

**NAME**

`fcntl` - manipulate file descriptor

**SYNOPSIS**

```
#include <unistd.h>
```

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

```
int fcntl(int fd, int cmd, ... /* arg */ );
```

**DESCRIPTION**

`fcntl()` performs one of the operations described below on the open file descriptor *fd*. The operation is determined by *cmd*.

`fcntl()` can take an optional third argument. Whether or not this argument is required is determined by *cmd*. The required argument type is indicated in parentheses after each *cmd* name (in most cases, the required type is *int*, and we identify the argument using the name *arg*), or *void* is specified if the argument is not required.

Certain of the operations below are supported only since a particular Linux kernel version. The preferred method of checking whether the host kernel supports a particular operation is to invoke `fcntl()` with the desired *cmd* value and then test whether the call failed with **EINVAL**, indicating that the kernel does not recognize this value.

**Duplicating a file descriptor****F\_DUPFD** (*int*)

Find the lowest numbered available file descriptor greater than or equal to *arg* and make it be a copy of *fd*. This is different from [dup2\(2\)](#), which uses exactly the descriptor specified.

On success, the new descriptor is returned.

See [dup\(2\)](#) for further details.

**F\_DUPFD\_CLOEXEC** (*int*; since Linux 2.6.24)

As for **F\_DUPFD**, but additionally set the close-on-exec flag for the duplicate descriptor. Specifying this flag permits a program to avoid an additional `fcntl()` **F\_SETFD** operation to set the **FD\_CLOEXEC** flag. For an explanation of why this flag is useful, see the description of **O\_CLOEXEC** in [open\(2\)](#).

**File descriptor flags**

The following commands manipulate the flags associated with a file descriptor. Currently, only one such flag is defined: **FD\_CLOEXEC**, the close-on-exec flag. If the **FD\_CLOEXEC** bit is 0, the file descriptor will remain open across an [execve\(2\)](#), otherwise it will be closed.

**F\_GETFD** (*void*)

Read the file descriptor flags; *arg* is ignored.

**F\_SETFD** (*int*)

Set the file descriptor flags to the value specified by *arg*.

In multithreaded programs, using `fcntl()` **F\_SETFD** to set the close-on-exec flag at the same time as another thread performs a [fork\(2\)](#) plus [execve\(2\)](#) is vulnerable to a race condition that may unintentionally leak the file descriptor to the program executed in the child process. See the discussion of the **O\_CLOEXEC** flag in [open\(2\)](#) for details and a remedy to the problem.

**File status flags**

Each open file description has certain associated status flags, initialized by [open\(2\)](#) and possibly modified by `fcntl()`. Duplicated file descriptors (made with [dup\(2\)](#), `fcntl(F_DUPFD)`, [fork\(2\)](#), etc.) refer to the same open file description, and thus share the same file status flags.

The file status flags and their semantics are described in [open\(2\)](#).

**F\_GETFL** (*void*)

Get the file access mode and the file status flags; *arg* is ignored.

**F\_SETFL** (*int*)

Set the file status flags to the value specified by *arg*. File access mode (**O\_RDONLY**, **O\_WRONLY**, **O\_RDWR**) and file creation flags (i.e., **O\_CREAT**, **O\_EXCL**, **O\_NOCTTY**, **O\_TRUNC**) in *arg* are ignored. On Linux this command can change only the **O\_APPEND**, **O\_ASYNC**, **O\_DIRECT**, **O\_NOATIME**, and **O\_NONBLOCK** flags. It is not possible to change the **O\_DSYNC** and **O\_SYNC** flags; see **BUGS**, below.

**Advisory record locking**

Linux implements traditional (process-associated) UNIX record locks, as standardized by POSIX. For a Linux-specific alternative with better semantics, see the discussion of open file description locks below.

**F\_SETLCK**, **F\_SETLKW**, and **F\_GETLCK** are used to acquire, release, and test for the existence of record locks (also known as byte-range, file-segment, or file-region locks). The third argument, *lock*, is a pointer to a structure that has at least the following fields (in unspecified order).

```
struct flock {
    short l_type; /* Type of lock: F_RDLCK,
                 F_WRLCK, F_UNLCK */
    short l_whence; /* How to interpret l_start:
                   SEEK_SET, SEEK_CUR, SEEK_END */
    off_t l_start; /* Starting offset for lock */
    off_t l_len; /* Number of bytes to lock */
    pid_t l_pid; /* PID of process blocking our lock
                (set by F_GETLCK and F_OFD_GETLCK) */
};
```

The *l\_whence*, *l\_start*, and *l\_len* fields of this structure specify the range of bytes we wish to lock. Bytes past the end of the file may be locked, but not bytes before the start of the file.

*l\_start* is the starting offset for the lock, and is interpreted relative to either: the start of the file (if *l\_whence* is **SEEK\_SET**); the current file offset (if *l\_whence* is **SEEK\_CUR**); or the end of the file (if *l\_whence* is **SEEK\_END**). In the final two cases, *l\_start* can be a negative number provided the offset does not lie before the start of the file.

*l\_len* specifies the number of bytes to be locked. If *l\_len* is positive, then the range to be locked covers bytes *l\_start* up to and including *l\_start+l\_len-1*. Specifying 0 for *l\_len* has the special meaning: lock all bytes starting at the location specified by *l\_whence* and *l\_start* through to the end of file, no matter how large the file grows.

POSIX.1-2001 allows (but does not require) an implementation to support a negative *l\_len* value; if *l\_len* is negative, the interval described by *lock* covers bytes *l\_start+l\_len* up to and including *l\_start-1*. This is supported by Linux since kernel versions 2.4.21 and 2.5.49.

The *l\_type* field can be used to place a read (**F\_RDLCK**) or a write (**F\_WRLCK**) lock on a file. Any number of processes may hold a read lock (shared lock) on a file region, but only one process may hold a write lock (exclusive lock). An exclusive lock excludes all other locks, both shared and exclusive. A single process can hold only one type of lock on a file region; if a new lock is applied to an already-locked region, then the existing lock is converted to the new lock type. (Such conversions may involve splitting, shrinking, or coalescing with an existing lock if the byte range specified by the new lock does not precisely coincide with the range of the existing lock.)

**F\_SETLCK** (*struct flock \**)

Acquire a lock (when *l\_type* is **F\_RDLCK** or **F\_WRLCK**) or release a lock (when *l\_type* is **F\_UNLCK**) on the bytes specified by the *l\_whence*, *l\_start*, and *l\_len* fields of *lock*. If a conflicting lock is held by another process, this call returns -1 and sets *errno* to

**EACCES** or **EAGAIN**. (The error returned in this case differs across implementations, so POSIX requires a portable application to check for both errors.)

**F\_SETLKW** (*struct flock* \*)

As for **F\_SETLK**, but if a conflicting lock is held on the file, then wait for that lock to be released. If a signal is caught while waiting, then the call is interrupted and (after the signal handler has returned) returns immediately (with return value -1 and *errno* set to **EINTR**; see [signal\(7\)](#)).

**F\_GETLK** (*struct flock* \*)

On input to this call, *lock* describes a lock we would like to place on the file. If the lock could be placed, **fcntl()** does not actually place it, but returns **F\_UNLCK** in the *l\_type* field of *lock* and leaves the other fields of the structure unchanged.

If one or more incompatible locks would prevent this lock being placed, then **fcntl()** returns details about one of those locks in the *l\_type*, *l\_whence*, *l\_start*, and *l\_len* fields of *lock*. If the conflicting lock is a traditional (process-associated) record lock, then the *l\_pid* field is set to the PID of the process holding that lock. If the conflicting lock is an open file description lock, then *l\_pid* is set to -1. Note that the returned information may already be out of date by the time the caller inspects it.

In order to place a read lock, *fd* must be open for reading. In order to place a write lock, *fd* must be open for writing. To place both types of lock, open a file read-write.

When placing locks with **F\_SETLKW**, the kernel detects *deadlocks*, whereby two or more processes have their lock requests mutually blocked by locks held by the other processes. For example, suppose process A holds a write lock on byte 100 of a file, and process B holds a write lock on byte 200. If each process then attempts to lock the byte already locked by the other process using **F\_SETLKW**, then, without deadlock detection, both processes would remain blocked indefinitely. When the kernel detects such deadlocks, it causes one of the blocking lock requests to immediately fail with the error **EDEADLK**; an application that encounters such an error should release some of its locks to allow other applications to proceed before attempting to regain the locks that it requires. Circular deadlocks involving more than two processes are also detected. Note, however, that there are limitations to the kernel's deadlock-detection algorithm; see **BUGS**.

As well as being removed by an explicit **F\_UNLCK**, record locks are automatically released when the process terminates.

Record locks are not inherited by a child created via [fork\(2\)](#), but are preserved across an [execve\(2\)](#).

Because of the buffering performed by the [stdio\(3\)](#) library, the use of record locking with routines in that package should be avoided; use [read\(2\)](#) and [write\(2\)](#) instead.

The record locks described above are associated with the process (unlike the open file description locks described below). This has some unfortunate consequences:

- \* If a process closes *any* file descriptor referring to a file, then all of the process's locks on that file are released, regardless of the file descriptor(s) on which the locks were obtained. This is bad: it means that a process can lose its locks on a file such as */etc/passwd* or */etc/mtab* when for some reason a library function decides to open, read, and close the same file.
- \* The threads in a process share locks. In other words, a multithreaded program can't use record locking to ensure that threads don't simultaneously access the same region of a file.

Open file description locks solve both of these problems.

### Open file description locks (non-POSIX)

Open file description locks are advisory byte-range locks whose operation is in most respects identical to the traditional record locks described above. This lock type is Linux-specific, and available since Linux 3.15. For an explanation of open file descriptions, see [open\(2\)](#).

The principal difference between the two lock types is that whereas traditional record locks are associated with a process, open file description locks are associated with the open file description on which they are acquired, much like locks acquired with [flock\(2\)](#). Consequently (and unlike traditional advisory record locks), open file description locks are inherited across [fork\(2\)](#) (and [clone\(2\)](#) with **CLONE\_FILES**), and are only automatically released on the last close of the open file description, instead of being released on any close of the file.

Open file description locks always conflict with traditional record locks, even when they are acquired by the same process on the same file descriptor.

Open file description locks placed via the same open file description (i.e., via the same file descriptor, or via a duplicate of the file descriptor created by [fork\(2\)](#), [dup\(2\)](#), [fcntl\(2\)](#) **F\_DUPFD**, and so on) are always compatible: if a new lock is placed on an already locked region, then the existing lock is converted to the new lock type. (Such conversions may result in splitting, shrinking, or coalescing with an existing lock as discussed above.)

On the other hand, open file description locks may conflict with each other when they are acquired via different open file descriptions. Thus, the threads in a multithreaded program can use open file description locks to synchronize access to a file region by having each thread perform its own [open\(2\)](#) on the file and applying locks via the resulting file descriptor.

As with traditional advisory locks, the third argument to [fcntl\(\)](#), *lock*, is a pointer to an *flock* structure. By contrast with traditional record locks, the *l\_pid* field of that structure must be set to zero when using the commands described below.

The commands for working with open file description locks are analogous to those used with traditional locks:

#### **F\_OFD\_SETLK** (*struct flock* \*)

Acquire an open file description lock (when *l\_type* is **F\_RDLCK** or **F\_WRLCK**) or release an open file description lock (when *l\_type* is **F\_UNLCK**) on the bytes specified by the *l\_whence*, *l\_start*, and *l\_len* fields of *lock*. If a conflicting lock is held by another process, this call returns -1 and sets *errno* to **EAGAIN**.

#### **F\_OFD\_SETLKW** (*struct flock* \*)

As for **F\_OFD\_SETLK**, but if a conflicting lock is held on the file, then wait for that lock to be released. If a signal is caught while waiting, then the call is interrupted and (after the signal handler has returned) returns immediately (with return value -1 and *errno* set to **EINTR**; see [signal\(7\)](#)).

#### **F\_OFD\_GETLK** (*struct flock* \*)

On input to this call, *lock* describes an open file description lock we would like to place on the file. If the lock could be placed, [fcntl\(\)](#) does not actually place it, but returns **F\_UNLCK** in the *l\_type* field of *lock* and leaves the other fields of the structure unchanged. If one or more incompatible locks would prevent this lock being placed, then details about one of these locks are returned via *lock*, as described above for **F\_GETLK**.

In the current implementation, no deadlock detection is performed for open file description locks. (This contrasts with process-associated record locks, for which the kernel does perform deadlock detection.)

### **Mandatory locking**

*Warning:* the Linux implementation of mandatory locking is unreliable. See **BUGS** below.

By default, both traditional (process-associated) and open file description record locks are advisory. Advisory locks are not enforced and are useful only between cooperating processes.

Both lock types can also be mandatory. Mandatory locks are enforced for all processes. If a process tries to perform an incompatible access (e.g., [read\(2\)](#) or [write\(2\)](#)) on a file region that has an incompatible mandatory lock, then the result depends upon whether the **O\_NONBLOCK** flag is enabled for its open file description. If the **O\_NONBLOCK** flag is not enabled, then the

system call is blocked until the lock is removed or converted to a mode that is compatible with the access. If the **O\_NONBLOCK** flag is enabled, then the system call fails with the error **EAGAIN**.

To make use of mandatory locks, mandatory locking must be enabled both on the filesystem that contains the file to be locked, and on the file itself. Mandatory locking is enabled on a filesystem using the **-o** mand option to [mount\(8\)](#), or the **MS\_MANDLOCK** flag for [mount\(2\)](#). Mandatory locking is enabled on a file by disabling group execute permission on the file and enabling the set-group-ID permission bit (see [chmod\(1\)](#) and [chmod\(2\)](#)).

Mandatory locking is not specified by POSIX. Some other systems also support mandatory locking, although the details of how to enable it vary across systems.

### Managing signals

**F\_GETOWN**, **F\_SETOWN**, **F\_GETOWN\_EX**, **F\_SETOWN\_EX**, **F\_GETSIG** and **F\_SETSIG** are used to manage I/O availability signals:

#### **F\_GETOWN** (*void*)

Return (as the function result) the process ID or process group currently receiving **SIGIO** and **SIGURG** signals for events on file descriptor *fd*. Process IDs are returned as positive values; process group IDs are returned as negative values (but see **BUGS** below). *arg* is ignored.

#### **F\_SETOWN** (*int*)

Set the process ID or process group ID that will receive **SIGIO** and **SIGURG** signals for events on file descriptor *fd* to the ID given in *arg*. A process ID is specified as a positive value; a process group ID is specified as a negative value. Most commonly, the calling process specifies itself as the owner (that is, *arg* is specified as [getpid\(2\)](#)).

If you set the **O\_ASYNC** status flag on a file descriptor by using the **F\_SETFL** command of [fcntl\(\)](#), a **SIGIO** signal is sent whenever input or output becomes possible on that file descriptor. **F\_SETSIG** can be used to obtain delivery of a signal other than **SIGIO**. If this permission check fails, then the signal is silently discarded.

Sending a signal to the owner process (group) specified by **F\_SETOWN** is subject to the same permissions checks as are described for [kill\(2\)](#), where the sending process is the one that employs **F\_SETOWN** (but see **BUGS** below).

If the file descriptor *fd* refers to a socket, **F\_SETOWN** also selects the recipient of **SIGURG** signals that are delivered when out-of-band data arrives on that socket. (**SIGURG** is sent in any situation where [select\(2\)](#) would report the socket as having an exceptional condition.)

The following was true in 2.6.x kernels up to and including kernel 2.6.11:

If a nonzero value is given to **F\_SETSIG** in a multithreaded process running with a threading library that supports thread groups (e.g., **NPTL**), then a positive value given to **F\_SETOWN** has a different meaning: instead of being a process ID identifying a whole process, it is a thread ID identifying a specific thread within a process. Consequently, it may be necessary to pass **F\_SETOWN** the result of [gettid\(2\)](#) instead of [getpid\(2\)](#) to get sensible results when **F\_SETSIG** is used. (In current Linux threading implementations, a main thread's thread ID is the same as its process ID. This means that a single-threaded program can equally use [gettid\(2\)](#) or [getpid\(2\)](#) in this scenario.) Note, however, that the statements in this paragraph do not apply to the **SIGURG** signal generated for out-of-band data on a socket: this signal is always sent to either a process or a process group, depending on the value given to **F\_SETOWN**.

The above behavior was accidentally dropped in Linux 2.6.12, and won't be restored. From Linux 2.6.32 onward, use **F\_SETOWN\_EX** to target **SIGIO** and **SIGURG** signals at a particular thread.

**F\_GETOWN\_EX** (struct *f\_owner\_ex* \*) (since Linux 2.6.32)

Return the current file descriptor owner settings as defined by a previous **F\_SETOWN\_EX** operation. The information is returned in the structure pointed to by *arg*, which has the following form:

```
struct f_owner_ex {
    int type;
    pid_t pid;
};
```

The *type* field will have one of the values **F\_OWNER\_TID**, **F\_OWNER\_PID**, or **F\_OWNER\_PGRP**. The *pid* field is a positive integer representing a thread ID, process ID, or process group ID. See **F\_SETOWN\_EX** for more details.

**F\_SETOWN\_EX** (struct *f\_owner\_ex* \*) (since Linux 2.6.32)

This operation performs a similar task to **F\_SETOWN**. It allows the caller to direct I/O availability signals to a specific thread, process, or process group. The caller specifies the target of signals via *arg*, which is a pointer to a *f\_owner\_ex* structure. The *type* field has one of the following values, which define how *pid* is interpreted:

**F\_OWNER\_TID**

Send the signal to the thread whose thread ID (the value returned by a call to [clone\(2\)](#) or [gettid\(2\)](#)) is specified in *pid*.

**F\_OWNER\_PID**

Send the signal to the process whose ID is specified in *pid*.

**F\_OWNER\_PGRP**

Send the signal to the process group whose ID is specified in *pid*. (Note that, unlike with **F\_SETOWN**, a process group ID is specified as a positive value here.)

**F\_GETSIG** (*void*)

Return (as the function result) the signal sent when input or output becomes possible. A value of zero means **SIGIO** is sent. Any other value (including **SIGIO**) is the signal sent instead, and in this case additional info is available to the signal handler if installed with **SA\_SIGINFO**. *arg* is ignored.

**F\_SETSIG** (*int*)

Set the signal sent when input or output becomes possible to the value given in *arg*. A value of zero means to send the default **SIGIO** signal. Any other value (including **SIGIO**) is the signal to send instead, and in this case additional info is available to the signal handler if installed with **SA\_SIGINFO**.

By using **F\_SETSIG** with a nonzero value, and setting **SA\_SIGINFO** for the signal handler (see [sigaction\(2\)](#)), extra information about I/O events is passed to the handler in a *siginfo\_t* structure. If the *si\_code* field indicates the source is **SI\_SIGIO**, the *si\_fd* field gives the file descriptor associated with the event. Otherwise, there is no indication which file descriptors are pending, and you should use the usual mechanisms ([select\(2\)](#), [poll\(2\)](#), [read\(2\)](#) with **O\_NONBLOCK** set etc.) to determine which file descriptors are available for I/O.

By selecting a real time signal (value  $\geq$  **SIGRTMIN**), multiple I/O events may be queued using the same signal numbers. (Queuing is dependent on available memory). Extra information is available if **SA\_SIGINFO** is set for the signal handler, as above.

Note that Linux imposes a limit on the number of real-time signals that may be queued

to a process (see [getrlimit\(2\)](#) and [signal\(7\)](#)) and if this limit is reached, then the kernel reverts to delivering **SIGIO**, and this signal is delivered to the entire process rather than to a specific thread.

Using these mechanisms, a program can implement fully asynchronous I/O without using [select\(2\)](#) or [poll\(2\)](#) most of the time.

The use of **O\_ASYNC** is specific to BSD and Linux. The only use of **F\_GETOWN** and **F\_SETOWN** specified in POSIX.1 is in conjunction with the use of the **SIGURG** signal on sockets. (POSIX does not specify the **SIGIO** signal.) **F\_GETOWN\_EX**, **F\_SETOWN\_EX**, **F\_GETSIG**, and **F\_SETSIG** are Linux-specific. POSIX has asynchronous I/O and the *aiocb* structure to achieve similar things; these are also available in Linux as part of the GNU C Library (Glibc).

### Leases

**F\_SETLEASE** and **F\_GETLEASE** (Linux 2.4 onward) are used (respectively) to establish a new lease, and retrieve the current lease, on the open file description referred to by the file descriptor *fd*. A file lease provides a mechanism whereby the process holding the lease (the lease holder) is notified (via delivery of a signal) when a process (the lease breaker) tries to [open\(2\)](#) or [truncate\(2\)](#) the file referred to by that file descriptor.

#### **F\_SETLEASE** (*int*)

Set or remove a file lease according to which of the following values is specified in the integer *arg*:

##### **F\_RDLCK**

Take out a read lease. This will cause the calling process to be notified when the file is opened for writing or is truncated. A read lease can be placed only on a file descriptor that is opened read-only.

##### **F\_WRLCK**

Take out a write lease. This will cause the caller to be notified when the file is opened for reading or writing or is truncated. A write lease may be placed on a file only if there are no other open file descriptors for the file.

##### **F\_UNLCK**

Remove our lease from the file.

Leases are associated with an open file description (see [open\(2\)](#)). This means that duplicate file descriptors (created by, for example, [fork\(2\)](#) or [dup\(2\)](#)) refer to the same lease, and this lease may be modified or released using any of these descriptors. Furthermore, the lease is released by either an explicit **F\_UNLCK** operation on any of these duplicate descriptors, or when all such descriptors have been closed.

Leases may be taken out only on regular files. An unprivileged process may take out a lease only on a file whose UID (owner) matches the filesystem UID of the process. A process with the **CAP\_LEASE** capability may take out leases on arbitrary files.

#### **F\_GETLEASE** (*void*)

Indicates what type of lease is associated with the file descriptor *fd* by returning either **F\_RDLCK**, **F\_WRLCK**, or **F\_UNLCK**, indicating, respectively, a read lease, a write lease, or no lease. *arg* is ignored.

When a process (the lease breaker) performs an [open\(2\)](#) or [truncate\(2\)](#) that conflicts with a lease established via **F\_SETLEASE**, the system call is blocked by the kernel and the kernel notifies the lease holder by sending it a signal (**SIGIO** by default). The lease holder should respond to receipt of this signal by doing whatever cleanup is required in preparation for the file to be accessed by another process (e.g., flushing cached buffers) and then either remove or downgrade its lease. A lease is removed by performing an **F\_SETLEASE** command specifying *arg* as **F\_UNLCK**. If the lease holder currently holds a write lease on the file, and the lease breaker is opening the file for reading, then it is sufficient for the lease holder to downgrade the lease to a



read lease. This is done by performing an **F\_SETLEASE** command specifying *arg* as **F\_RDLCK**.

If the lease holder fails to downgrade or remove the lease within the number of seconds specified in */proc/sys/fs/lease-break-time*, then the kernel forcibly removes or downgrades the lease holder's lease.

Once a lease break has been initiated, **F\_GETLEASE** returns the target lease type (either **F\_RDLCK** or **F\_UNLCK**, depending on what would be compatible with the lease breaker) until the lease holder voluntarily downgrades or removes the lease or the kernel forcibly does so after the lease break timer expires.

Once the lease has been voluntarily or forcibly removed or downgraded, and assuming the lease breaker has not unblocked its system call, the kernel permits the lease breaker's system call to proceed.

If the lease breaker's blocked **open(2)** or **truncate(2)** is interrupted by a signal handler, then the system call fails with the error **EINTR**, but the other steps still occur as described above. If the lease breaker is killed by a signal while blocked in **open(2)** or **truncate(2)**, then the other steps still occur as described above. If the lease breaker specifies the **O\_NONBLOCK** flag when calling **open(2)**, then the call immediately fails with the error **EWOULDBLOCK**, but the other steps still occur as described above.

The default signal used to notify the lease holder is **SIGIO**, but this can be changed using the **F\_SETSIG** command to **fcntl()**. If a **F\_SETSIG** command is performed (even one specifying **SIGIO**), and the signal handler is established using **SA\_SIGINFO**, then the handler will receive a *siginfo\_t* structure as its second argument, and the *si\_fd* field of this argument will hold the descriptor of the leased file that has been accessed by another process. (This is useful if the caller holds leases against multiple files).

### File and directory change notification (**dnotify**)

#### **F\_NOTIFY** (*int*)

(Linux 2.4 onward) Provide notification when the directory referred to by *fd* or any of the files that it contains is changed. The events to be notified are specified in *arg*, which is a bit mask specified by ORing together zero or more of the following bits:

#### **DN\_ACCESS**

A file was accessed (**read(2)**, **pread(2)**, **readv(2)**, and similar)

#### **DN\_MODIFY**

A file was modified (**write(2)**, **pwrite(2)**, **writew(2)**, **truncate(2)**, **ftruncate(2)**, and similar).

#### **DN\_CREATE**

A file was created (**open(2)**, **creat(2)**, **mknod(2)**, **mkdir(2)**, **link(2)**, **symlink(2)**, **rename(2)** into this directory).

#### **DN\_DELETE**

A file was unlinked (**unlink(2)**, **rename(2)** to another directory, **rmdir(2)**).

#### **DN\_RENAME**

A file was renamed within this directory (**rename(2)**).

#### **DN\_ATTRIB**

The attributes of a file were changed (**chown(2)**, **chmod(2)**, **utime(2)**, **utimensat(2)**, and similar).

(In order to obtain these definitions, the **\_GNU\_SOURCE** feature test macro must be defined before including *any* header files.)

Directory notifications are normally one-shot, and the application must reregister to receive further notifications. Alternatively, if **DN\_MULTISHOT** is included in *arg*, then notification will remain in effect until explicitly removed.

A series of **F\_NOTIFY** requests is cumulative, with the events in *arg* being added to the



set already monitored. To disable notification of all events, make an **F\_NOTIFY** call specifying *arg* as 0.

Notification occurs via delivery of a signal. The default signal is **SIGIO**, but this can be changed using the **F\_SETSIG** command to **fcntl()**. (Note that **SIGIO** is one of the nonqueuing standard signals; switching to the use of a real-time signal means that multiple notifications can be queued to the process.) In the latter case, the signal handler receives a *siginfo\_t* structure as its second argument (if the handler was established using **SA\_SIGINFO**) and the *si\_fd* field of this structure contains the file descriptor which generated the notification (useful when establishing notification on multiple directories).

Especially when using **DN\_MULTISHOT**, a real time signal should be used for notification, so that multiple notifications can be queued.

**NOTE:** New applications should use the *inotify* interface (available since kernel 2.6.13), which provides a much superior interface for obtaining notifications of filesystem events. See [inotify\(7\)](#).

### Changing the capacity of a pipe

**F\_SETPIPE\_SZ** (*int*; since Linux 2.6.35)

Change the capacity of the pipe referred to by *fd* to be at least *arg* bytes. An unprivileged process can adjust the pipe capacity to any value between the system page size and the limit defined in */proc/sys/fs/pipe-max-size* (see [proc\(5\)](#)). Attempts to set the pipe capacity below the page size are silently rounded up to the page size. Attempts by an unprivileged process to set the pipe capacity above the limit in */proc/sys/fs/pipe-max-size* yield the error **EPERM**; a privileged process (**CAP\_SYS\_RESOURCE**) can override the limit. When allocating the buffer for the pipe, the kernel may use a capacity larger than *arg*, if that is convenient for the implementation. The actual capacity that is set is returned as the function result. Attempting to set the pipe capacity smaller than the amount of buffer space currently used to store data produces the error **EBUSY**.

**F\_GETPIPE\_SZ** (*void*; since Linux 2.6.35)

Return (as the function result) the capacity of the pipe referred to by *fd*.

### RETURN VALUE

For a successful call, the return value depends on the operation:

**F\_DUPFD** The new descriptor.

**F\_GETFD** Value of file descriptor flags.

**F\_GETFL** Value of file status flags.

**F\_GETLEASE**

Type of lease held on file descriptor.

**F\_GETOWN**

Value of descriptor owner.

**F\_GETSIG** Value of signal sent when read or write becomes possible, or zero for traditional **SIGIO** behavior.

**F\_GETPIPE\_SZ, F\_SETPIPE\_SZ**

The pipe capacity.

All other commands

Zero.

On error, -1 is returned, and *errno* is set appropriately.

### ERRORS

**EACCES** or **EAGAIN**

Operation is prohibited by locks held by other processes.

**EAGAIN**

The operation is prohibited because the file has been memory-mapped by another process.

**EBADF**

*fd* is not an open file descriptor, or the command was **F\_SETLK** or **F\_SETLKW** and the file descriptor open mode doesn't match with the type of lock requested.

**EDEADLK**

It was detected that the specified **F\_SETLKW** command would cause a deadlock.

**EFAULT**

*lock* is outside your accessible address space.

**EINTR**

For **F\_SETLKW**, the command was interrupted by a signal; see [signal\(7\)](#). For **F\_GETLK** and **F\_SETLK**, the command was interrupted by a signal before the lock was checked or acquired. Most likely when locking a remote file (e.g., locking over NFS), but can sometimes happen locally.

**EINVAL**

The value specified in *cmd* is not recognized by this kernel.

**EINVAL**

For **F\_DUPFD**, *arg* is negative or is greater than the maximum allowable value. For **F\_SETSIG**, *arg* is not an allowable signal number.

**EINVAL**

*cmd* is **F\_OFD\_SETLK**, **F\_OFD\_SETLKW**, or **F\_OFD\_GETLK**, and *l\_pid* was not specified as zero.

**EMFILE**

For **F\_DUPFD**, the process already has the maximum number of file descriptors open.

**ENOLCK**

Too many segment locks open, lock table is full, or a remote locking protocol failed (e.g., locking over NFS).

**ENOTDIR**

**F\_NOTIFY** was specified in *cmd*, but *fd* does not refer to a directory.

**EPERM**

Attempted to clear the **O\_APPEND** flag on a file that has the append-only attribute set.

**CONFORMING TO**

SVr4, 4.3BSD, POSIX.1-2001. Only the operations **F\_DUPFD**, **F\_GETFD**, **F\_SETFD**, **F\_GETFL**, **F\_SETFL**, **F\_GETLK**, **F\_SETLK**, and **F\_SETLKW** are specified in POSIX.1-2001.

**F\_GETOWN** and **F\_SETOWN** are specified in POSIX.1-2001. (To get their definitions, define either **\_BSD\_SOURCE**, or **\_XOPEN\_SOURCE** with the value 500 or greater, or **\_POSIX\_C\_SOURCE** with the value 200809L or greater.)

**F\_DUPFD\_CLOEXEC** is specified in POSIX.1-2008. (To get this definition, define **\_POSIX\_C\_SOURCE** with the value 200809L or greater, or **\_XOPEN\_SOURCE** with the value 700 or greater.)

**F\_GETOWN\_EX**, **F\_SETOWN\_EX**, **F\_SETPIPE\_SZ**, **F\_GETPIPE\_SZ**, **F\_GETSIG**, **F\_SETSIG**, **F\_NOTIFY**, **F\_GETLEASE**, and **F\_SETLEASE** are Linux-specific. (Define the **\_GNU\_SOURCE** macro to obtain these definitions.)

**F\_OFD\_SETLK**, **F\_OFD\_SETLKW**, and **F\_OFD\_GETLK** are Linux-specific (and one must define **\_GNU\_SOURCE** to obtain their definitions), but work is being done to have them

included in the next version of POSIX.1.

## NOTES

The errors returned by `dup2(2)` are different from those returned by `F_DUPFD`.

### File locking

The original Linux `fcntl()` system call was not designed to handle large file offsets (in the *flock* structure). Consequently, an `fcntl64()` system call was added in Linux 2.4. The newer system call employs a different structure for file locking, *flock64*, and corresponding commands, `F_GETLK64`, `F_SETLK64`, and `F_SETLKW64`. However, these details can be ignored by applications using glibc, whose `fcntl()` wrapper function transparently employs the more recent system call where it is available.

The errors returned by `dup2(2)` are different from those returned by `F_DUPFD`.

### Record locks

Since kernel 2.0, there is no interaction between the types of lock placed by `flock(2)` and `fcntl()`.

Several systems have more fields in *struct flock* such as, for example, *l\_sysid*. Clearly, *l\_pid* alone is not going to be very useful if the process holding the lock may live on a different machine.

The original Linux `fcntl()` system call was not designed to handle large file offsets (in the *flock* structure). Consequently, an `fcntl64()` system call was added in Linux 2.4. The newer system call employs a different structure for file locking, *flock64*, and corresponding commands, `F_GETLK64`, `F_SETLK64`, and `F_SETLKW64`. However, these details can be ignored by applications using glibc, whose `fcntl()` wrapper function transparently employs the more recent system call where it is available.

### Record locking and NFS

Before Linux 3.12, if an NFSv4 client loses contact with the server for a period of time (defined as more than 90 seconds with no communication), it might lose and regain a lock without ever being aware of the fact. (The period of time after which contact is assumed lost is known as the NFSv4 leasetime. On a Linux NFS server, this can be determined by looking at `/proc/fs/nfsd/nfsv4leasetime`, which expresses the period in seconds. The default value for this file is 90.) This scenario potentially risks data corruption, since another process might acquire a lock in the intervening period and perform file I/O.

Since Linux 3.12, if an NFSv4 client loses contact with the server, any I/O to the file by a process which thinks it holds a lock will fail until that process closes and reopens the file. A kernel parameter, `nfs.recover_lost_locks`, can be set to 1 to obtain the pre-3.12 behavior, whereby the client will attempt to recover lost locks when contact is reestablished with the server. Because of the attendant risk of data corruption, this parameter defaults to 0 (disabled).

## BUGS

### F\_SETFL

It is not possible to use `F_SETFL` to change the state of the `O_DSYNC` and `O_SYNC` flags. Attempts to change the state of these flags are silently ignored.

### F\_GETOWN

A limitation of the Linux system call conventions on some architectures (notably i386) means that if a (negative) process group ID to be returned by `F_GETOWN` falls in the range -1 to -4095, then the return value is wrongly interpreted by glibc as an error in the system call; that is, the return value of `fcntl()` will be -1, and `errno` will contain the (positive) process group ID. The Linux-specific `F_GETOWN_EX` operation avoids this problem. Since glibc version 2.11, glibc makes the kernel `F_GETOWN` problem invisible by implementing `F_GETOWN` using `F_GETOWN_EX`.

### F\_SETOWN

In Linux 2.4 and earlier, there is bug that can occur when an unprivileged process uses `F_SETOWN` to specify the owner of a socket file descriptor as a process (group) other than the caller. In this case, `fcntl()` can return -1 with `errno` set to `EPERM`, even when the owner

process (group) is one that the caller has permission to send signals to. Despite this error return, the file descriptor owner is set, and signals will be sent to the owner.

### Deadlock detection

The deadlock-detection algorithm employed by the kernel when dealing with **F\_SETLKW** requests can yield both false negatives (failures to detect deadlocks, leaving a set of deadlocked processes blocked indefinitely) and false positives (**EDEADLK** errors when there is no deadlock). For example, the kernel limits the lock depth of its dependency search to 10 steps, meaning that circular deadlock chains that exceed that size will not be detected. In addition, the kernel may falsely indicate a deadlock when two or more processes created using the [clone\(2\)](#) **CLONE\_FILES** flag place locks that appear (to the kernel) to conflict.

### Mandatory locking

The Linux implementation of mandatory locking is subject to race conditions which render it unreliable: a [write\(2\)](#) call that overlaps with a lock may modify data after the mandatory lock is acquired; a [read\(2\)](#) call that overlaps with a lock may detect changes to data that were made only after a write lock was acquired. Similar races exist between mandatory locks and [mmap\(2\)](#). It is therefore inadvisable to rely on mandatory locking.

### SEE ALSO

[dup2\(2\)](#), [flock\(2\)](#), [open\(2\)](#), [socket\(2\)](#), [lockf\(3\)](#), [capabilities\(7\)](#), [feature\\_test\\_macros\(7\)](#)

*locks.txt*, *mandatory-locking.txt*, and *dnotify.txt* in the Linux kernel source directory *Documentation/filesystems/* (on older kernels, these files are directly under the *Documentation/* directory, and *mandatory-locking.txt* is called *mandatory.txt*)

### COLOPHON

This page is part of release 3.74 of the Linux *man-pages* project. A description of the project, information about reporting bugs, and the latest version of this page, can be found at <http://www.kernel.org/doc/man-pages/>.