

*fork()* or *vfork()* is followed by an *exec()*. This is illustrated by the final pair of data rows in Table 28-3, where each child performs an *exec()*, rather than immediately exiting. The program execed was the *true* command (*/bin/true*, chosen because it produces no output). In this case, we see that the relative differences between *fork()* and *vfork()* are much lower.

In fact, the data shown in Table 28-3 doesn't reveal the full cost of an *exec()*, because the child execs the same program in each loop of the test. As a result, the cost of disk I/O to read the program into memory is essentially eliminated, because the program will be read into the kernel buffer cache on the first *exec()*, and then remain there. If each loop of the test execed a different program (e.g., a differently named copy of the same program), then we would observe a greater cost for an *exec()*.

## 28.4 Effect of *exec()* and *fork()* on Process Attributes

A process has numerous attributes, some of which we have already described in earlier chapters, and others that we explore in later chapters. Regarding these attributes, two questions arise:

- What happens to these attributes when a process performs an *exec()*?
- Which attributes are inherited by a child when a *fork()* is performed?

Table 28-4 summarizes the answers to these questions. The *exec()* column indicates which attributes are preserved during an *exec()*. The *fork()* column indicates which attributes are inherited (or in some cases, shared) by a child after *fork()*. Other than the attributes indicated as being Linux-specific, all listed attributes appear in standard UNIX implementations, and their handling during *exec()* and *fork()* conforms to the requirements of SUSv3.

**Table 28-4:** Effect of *exec()* and *fork()* on process attributes

Process attribute	<i>exec()</i>	<i>fork()</i>	Interfaces affecting attribute; additional notes
<b>Process address space</b>			
Text segment	No	Shared	Child process shares text segment with parent.
Stack segment	No	Yes	Function entry/exit; <i>alloca()</i> , <i>longjmp()</i> , <i>siglongjmp()</i> .
Data and heap segments	No	Yes	<i>brk()</i> , <i>sbrk()</i> .
Environment variables	See notes	Yes	<i>putenv()</i> , <i>setenv()</i> ; direct modification of <i>environ</i> . Overwritten by <i>execle()</i> and <i>execve()</i> and preserved by remaining <i>exec()</i> calls.
Memory mappings	No	Yes; see notes	<i>mmap()</i> , <i>munmap()</i> . A mapping's <i>MAP_NORESERVE</i> flag is inherited across <i>fork()</i> . Mappings that have been marked with <i>madvice(MADV_DONTFORK)</i> are not inherited across <i>fork()</i> .
Memory locks	No	No	<i>mlock()</i> , <i>munlock()</i> .

**Table 28-4:** Effect of *exec()* and *fork()* on process attributes (continued)

Process attribute	<i>exec()</i>	<i>fork()</i>	Interfaces affecting attribute; additional notes
<b>Process identifiers and credentials</b>			
Process ID	Yes	No	
Parent process ID	Yes	No	
Process group ID	Yes	Yes	<i>setpgid()</i> .
Session ID	Yes	Yes	<i>setsid()</i> .
Real IDs	Yes	Yes	<i>setuid()</i> , <i>setgid()</i> , and related calls.
Effective and saved set IDs	See notes	Yes	<i>setuid()</i> , <i>setgid()</i> , and related calls. Chapter 9 explains how <i>exec()</i> affects these IDs.
Supplementary group IDs	Yes	Yes	<i>setgroups()</i> , <i>initgroups()</i> .
<b>Files, file I/O, and directories</b>			
Open file descriptors	See notes	Yes	<i>open()</i> , <i>close()</i> , <i>dup()</i> , <i>pipe()</i> , <i>socket()</i> , and so on. File descriptors are preserved across <i>exec()</i> unless marked close-on-exec. Descriptors in child and parent refer to same open file descriptions; see Section 5.4.
Close-on-exec flag	Yes (if off)	Yes	<i>fcntl(F_SETFD)</i> .
File offsets	Yes	Shared	<i>lseek()</i> , <i>read()</i> , <i>write()</i> , <i>readv()</i> , <i>writev()</i> . Child shares file offsets with parent.
Open file status flags	Yes	Shared	<i>open()</i> , <i>fctl(F_SETFL)</i> . Child shares open file status flags with parent.
Asynchronous I/O operations	See notes	No	<i>aio_read()</i> , <i>aio_write()</i> , and related calls. Outstanding operations are canceled during an <i>exec()</i> .
Directory streams	No	Yes; see notes	<i>opendir()</i> , <i>readdir()</i> . SUSv3 states that child gets a copy of parent's directory streams, but these copies may or may not share the directory stream position. On Linux, the directory stream position is not shared.
<b>File system</b>			
Current working directory	Yes	Yes	<i>chdir()</i> .
Root directory	Yes	Yes	<i>chroot()</i> .
File mode creation mask	Yes	Yes	<i>umask()</i> .
<b>Signals</b>			
Signal dispositions	See notes	Yes	<i>signal()</i> , <i>sigaction()</i> . During an <i>exec()</i> , signals with dispositions set to default or ignore are unchanged; caught signals revert to their default dispositions. See Section 27.5.
Signal mask	Yes	Yes	Signal delivery, <i>sigprocmask()</i> , <i>sigaction()</i> .
Pending signal set	Yes	No	Signal delivery; <i>raise()</i> , <i>kill()</i> , <i>sigqueue()</i> .
Alternate signal stack	No	Yes	<i>sigaltstack()</i> .

(continued)

**Table 28-4:** Effect of *exec()* and *fork()* on process attributes (continued)

Process attribute	<i>exec()</i>	<i>fork()</i>	Interfaces affecting attribute; additional notes
<b>Timers</b>			
Interval timers	Yes	No	<i>setitimer()</i> .
Timers set by <i>alarm()</i>	Yes	No	<i>alarm()</i> .
POSIX timers	No	No	<i>timer_create()</i> and related calls.
<b>POSIX threads</b>			
Threads	No	See notes	During <i>fork()</i> , only calling thread is replicated in child.
Thread cancelability state and type	No	Yes	After an <i>exec()</i> , the cancelability type and state are reset to PTHREAD_CANCEL_ENABLE and PTHREAD_CANCEL_DEFERRED, respectively
Mutexes and condition variables	No	Yes	See Section 33.3 for details of the treatment of mutexes and other thread resources during <i>fork()</i> .
<b>Priority and scheduling</b>			
Nice value	Yes	Yes	<i>nice()</i> , <i>setpriority()</i> .
Scheduling policy and priority	Yes	Yes	<i>sched_setscheduler()</i> , <i>sched_setparam()</i> .
<b>Resources and CPU time</b>			
Resource limits	Yes	Yes	<i>setrlimit()</i> .
Process and child CPU times	Yes	No	As returned by <i>times()</i> .
Resource usages	Yes	No	As returned by <i>getrusage()</i> .
<b>Interprocess communication</b>			
System V shared memory segments	No	Yes	<i>shmat()</i> , <i>shmdt()</i> .
POSIX shared memory	No	Yes	<i>shm_open()</i> and related calls.
POSIX message queues	No	Yes	<i>mq_open()</i> and related calls. Descriptors in child and parent refer to same open message queue descriptions. A child doesn't inherit its parent's message notification registrations.
POSIX named semaphores	No	Shared	<i>sem_open()</i> and related calls. Child shares references to same semaphores as parent.
POSIX unnamed semaphores	No	See notes	<i>sem_init()</i> and related calls. If semaphores are in a shared memory region, then child shares semaphores with parent; otherwise, child has its own copy of the semaphores.
System V semaphore adjustments	Yes	No	<i>semop()</i> . See Section 47.8.
File locks	Yes	See notes	<i>flock()</i> . Child inherits a reference to the same lock as parent.
Record locks	See notes	No	<i>fcntl(F_SETLK)</i> . Locks are preserved across <i>exec()</i> unless a file descriptor referring to the file is marked close-on-exec; see Section 55.3.5.

**Table 28-4:** Effect of *exec()* and *fork()* on process attributes (continued)

Process attribute	<i>exec()</i>	<i>fork()</i>	Interfaces affecting attribute; additional notes
<b>Miscellaneous</b>			
Locale settings	No	Yes	<i>setlocale()</i> . As part of C run-time initialization, the equivalent of <i>setlocale(LC_ALL, "C")</i> is executed after a new program is execed.
Floating-point environment	No	Yes	When a new program is execed, the state of the floating-point environment is reset to the default; see <i>fenv(3)</i> .
Controlling terminal	Yes	Yes	
Exit handlers	No	Yes	<i>atexit()</i> , <i>on_exit()</i> .
<b>Linux-specific</b>			
File-system IDs	See notes	Yes	<i>setsuid()</i> , <i>setsgid()</i> . These IDs are also changed any time the corresponding effective IDs are changed.
<i>timerfd</i> timers	Yes	See notes	<i>timerfd_create()</i> ; child inherits file descriptors referring to same timers as parent.
Capabilities	See notes	Yes	<i>capset()</i> . The handling of capabilities during an <i>exec()</i> is described in Section 39.5.
Capability bounding set	Yes	Yes	
Capabilities <i>securebits</i> flags	See notes	Yes	All <i>securebits</i> flags are preserved during an <i>exec()</i> except <i>SECBIT_KEEP_CAPS</i> , which is always cleared.
CPU affinity	Yes	Yes	<i>sched_setaffinity()</i> .
<i>SCHED_RESET_ON_FORK</i>	Yes	No	See Section 35.3.2.
Allowed CPUs	Yes	Yes	See <i>cpuset(7)</i> .
Allowed memory nodes	Yes	Yes	See <i>cpuset(7)</i> .
Memory policy	Yes	Yes	See <i>set_mempolicy(2)</i> .
File leases	Yes	See notes	<i>fcntl(F_SETLEASE)</i> . Child inherits a reference to the same lease as parent.
Directory change notifications	Yes	No	The <i>dnotify</i> API, available via <i>fcntl(F_NOTIFY)</i> .
<i>prctl(PR_SET_DUMPABLE)</i>	See notes	Yes	During an <i>exec()</i> , the <i>PR_SET_DUMPABLE</i> flag is set, unless execing a set-user-ID or set-group-ID program, in which case it is cleared.
<i>prctl(PR_SET_PDEATHSIG)</i>	Yes	No	
<i>prctl(PR_SET_NAME)</i>	No	Yes	
<i>oom_adj</i>	Yes	Yes	See Section 49.9.
<i>coredump_filter</i>	Yes	Yes	See Section 22.1.

## 28.5 Summary

When process accounting is enabled, the kernel writes an accounting record to a file for each process that terminates on the system. This record contains statistics on the resources used by the process.

Like *fork()*, the Linux-specific *clone()* system call creates a new process, but allows finer control over which attributes are shared between the parent and child. This system call is used primarily for implementing threading libraries.

We compared the speed of process creation using *fork()*, *vfork()*, and *clone()*. Although *vfork()* is faster than *fork()*, the time difference between these system calls is small by comparison with the time required for a child process to do a subsequent *exec()*.

When a child process is created via *fork()*, it inherits copies of (or in some cases shares) certain process attributes from its parent, while other process attributes are not inherited. For example, a child inherits copies of its parent's file descriptor table and signal dispositions, but doesn't inherit its parent's interval timers, record locks, or set of pending signals. Correspondingly, when a process performs an *exec()*, certain process attributes remain unchanged, while others are reset to defaults. For example, the process ID remains the same, file descriptors remain open (unless marked close-on-exec), interval timers are preserved, and pending signals remain pending, but handled signals are reset to their default disposition and shared memory segments are detached.

### Further information

Refer to the sources of further information listed in Section 24.6. Chapter 17 of [Frisch, 2002] describes the administration of process accounting, as well as some of the variations across UNIX implementations. [Bovet & Cesati, 2005] describes the implementation of the *clone()* system call.

## 28.6 Exercise

- 28-1. Write a program to see how fast the *fork()* and *vfork()* system calls are on your system. Each child process should immediately exit, and the parent should *wait()* on each child before creating the next. Compare the relative differences for these two system calls with those of Table 28-3. The shell built-in command *time* can be used to measure the execution time of a program.