_t n)
3  {
4      size_t nleft;
5      ssize_t nwritten;
6
7      nleft = n;
8      while (nleft > 0) {
9          if ((nwritten = write(fd, ptr, nleft)) < 0) {
10             if (nleft == n)
11                 return(-1); /* error, return -1 */
12             else
13                 break; /* error, return amount written so far */
14         } else if (nwritten == 0) {
15             break;
16         }
17         nleft -= nwritten;
18         ptr += nwritten;
19     }
20     return(n - nleft); /* return >= 0 */
21 }
```



# MEMORY-MAPPED I/O

---

# 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`.



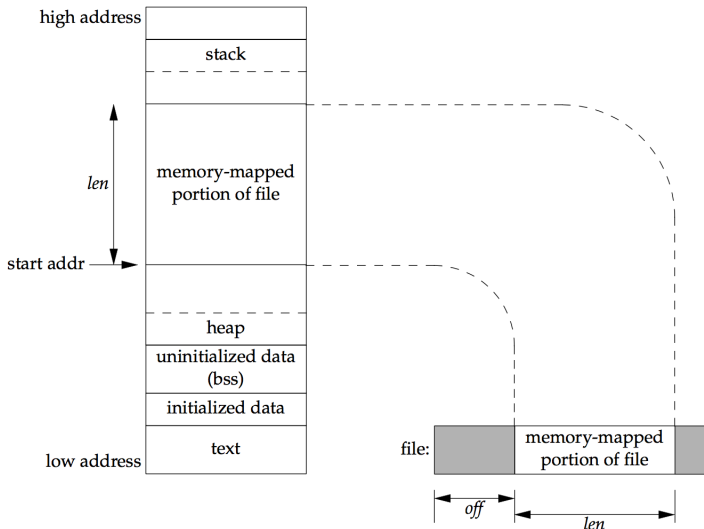
# Memory-Mapped I/O

```
#include <sys/mman.h>
void *mmap(void *addr, size_t len, int prot, int flag, int fd, off_t off );
// Returns: starting address of mapped region if OK, MAP_FAILED on error
```

- addr: specify the address where we want the mapped region to start (Normally set to 0 to let system choose)
- fd: file descriptor specifying the file that is to be mapped (file must be opened before mapping)
- len: number of bytes to map
- off: starting offset in the file of the bytes to map
- prot: READ, WRITE, and EXEC can be ORed
  - PROT\_READ: Region can be read
  - PROT\_WRITE: Region can be written
  - PROT\_EXEC: Region can be executed
  - PROT\_NONE: Region cannot be accessed



# Memory-Mapped I/O



**Figure:** Example of a memory-mapped file





# Memory-Mapped I/O

`flag` argument affects various attributes of the mapped region

- `MAP_FIXED`: The return value must equal `addr`. Use of this flag is discouraged (Portability issue)
- `MAP_SHARED`: store operations into the mapped region by this process
- `MAP_PRIVATE`: store operation 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



# Memory-Mapped I/O

The value of `off` and the value of `addr` are usually required to be multiples of the system's virtual memory page size. What happens if the length of the mapped region isn't a multiple of the page size?

- File size 12 bytes
- Page size is 512 bytes
- system provides a mapped region of 512 bytes
- final 500 bytes of the this region are set to 0
- any changes to the final 500 bytes are not reflected to the file



# Memory-Mapped I/O cnt'd

- A memory-mapped region is inherited by a child across a fork
- But, is not inherited by the new program across an exec.

Permission can be changed by calling mprotect

```
#include <sys/mman.h>
int mprotect(void *addr, size_t len, int prot);
// Returns: 0 if OK, 1 on error
```



# Memory-Mapped I/O cnt'd

When we modify pages that we've mapped into our address space using the MAP\_SHARED flag, the changes aren't written back to the file immediately. kernel daemons decide when dirty pages are written

back based on

- system load
- configuration parameters meant to limit data loss in the event of a system failure

**When the changes are written back, they are written in units of pages.**



# Memory-Mapped I/O cnt'd

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.

```
#include <sys/mman.h>
int msync(void *addr, size_t len, int flags);
// Returns: 0 if OK, 1 on error
```

Two values for flags

- `MS_SYNC`: wait for the writes to complete before returning
- `MS_ASYNC`: simply schedule the pages to be written



# Memory-Mapped I/O cnt'd

A memory-mapped region is automatically unmapped when the process terminates or we can unmap a region directly by calling the `munmap` function.

```
#include <sys/mman.h>
int munmap(void *addr, size_t len);
// Returns: 0 if OK, 1 on error
```



# Example: CP with MMAP codes/mcopy2.c |

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <string.h>
5  #include <fcntl.h>
6  #include <sys/mman.h>
7  #include <sys/stat.h>
8
9  #define COPYINCR (1024*1024*1024) /* 1 GB */
10 #define FILE_MODE (S_IRUSR | S_IWUSR | S_IRGRP | S_IROTH)
11
12 int
13 main(int argc, char *argv[])
14 {
15     int fdin, fdout;
16     void *src, *dst;
17     size_t copysz;
18     struct stat sbuf;
```



## Example: CP with MMAP codes/mcopy2.c II

```
19  off_t fsz = 0;
20
21  if (argc != 3){
22      fprintf(stderr, "usage: %s <fromfile> <tofile>", argv[0]);
23      exit(1);
24  }
25
26  if ((fdin = open(argv[1], O_RDONLY)) < 0){
27      fprintf(stderr, "can't open %s for reading", argv[1]);
28      exit(1);
29  }
30
31  if ((fdout = open(argv[2], O_RDWR | O_CREAT | O_TRUNC,
32      FILE_MODE)) < 0){
33      fprintf(stderr, "can't creat %s for writing", argv[2]);
34      exit(1);
35  }
36
37  if (fstat(fdin, &sbuf) < 0){ /* need size of input file */
```





## Example: CP with MMAP codes/mcopy2.c III

```
38     fprintf(stderr, "fstat error");
39     exit(1);
40 }
41
42 if (ftruncate(fdout, sbuf.st_size) < 0){ /* set output file size */
43     fprintf(stderr, "ftruncate error");
44     exit(1);
45 }
46
47 while (fsz < sbuf.st_size) {
48     if ((sbuf.st_size - fsz) > COPYINCR)
49         copysz = COPYINCR;
50     else
51         copysz = sbuf.st_size - fsz;
52
53     if ((src = mmap(0, copysz, PROT_READ, MAP_SHARED,
54         fdin, fsz)) == MAP_FAILED){
55         fprintf(stderr, "mmap error for input");
56         exit(1);
```



## Example: CP with MMAP codes/mcopy2.c IV

```
57     }
58     if ((dst = mmap(o, copysz, PROT_READ | PROT_WRITE,
59         MAP_SHARED, fdout, fsz)) == MAP_FAILED){
60         fprintf(stderr, "mmap error for output");
61         exit(1);
62     }
63
64     memcpy(dst, src, copysz); /* does the file copy */
65     munmap(src, copysz);
66     munmap(dst, copysz);
67     fsz += copysz;
68 }
69 exit(o);
70 }
```

```
dd if=/dev/urandom of=testfile count=1000000
time ./mcopy2 testfile mcopyout
time cp testfile mcopyout
```

