

**NAME**

ioctl\_userfaultfd – create a file descriptor for handling page faults in user space

**LIBRARY**

Standard C library (*libc*, *-lc*)

**SYNOPSIS**

```
#include <linux/userfaultfd.h> /* Definition of UFFD* constants */
#include <sys/ioctl.h>

int ioctl(int fd, int cmd, ...);
```

**DESCRIPTION**

Various **ioctl(2)** operations can be performed on a userfaultfd object (created by a call to **userfaultfd(2)**) using calls of the form:

```
ioctl(fd, cmd, argp);
```

In the above, *fd* is a file descriptor referring to a userfaultfd object, *cmd* is one of the commands listed below, and *argp* is a pointer to a data structure that is specific to *cmd*.

The various **ioctl(2)** operations are described below. The **UFFDIO\_API**, **UFFDIO\_REGISTER**, and **UFFDIO\_UNREGISTER** operations are used to *configure* userfaultfd behavior. These operations allow the caller to choose what features will be enabled and what kinds of events will be delivered to the application. The remaining operations are *ange* operations. These operations enable the calling application to resolve page-fault events.

**UFFDIO\_API**

(Since Linux 4.3.) Enable operation of the userfaultfd and perform API handshake.

The *argp* argument is a pointer to a *uffdio\_api* structure, defined as:

```
struct uffdio_api {
    __u64 api;          /* Requested API version (input) */
    __u64 features;     /* Requested features (input/output) */
    __u64 ioctls;      /* Available ioctl() operations (output) */
};
```

The *api* field denotes the API version requested by the application.

The kernel verifies that it can support the requested API version, and sets the *features* and *ioctls* fields to bit masks representing all the available features and the generic **ioctl(2)** operations available.

Before Linux 4.11, the *features* field must be initialized to zero before the call to **UFFDIO\_API**, and zero (i.e., no feature bits) is placed in the *features* field by the kernel upon return from **ioctl(2)**.

Starting from Linux 4.11, the *features* field can be used to ask whether particular features are supported and explicitly enable userfaultfd features that are disabled by default. The kernel always reports all the available features in the *features* field.

To enable userfaultfd features the application should set a bit corresponding to each feature it wants to enable in the *features* field. If the kernel supports all the requested features it will enable them. Otherwise it will zero out the returned *uffdio\_api* structure and return **EINVAL**.

The following feature bits may be set:

**UFFD\_FEATURE\_EVENT\_FORK** (since Linux 4.11)

When this feature is enabled, the userfaultfd objects associated with a parent process are duplicated into the child process during **fork(2)** and a **UFFD\_EVENT\_FORK** event is delivered to the userfaultfd monitor

**UFFD\_FEATURE\_EVENT\_REMAP** (since Linux 4.11)

If this feature is enabled, when the faulting process invokes **mremap(2)**, the userfaultfd monitor will receive an event of type **UFFD\_EVENT\_REMAP**.

**UFFD\_FEATURE\_EVENT\_REMOVE** (since Linux 4.11)

If this feature is enabled, when the faulting process calls **madvise(2)** with the **MADV\_DONT-NEED** or **MADV\_REMOVE** advice value to free a virtual memory area the userfaultfd monitor will receive an event of type **UFFD\_EVENT\_REMOVE**.

**UFFD\_FEATURE\_EVENT\_UNMAP** (since Linux 4.11)

If this feature is enabled, when the faulting process unmaps virtual memory either explicitly with **munmap(2)**, or implicitly during either **mmap(2)** or **mremap(2)**, the userfaultfd monitor will receive an event of type **UFFD\_EVENT\_UNMAP**.

**UFFD\_FEATURE\_MISSING\_HUGETLBFS** (since Linux 4.11)

If this feature bit is set, the kernel supports registering userfaultfd ranges on hugetlbfs virtual memory areas

**UFFD\_FEATURE\_MISSING\_SHMEM** (since Linux 4.11)

If this feature bit is set, the kernel supports registering userfaultfd ranges on shared memory areas. This includes all kernel shared memory APIs: System V shared memory, **tmpfs(5)**, shared mappings of */dev/zero*, **mmap(2)** with the **MAP\_SHARED** flag set, **memfd\_create(2)**, and so on.

**UFFD\_FEATURE\_SIGBUS** (since Linux 4.14)

If this feature bit is set, no page-fault events (**UFFD\_EVENT\_PAGEFAULT**) will be delivered. Instead, a **SIGBUS** signal will be sent to the faulting process. Applications using this feature will not require the use of a userfaultfd monitor for processing memory accesses to the regions registered with userfaultfd.

**UFFD\_FEATURE\_THREAD\_ID** (since Linux 4.14)

If this feature bit is set, *uffd\_msg.pagefault.feat.ptid* will be set to the faulted thread ID for each page-fault message.

**UFFD\_FEATURE\_MINOR\_HUGETLBFS** (since Linux 5.13)

If this feature bit is set, the kernel supports registering userfaultfd ranges in minor mode on hugetlbfs-backed memory areas.

**UFFD\_FEATURE\_MINOR\_SHMEM** (since Linux 5.14)

If this feature bit is set, the kernel supports registering userfaultfd ranges in minor mode on shmem-backed memory areas.

The returned *ioctls* field can contain the following bits:

**1 << \_UFFDIO\_API**

The **UFFDIO\_API** operation is supported.

**1 << \_UFFDIO\_REGISTER**

The **UFFDIO\_REGISTER** operation is supported.

**1 << \_UFFDIO\_UNREGISTER**

The **UFFDIO\_UNREGISTER** operation is supported.

This **ioctl(2)** operation returns 0 on success. On error, -1 is returned and *errno* is set to indicate the error. Possible errors include:

**EFAULT**

*argp* refers to an address that is outside the calling process's accessible address space.

**EINVAL**

The userfaultfd has already been enabled by a previous **UFFDIO\_API** operation.

**EINVAL**

The API version requested in the *api* field is not supported by this kernel, or the *features* field passed to the kernel includes feature bits that are not supported by the current kernel version.

**UFFDIO\_REGISTER**

(Since Linux 4.3.) Register a memory address range with the userfaultfd object. The pages in the range must be "compatible". Please refer to the list of register modes below for the compatible memory backends

for each mode.

The *argp* argument is a pointer to a *uffdio\_register* structure, defined as:

```
struct uffdio_range {
    __u64 start;    /* Start of range */
    __u64 len;      /* Length of range (bytes) */
};

struct uffdio_register {
    struct uffdio_range range;
    __u64 mode;     /* Desired mode of operation (input) */
    __u64 ioctls;   /* Available ioctl() operations (output) */
};
```

The *range* field defines a memory range starting at *start* and continuing for *len* bytes that should be handled by the userfaultfd.

The *mode* field defines the mode of operation desired for this memory region. The following values may be bitwise ORed to set the userfaultfd mode for the specified range:

#### **UFFDIO\_REGISTER\_MODE\_MISSING**

Track page faults on missing pages. Since Linux 4.3, only private anonymous ranges are compatible. Since Linux 4.11, hugetlbfs and shared memory ranges are also compatible.

#### **UFFDIO\_REGISTER\_MODE\_WP**

Track page faults on write-protected pages. Since Linux 5.7, only private anonymous ranges are compatible.

#### **UFFDIO\_REGISTER\_MODE\_MINOR**

Track minor page faults. Since Linux 5.13, only hugetlbfs ranges are compatible. Since Linux 5.14, compatibility with shmем ranges was added.

If the operation is successful, the kernel modifies the *ioctls* bit-mask field to indicate which **ioctl(2)** operations are available for the specified range. This returned bit mask can contain the following bits:

#### **1 << \_UFFDIO\_COPY**

The **UFFDIO\_COPY** operation is supported.

#### **1 << \_UFFDIO\_WAKE**

The **UFFDIO\_WAKE** operation is supported.

#### **1 << \_UFFDIO\_WRITEPROTECT**

The **UFFDIO\_WRITEPROTECT**

#### **1 << \_UFFDIO\_ZEROPAGE**

The **UFFDIO\_ZEROPAGE** operation is supported.

#### **1 << \_UFFDIO\_CONTINUE**

The **UFFDIO\_CONTINUE** operation is supported.

This **ioctl(2)** operation returns 0 on success. On error, -1 is returned and *errno* is set to indicate the error. Possible errors include:

#### **EBUSY**

A mapping in the specified range is registered with another userfaultfd object.

#### **EFAULT**

*argp* refers to an address that is outside the calling process's accessible address space.

#### **EINVAL**

An invalid or unsupported bit was specified in the *mode* field; or the *mode* field was zero.

**EINVAL**

There is no mapping in the specified address range.

**EINVAL**

*range.start* or *range.len* is not a multiple of the system page size; or, *range.len* is zero; or these fields are otherwise invalid.

**EINVAL**

There as an incompatible mapping in the specified address range.

**UFFDIO\_UNREGISTER**

(Since Linux 4.3.) Unregister a memory address range from userfaultfd. The pages in the range must be "compatible" (see the description of **UFFDIO\_REGISTER**.)

The address range to unregister is specified in the *uffdio\_range* structure pointed to by *argp*.

This **ioctl**(2) operation returns 0 on success. On error,  $-1$  is returned and *errno* is set to indicate the error. Possible errors include:

**EINVAL**

Either the *start* or the *len* field of the *uffdio\_range* structure was not a multiple of the system page size; or the *len* field was zero; or these fields were otherwise invalid.

**EINVAL**

There as an incompatible mapping in the specified address range.

**EINVAL**

There was no mapping in the specified address range.

**UFFDIO\_COPY**

(Since Linux 4.3.) Atomically copy a continuous memory chunk into the userfault registered range and optionally wake up the blocked thread. The source and destination addresses and the number of bytes to copy are specified by the *src*, *dst*, and *len* fields of the *uffdio\_copy* structure pointed to by *argp*:

```
struct uffdio_copy {
    __u64 dst;      /* Destination of copy */
    __u64 src;      /* Source of copy */
    __u64 len;      /* Number of bytes to copy */
    __u64 mode;     /* Flags controlling behavior of copy */
    __s64 copy;     /* Number of bytes copied, or negated error */
};
```

The following value may be bitwise ORed in *mode* to change the behavior of the **UFFDIO\_COPY** operation:

**UFFDIO\_COPY\_MODE\_DONTWAKE**

Do not wake up the thread that waits for page-fault resolution

**UFFDIO\_COPY\_MODE\_WP**

Copy the page with read-only permission. This allows the user to trap the next write to the page, which will block and generate another write-protect userfault message. This is used only when both **UFFDIO\_REGISTER\_MODE\_MISSING** and **UFFDIO\_REGISTER\_MODE\_WP** modes are enabled for the registered range.

The *copy* field is used by the kernel to return the number of bytes that was actually copied, or an error (a negated *errno*-style value). If the value returned in *copy* doesn't match the value that was specified in *len*, the operation fails with the error **EAGAIN**. The *copy* field is output-only; it is not read by the **UFFDIO\_COPY** operation.

This **ioctl**(2) operation returns 0 on success. In this case, the entire area was copied. On error,  $-1$  is returned and *errno* is set to indicate the error. Possible errors include:

**EAGAIN**

The number of bytes copied (i.e., the value returned in the *copy* field) does not equal the value that was specified in the *len* field.

**EINVAL**

Either *dst* or *len* was not a multiple of the system page size, or the range specified by *src* and *len* or *dst* and *len* was invalid.

**EINVAL**

An invalid bit was specified in the *mode* field.

**ENOENT** (since Linux 4.11)

The faulting process has changed its virtual memory layout simultaneously with an outstanding **UFFDIO\_COPY** operation.

**ENOSPC** (from Linux 4.11 until Linux 4.13)

The faulting process has exited at the time of a **UFFDIO\_COPY** operation.

**ESRCH** (since Linux 4.13)

The faulting process has exited at the time of a **UFFDIO\_COPY** operation.

**UFFDIO\_ZEROPAGE**

(Since Linux 4.3.) Zero out a memory range registered with userfaultfd.

The requested range is specified by the *range* field of the *uffdio\_zeropage* structure pointed to by *argp*:

```
struct uffdio_zeropage {
    struct uffdio_range range;
    __u64 mode;          /* Flags controlling behavior of copy */
    __s64 zeropage;      /* Number of bytes zeroed, or negated error */
};
```

The following value may be bitwise ORed in *mode* to change the behavior of the **UFFDIO\_ZEROPAGE** operation:

**UFFDIO\_ZEROPAGE\_MODE\_DONTWAKE**

Do not wake up the thread that waits for page-fault resolution.

The *zeropage* field is used by the kernel to return the number of bytes that was actually zeroed, or an error in the same manner as **UFFDIO\_COPY**. If the value returned in the *zeropage* field doesn't match the value that was specified in *range.len*, the operation fails with the error **EAGAIN**. The *zeropage* field is output-only; it is not read by the **UFFDIO\_ZEROPAGE** operation.

This **ioctl**(2) operation returns 0 on success. In this case, the entire area was zeroed. On error, -1 is returned and *errno* is set to indicate the error. Possible errors include:

**EAGAIN**

The number of bytes zeroed (i.e., the value returned in the *zeropage* field) does not equal the value that was specified in the *range.len* field.

**EINVAL**

Either *range.start* or *range.len* was not a multiple of the system page size; or *range.len* was zero; or the range specified was invalid.

**EINVAL**

An invalid bit was specified in the *mode* field.

**ESRCH** (since Linux 4.13)

The faulting process has exited at the time of a **UFFDIO\_ZEROPAGE** operation.

**UFFDIO\_WAKE**

(Since Linux 4.3.) Wake up the thread waiting for page-fault resolution on a specified memory address range.

The **UFFDIO\_WAKE** operation is used in conjunction with **UFFDIO\_COPY** and

**UFFDIO\_ZEROPAGE** operations that have the **UFFDIO\_COPY\_MODE\_DONTWAKE** or **UFFDIO\_ZEROPAGE\_MODE\_DONTWAKE** bit set in the *mode* field. The userfault monitor can perform several **UFFDIO\_COPY** and **UFFDIO\_ZEROPAGE** operations in a batch and then explicitly wake up the faulting thread using **UFFDIO\_WAKE**.

The *argp* argument is a pointer to a *uffdio\_range* structure (shown above) that specifies the address range.

This **ioctl(2)** operation returns 0 on success. On error, -1 is returned and *errno* is set to indicate the error. Possible errors include:

#### **EINVAL**

The *start* or the *len* field of the *uffdio\_range* structure was not a multiple of the system page size; or *len* was zero; or the specified range was otherwise invalid.

#### **UFFDIO\_WRITEPROTECT (Since Linux 5.7)**

Write-protect or write-unprotect a userfaultfd-registered memory range registered with mode **UFFDIO\_REGISTER\_MODE\_WP**.

The *argp* argument is a pointer to a *uffdio\_range* structure as shown below:

```
struct uffdio_writeprotect {
    struct uffdio_range range; /* Range to change write permission */
    __u64 mode;                /* Mode to change write permission */
};
```

There are two mode bits that are supported in this structure:

#### **UFFDIO\_WRITEPROTECT\_MODE\_WP**

When this mode bit is set, the **ioctl** will be a write-protect operation upon the memory range specified by *range*. Otherwise it will be a write-unprotect operation upon the specified range, which can be used to resolve a userfaultfd write-protect page fault.

#### **UFFDIO\_WRITEPROTECT\_MODE\_DONTWAKE**

When this mode bit is set, do not wake up any thread that waits for page-fault resolution after the operation. This can be specified only if **UFFDIO\_WRITEPROTECT\_MODE\_WP** is not specified.

This **ioctl(2)** operation returns 0 on success. On error, -1 is returned and *errno* is set to indicate the error. Possible errors include:

#### **EINVAL**

The *start* or the *len* field of the *uffdio\_range* structure was not a multiple of the system page size; or *len* was zero; or the specified range was otherwise invalid.

#### **EAGAIN**

The process was interrupted; retry this call.

#### **ENOENT**

The range specified in *range* is not valid. For example, the virtual address does not exist, or not registered with userfaultfd write-protect mode.

#### **EFAULT**

Encountered a generic fault during processing.

#### **UFFDIO\_CONTINUE**

(Since Linux 5.13.) Resolve a minor page fault by installing page table entries for existing pages in the page cache.

The *argp* argument is a pointer to a *uffdio\_continue* structure as shown below:

```
struct uffdio_continue {
    struct uffdio_range range;
                                /* Range to install PTEs for and continue */
    __u64 mode;                 /* Flags controlling the behavior of continue */
    __s64 mapped;               /* Number of bytes mapped, or negated error */
};
```

```
} ;
```

The following value may be bitwise ORed in *mode* to change the behavior of the **UFFDIO\_CONTINUE** operation:

#### **UFFDIO\_CONTINUE\_MODE\_DONTWAKE**

Do not wake up the thread that waits for page-fault resolution.

The *mapped* field is used by the kernel to return the number of bytes that were actually mapped, or an error in the same manner as **UFFDIO\_COPY**. If the value returned in the *mapped* field doesn't match the value that was specified in *range.len*, the operation fails with the error **EAGAIN**. The *mapped* field is output-only; it is not read by the **UFFDIO\_CONTINUE** operation.

This **ioctl(2)** operation returns 0 on success. In this case, the entire area was mapped. On error, -1 is returned and *errno* is set to indicate the error. Possible errors include:

#### **EAGAIN**

The number of bytes mapped (i.e., the value returned in the *mapped* field) does not equal the value that was specified in the *range.len* field.

#### **EINVAL**

Either *range.start* or *range.len* was not a multiple of the system page size; or *range.len* was zero; or the range specified was invalid.

#### **EINVAL**

An invalid bit was specified in the *mode* field.

#### **EEXIST**

One or more pages were already mapped in the given range.

#### **ENOENT**

The faulting process has changed its virtual memory layout simultaneously with an outstanding **UFFDIO\_CONTINUE** operation.

#### **ENOMEM**

Allocating memory needed to setup the page table mappings failed.

#### **EFAULT**

No existing page could be found in the page cache for the given range.

#### **ESRCH**

The faulting process has exited at the time of a **UFFDIO\_CONTINUE** operation.

### **RETURN VALUE**

See descriptions of the individual operations, above.

### **ERRORS**

See descriptions of the individual operations, above. In addition, the following general errors can occur for all of the operations described above:

#### **EFAULT**

*argp* does not point to a valid memory address.

#### **EINVAL**

(For all operations except **UFFDIO\_API**.) The userfaultfd object has not yet been enabled (via the **UFFDIO\_API** operation).

### **STANDARDS**

These **ioctl(2)** operations are Linux-specific.

### **BUGS**

In order to detect available userfault features and enable some subset of those features the userfaultfd file descriptor must be closed after the first **UFFDIO\_API** operation that queries features availability and reopened before the second **UFFDIO\_API** operation that actually enables the desired features.

**EXAMPLES**

See **userfaultfd(2)**.

**SEE ALSO**

**ioctl(2)**, **mmap(2)**, **userfaultfd(2)**

*Documentation/admin-guide/mm/userfaultfd.rst* in the Linux kernel source tree