

Изучение создания процесса UNIX

Анализ жизненного цикла процесса, запускаемого операционной системой UNIX



Шон Волберг

Опубликовано 01.03.2007

Одна из многочисленных обязанностей системного администратора – это обеспечивать правильный запуск программ пользователей. Эта задача усложняется наличием в системе других одновременно выполняющихся программ. По разным причинам эти программы могут не работать, зависать или быть неисправными. Понимание процесса создания, управления и уничтожения этих заданий в операционной системе UNIX® является важнейшим шагом к созданию более надежной системы.

Разработчики изучают, как ядро управляет процессами, еще и потому, что приложения, которые нормально работают с другими составляющими системы, требуют меньше ресурсов и не так часто вызывают проблемы у системных администраторов. Приложение, которое постоянно нужно перезапускать, потому что оно создает процессы-зомби (описываются дальше), естественно не желательно. Понимание системы UNIX означает, что управляемые процессы позволяют разработчикам создавать программы, которые спокойно выполняются в фоновом режиме. Необходимость в сеансе работы с терминалом, который должен отображаться на чем-то экране, отпадает.

Основным компоновочным блоком управления этих программ является процесс. Процесс – это имя, присвоенное программе, выполняемой операционной системой. Если вы знаете команду `ps`, вам должен быть знаком список процессов, такой как в [листинге 1](#).

Листинг1. Вывод команды `ps`

```
1
2  sunbox#ps -ef
3      UID    PID  PPID    C   STIME TTY          TIME CMD
4      root     0     0    0  20:15:23 ?           0:14 sched
5      root     1     0    0  20:15:24 ?           0:00 /sbin/init
6      root     2     0    0  20:15:24 ?           0:00 pageout
7      root     3     0    0  20:15:24 ?           0:00 fsflush
8  daemon  240     1    0  20:16:37 ?           0:00 /usr/lib/nfs/statd
9  ...
```

Для рассмотрения важны первые три колонки. В первой находится список пользователей, от имени которых работают процессы, во второй перечисляются ID процессов, в третьей - ID родительских процессов. Последняя колонка содержит описание процесса, как правило, имя запущенной программы. Каждому процессу присвоен идентификатор, который называется идентификатор процесса (PID). Также у процесса есть родитель, в большинстве случаев указывается PID процесса, который запустил данный процесс.

Существование родительского PID (PPID), означает, что один процесс создается другим процессом. Исходный процесс, который запускается в системе, называется `init`, и ему всегда присваивается PID 1. `init` - это первый действительный процесс, запускаемый ядром при загрузке. Основная задача `init` - запуск всей системы. `init` и другие процессы с PPID 0 являются процессами ядра.

Использование системного вызова `fork`

Системный вызов `fork(2)` создает новый процесс. В [листинге 2](#) показан `fork` используемый в простом примере С-кода.

Листинг 2. Простое применение `fork(2)`

```
1
2  sunbox$ cat fork1.c
3  #include <unistd.h>
4  #include <stdio.h>
5
6  int main (void) {
7
8      pid_t p; /* fork returns type pid_t */
9      p = fork();
10     printf("fork returned %d\n", p);
11 }
12
13 sunbox$ gcc fork1.c -o fork1
14 sunbox$ ./fork1
15 fork returned 0
16 fork returned 698
```

Код в `fork1.c` просто вызывает `fork` и отображает целочисленный результат выполнения `fork` через вызов `printf`. Делается только один вызов, но вывод отображается дважды. Это происходит потому, что новый процесс создается в рамках вызова `fork`. После вызова возвращаются два отдельных процесса. Это часто называют "вызванный единожды, возвращается дважды."

Возвращаемые `fork` значения очень интересны. Одно из них - 0; другое – ненулевое значение. Процесс, который получает 0, называется *порожденным процессом*, а ненулевое значение достается исходному процессу, который является *родительским процессом*. Вы используете возвращаемые значения, для того чтобы определить, где какой процесс. Поскольку оба процесса возобновляют выполнение в одной и той же области, единственный возможный дифференциатор это возвращаемые значения `fork`.

Логическим основанием для нулевого и ненулевого возвращаемого значения служит то, что порожденный процесс всегда может вычислить своего родителя с помощью запроса `getppid(2)`, однако родителю намного сложнее определить всех своих потомков. Таким образом, родитель узнает о своем новом потомке, и потомок при необходимости может отыскать своего родителя.

Теперь, зная о возвращаемом значении `fork`, код может различать порожденный и родительский процессы и вести себя соответственно. В [листинге 3](#) показана программа, которая отображает разные выводы, основанные на результатах `fork`.

Листинг 3. Более полный пример использования `fork`

```
1  sunbox$ cat fork2.c
2  #include <unistd.h>
3  #include <stdio.h>
4
5  int main (void) {
6
7      pid_t p;
8
9      printf("Original program, pid=%d\n", getpid());
10     p = fork();
11     if (p == 0) {
12         printf("In child process, pid=%d, ppid=%d\n",
13               getpid(), getppid());
14     } else {
15         printf("In parent, pid=%d, fork returned=%d\n",
16               getpid(), p);
17     }
```

```

15         }
16     }
17
18     sunbox$ gcc fork2.c -o fork2
19     sunbox$ ./fork2
20     Original program, pid=767
21     In child process, pid=768, ppid=767
22     In parent, pid=767, fork returned=768
23
24

```

В [листинге 3](#) распечатываются PID распечатываются на каждом шаге, и код проверяет возвращаемые `fork` значения, для того чтобы определить, какой процесс родительский, а какой порожденный (дочерний). Сравнивая распечатанные PID, вы можете увидеть, что исходный процесс – это родительский процесс (PID 767), и порождаемый процесс (PID 768) знает, кто его родитель. Обратите внимание на то, как потомок находит своего родителя при помощи `getppid`, и как родитель использует результат `fork` для поиска своих потомков.

Теперь, когда вы разобрались с методом дублирования процессов, давайте рассмотрим, как выполнять разные процессы. `fork` - это только половина уравнения. Семейство системных вызовов `exec` запускает определенную программу.

Использование семейства системных вызовов `exec`

Задачей `exec` является замена текущего процесса на новый процесс. Отметьте использование слова *заменить*. Как только вы вызываете `exec`, текущий процесс завершается и начинается новый. Если вы хотите создать отдельный процесс, сначала вы должны вызвать `fork`, затем вызвать `exec` для новой программы в дочернем процессе. В [листинге 4](#) показан этот сценарий.

Листинг 4. Запуск разных программ посредством соединения `fork` с `exec`

```

1
2     sunbox$ cat execl.c
3     #include <unistd.h>
4     #include <stdio.h>
5
6     int main (void) {
7
8         /* Определить массив с завершающим нулем команды для запуска
9         следующим за любым параметром, в этом случае никаким */
10        char *arg[] = { "/usr/bin/ls", 0 };
11
12        /* fork и exec в порожденном процессе */
13        if (fork() == 0) {
14            printf("In child process:\n");
15            execl(arg[0], arg);
16            printf("I will never be called\n");
17        }
18        printf("Execution continues in parent process\n");
19    }
20
21     sunbox$ gcc execl.c -o execl
22     sunbox$ ./execl
23     В порожденном процессе:
24     fork1.c      execl      fork2      execl.c      fork1
25     fork2.c
26     Выполнение продолжается в родительском процессе
27

```

Код в [листинге 4](#) прежде всего, определяет массив, первый элемент которого является путем к исполняемой программе, а остальные элементы представляют собой параметры командной строки. Массив заканчивается нулевым символом. После возврата от системного вызова `fork` дочерний процесс должен запустить новую программу с помощью `execv`.

При вызове `execv` в первую очередь получает указатель на строку с именем программы для запуска, а затем указатель на массив параметров, который вы задали ранее. Первый элемент массива фактически является именем программы, следовательно, параметры начинают перечисляться со второго элемента. Запомните, что порожденный процесс никогда не возвращается после вызова `execv`. Это означает, что выполняемый процесс заменяется на новый. Есть и другие системные вызовы `exec`. Они отличаются способом приема параметров и вопросом относительно необходимости передачи переменных окружения. `execv(2)` - это один из самых простых способов замены текущего образа процесса, поскольку он использует массив с завершающим null-символом и не требует информации об окружении. Другие варианты: `execl(2)`, который принимает параметры в отдельные аргументы или `execvp(2)`, который также принимает массив переменных окружения с завершающим null-символом. Причем, не каждая операционная система поддерживает все варианты. Выбор зависит от платформы, способа программирования и от того, нужно ли вам определять какие-либо переменные окружения.

Что происходит с открытыми файлами, когда вызывается `fork`?

Когда процесс дублируется, ядро создает копии всех дескрипторов открытых файлов. Дескриптор файла - это целое число, которое является ссылкой на открытый файл или устройство, и используется для чтения и записи. Если открытие файла находится в программе перед `fork`, то что произойдет, если оба процесса попытаются читать или записывать? Перезапишет один из процессов данные другого? Будут прочитаны две копии файла? Эти вопросы изучаются в [листинге 5](#), посредством открытия двух файлов - одного для чтения, другого для записи, а также при помощи одновременного чтения и записи родителем и потомком.

Листинг 5. Два процесса одновременно записывают и читают один и тот же файл.

```

1  #include <stdio.h>
2  #include <strings.h>
3  #include <unistd.h>
4  #include <sys/types.h>
5  #include <sys/stat.h>
6  #include <fcntl.h>
7
8  int main(void) {
9      int fd_in, fd_out;
10     char buf[1024];
11
12     memset(buf, 0, 1024); /* пустой буфер */
13     fd_in = open("/tmp/infile", O_RDONLY);
14     fd_out = open("/tmp/outfile", O_WRONLY|O_CREAT);
15
16     fork(); /* Потомок против родителя значения не имеет */
17
18     while (read(fd_in, buf, 2) > 0) { /* Цикл через infile */
19         printf("%d: %s", getpid(), buf);

```

```

18             /* Написать строку */
19             sprintf(buf, "%d Hello, world!\n\r", getpid());
20             write(fd_out, buf, strlen(buf));
21             sleep(1);
22             memset(buf, 0, 1024); /* пустой буфер*/
23         }
24         sleep(10);
25     }
26
27 sunbox$ gcc fdtest1.c -o fdtest1
28 sunbox$ ./fdtest1
29 2875: 1
30 2874: 2
31 2875: 3
32 2874: 4
33 2875: 5
34 2874: 6
35 2874: 7
36 sunbox$ cat /tmp/outfile
37 2875 Hello, world!
38 2874 Hello, world!
39 2875 Hello, world!
40 2874 Hello, world!
41 2875 Hello, world!
42 2874 Hello, world!
43 2875 Hello, world!
44 2874 Hello, world!
45 2874 Hello, world!
46

```

[Листинг 5](#) - это простая программа, открывающая файл и вызывающая `fork`.

Каждый процесс использует для чтения один и тот же дескриптор файла (это просто текстовый файл с числами от 1 до 7), печатая то, что было прочитано вместе с PID. После прочтения строки PID записывается в выходной файл (out file). Цикл завершается, когда в файле больше не остается непрочитанных символов.

Результаты работы кода [листинга 5](#) показывают, что, таким образом, когда один процесс читает из файла, указатель файла смещается для обоих процессов.

Также, когда файл записывается, каждый следующий символ добавляется в конец файла. Это имеет смысл, поскольку ядро отслеживает информацию об открытом файле. Дескриптор файла это просто идентификатор для процесса.

Вы должны быть также знаете о том, что стандартный вывод (экран) тоже имеет дескриптор файла. Это видно во время вызова `fork`: оба процесса могут выводить информацию на экран.

Смерть родителя или потомка

В какой-то момент процессы должны завершаться. Вопрос только в том, какой завершится первым: родитель или потомок.

Родительский процесс завершается раньше потомка

Если родительский процесс умирает раньше своих потомков, осиротевшие потомки должны знать, кто их родитель. Вспомните о том, что у каждого процесса есть родитель, и вы можете полностью отследить дерево процессов, начиная с PID 1, имеющего также название `init`. Когда родитель умирает, `init` усыновляет всех его потомков, как показано в [листинге 6](#).

Листинг 6. Родительский процесс умирает раньше потомка

```
1
2
3     #include <unistd.h>
4     #include <stdio.h>
5
6     int main(void) {
7         int i;
8         if (fork()) {
9             /* Родитель */
10            sleep(2);
11            _exit(0);
12        }
13        for (i=0; i < 5; i++) {
14            printf("My parent is   %d\n",  getppid());
15            sleep(1);
16        }
17    }
18
19    sunbox$ gcc diel.c -o diel
20    sunbox$ ./diel
21    My parent is   2920
22    My parent is   2920
23    sunbox$ My parent is 1
24    My parent is   1
25    My parent is   1
26    My parent is   1
27    My parent is   1
28
```

В этом примере родительский процесс вызывает `fork`, ждет две секунды и завершается. Порожденный процесс продолжается, распечатывая PID своего родителя в течение пяти секунд. Вы можете видеть, что когда родитель умирает, PPID изменяется на 1. Также интересно возвращение управления командному процессору. Поскольку порожденный процесс выполняется в фоне, как только родитель умирает, управление возвращается к командному процессору.

Потомок умирает раньше родителя

[Листинг 7](#) описывает процесс противоположный [листингу 6](#): смерть потомка раньше родителя. Для того чтобы лучше показать, что происходит, непосредственно из процесса ничего не распечатывается. Вместо этого, интересная информация находится в списке процессов.

Листинг 7. Порожденный процесс умирает раньше родительского

```
1     sunbox$ cat die2.c
2     #include <unistd.h>
3     #include <stdio.h>
4
5     int main(void) {
6         int i;
7         if (!fork()) {
8             /* Потомок немедленно завершается */
9             _exit(0);
10        }
11        /* Родитель ждет около минуты */
12        sleep(60);
13    }
14
15    sunbox$ gcc die2.c -o die2
16    sunbox$ ./die2 &
17    [1] 2934
18    sunbox$ ps -ef | grep 2934
19    sean  2934  2885    0 21:43:05 pts/1        0:00 ./die2
```

```

18      sean  2935  2934  0      - ?      0:00 <defunct>
19  sunbox$ ps -ef | grep 2934
20  [1]+  Exit 199      ./die2
21
22
23

```

die2 выполняется в фоновом режиме, используя оператор `&`, после этого на экран выводится список процессов, отображая только выполняемый процесс и его потомков. PID 2934 – родительский процесс, PID 2935 – процесс, который создается и немедленно завершается. Несмотря на преждевременный выход, порожденный процесс все еще находится в таблице процессов, уже как *умерший* процесс, который еще называется *зомби*. Когда через 60 секунд родитель умирает, оба процесса завершаются.

Когда порожденный процесс умирает, его родитель информируется при помощи сигнала, который называется `SIGCHLD`. Точный механизм всего этого сейчас не имеет значения. Что действительно важно, так это то, что родитель должен как-то узнать о смерти потомка. С момента смерти потомка и до того момента как родитель получает сигнал, потомок находится в состоянии зомби. Зомби не выполняется и не потребляет ресурсов CPU; он только занимает пространство в таблице процессов. Когда родитель умирает, ядро наконец-то может убрать потомков вместе с родителем. Значит, единственный способ избавиться от зомби - это убить родителя. Лучший способ справиться с зомби – гарантировать, что они не окажутся на первом месте. Код в [листинге 8](#) описывает обработчик сигналов, для работы с входящим сигналом `SIGCHLD`.

Листинг 8. Обработчик сигналов в действии

```

1
2
3  #include <unistd.h>
4  #include <stdio.h>
5  #include <sys/types.h>
6  #include <sys/wait.h>
7
8  void sighandler(int sig) {
9      printf("In signal handler for signal  %d\n",  sig);
10     /* wait() это основное для подтверждения SIGCHLD */
11     wait(0);
12 }
13
14 int main(void) {
15     int i;
16     /* Установить обработчик сигнала к SIGCHLD */
17     sigset(SIGCHLD,  &sighandler);
18     if (!fork()) {
19         /* Потомок */
20         _exit(0);
21     }
22     sleep(60);
23 }
24 sunbox$ gcc die3.c -o die3
25 sunbox$ ./die3 &
26 [1] 3116
27 sunbox$ In signal handler for signal  18
28 ps -ef | grep 3116
    sean  3116  2885  0 22:37:26 pts/1      0:00 ./die3

```

Листинг 8 немного сложнее, чем предыдущий пример, поскольку там есть функция `sigset`, которая устанавливает указатель функции на обработчик сигнала. Всякий раз, когда процесс получает обработанный сигнал, вызывается функция, заданная через `sigset`. Для сигнала `SIGCHLD`, приложение должно вызвать функцию `wait(3c)` для того, чтобы подождать завершения порожденного процесса. Поскольку процесс уже завершен, это необходимо для того, чтобы ядро получило подтверждение о смерти потомков. На самом деле, родителю следовало бы сделать больше, чем просто подтвердить сигнал. Ему следовало бы также очистить данные потомка.

После выполнения `die3`, проверяется список процессов. Обработчик сигнала получает значение 18 (`SIGCHLD`), подтверждение о завершении потомка сделано, и родитель возвращается в состояние ожидания `sleep(60)`.

Краткие выводы

Процессы UNIX создаются, когда один процесс вызывает `fork`, который разделяет выполняемый процесс на два. После этого процесс может выполнить один из системных вызовов в семействе `exec`, который заменяет текущий образ на новый. Когда родительский процесс умирает, всех его потомков усыновляет `init`, имеющий PID 1. Если потомок умирает раньше родителя, родительскому процессу передается сигнал, а потомок переходит в состояние зомби до тех пор, пока сигнал не подтвердится или родительский процесс не будет убит.

Теперь, когда вы знаете, как создаются и уничтожаются процессы, вам будет проще разобраться в процессах, работающих в вашей системе. Особенно это касается процессов со сложной структурой, которые усложняются множественным взовом других процессов, таких как Apache. Возможность прослеживать дерево процессов для какого-то отдельного процесса, позволяет отследить любое приложение до процесса