

Каналы (Pipes)

4 апреля 2025 г.

1 Каналы между процессами

1.1 Механизм связи процессов

Мы умеем соединять процессы конвейером в терминале и направлять вывод одного процесса на вход другому. В прошлый раз мы выясняли, что эти процессы связаны неименованным каналом, который представляет буфер в ядре, в который первый из запускаемых процессов пишет, а второй читает.

1.2 Завершение работы процессов при конвейеризации

Давайте разберемся, как получается, что когда мы соединяем два процесса, то они все как будто знают, когда закончить работу.

Допустим, отработала программа `echo`, высказав в канал все, что могла сказать, и программа `wc` знает, что ее ввод закончился:

```
1 src$ echo "hello" | wc -l
2 1
```

Запустим программу `yes`, которая печатает бесконечный вывод, и направим ее вывод программе `head`, которая возьмет только первые 10 строк. Программа `yes` тоже знает, что пора заканчиваться, и не надо больше печатать бесконечный поток строчек “y”:

```
1 src$ yes | head
2 y
3 y
4 y
5 y
6 y
7 y
8 y
9 y
10 y
11 y
```

Оказывается, у именованных каналов есть специальное поведение: ядро считает, сколько процессов и файловых дескрипторов во всех процессах в системе указывают на вход именованного канала, и когда таких не остается, то любые попытки чтения из читающего конца канала возвращают конец файла. Таким образом, если мы попытаемся запустить этот конвейер и посмотреть, какие операции `read` делает тот самый `word count`, то мы увидим, что он прочитал из своего `stdin` (именованного канала) те 6 байт, которые написала `echo`, а после этого программа `echo` завершилась, владельцев пишущего конца пайпа не осталось, и следующий `read` в этом процессе вернул нуль, получился конец файла, и программа `wc` завершилась.

Обратно, если исходная программа дописала свой вывод и завершилась, то следующая получит EOF и поймет, что пора обрабатывать данные. Попробуем использовать это в коде. Построим тот самый конвейер, используя системные вызовы. Используем системный вызов `pipe` чтобы добыть себе 2 конца — пишущий и читающий от одного неименованного канала. После этого мы создадим дочерний процесс и начнем этому дочернему процессу писать какой-то вывод. После этого дождемся его завершения и завершимся. А в дочернем процессе, соответственно, будем читать из читающего конца этого пайпа в какой-то буфер и печатать то, что мы прочитали на экран.

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <sys/wait.h>
4
5 int main() {
6     pid_t child;
7     int fds[2];
8     pipe(fds);
9
10    if (child = fork()) {
11        // parent process
12        for (int i = 0; i < 5; ++i) {
13            write(fds[1], "hello", 5);
14            sleep(1);
15        }
16        wait(NULL);
17    } else {
18        // child process
19        while (1) {
20            char buf[10] = {0};
21            if (read(fds[0], buf, sizeof(buf)) <= 0) {
22                return 0;
23            }
24            puts(buf);
25        }
26    }
27 }
```

Попробуем это запустить. Родительский процесс создает неименованный канал с помощью системного вызова `pipe`, который записывает в массив `fds` два дескриптора: для чтения (`fds[0]`) и для записи (`fds[1]`). Затем создается дочерний процесс с помощью `fork()`. Родительский процесс пишет строку “hello” в пайп пять раз с интервалом в 1 секунду, а затем ждет завершения дочернего процесса через вызов `wait()`. Дочерний процесс входит в бесконечный цикл, где постоянно пытается прочитать данные из пайпа. Проблема в том, что в этом коде возникает классический deadlock (взаимная блокировка): после того, как родительский процесс записал все 5 сообщений и вызвал `wait(NULL)`, он ожидает завершения дочернего процесса. Но дочерний процесс продолжает работать в бесконечном цикле, ожидая новых данных из пайпа. Поскольку родитель не закрывает дескриптор записи в пайп (`fds[1]`), дочерний процесс никогда не получит EOF при чтении из пайпа и будет вечно ждать данных. В результате `read` в дочернем процессе никогда не завершится, а, следовательно, и `wait` в родительском процессе тоже будет ждать вечно. Попробуем это исправить. Закроем пишущий конец пайпа. Теперь наша ОС должна догадаться, что больше никаких данных из пайпа ожидать не стоит.

```

1 #include <unistd.h>
2 #include <stdio.h>
3 #include <sys/wait.h>
4
5 int main() {
6     pid_t child;
7     int fds[2];
8     pipe(fds);
9
10    if (child = fork()) {
11        // parent process
12        for (int i = 0; i < 5; ++i) {
13            write(fds[1], "hello", 5);
14            sleep(1);
15        }
16        close(fds[1]); // closing the descriptor after usage
17        wait(NULL);
18    } else {
19        // child process
20        while (1) {
21            char buf[10] = {0};
22            if (read(fds[0], buf, sizeof(buf)) <= 0) {
23                return 0;
24            }
25            puts(buf);
26        }
27    }
28 }

```

После запуска увидим, что мы закрыли конец пайпа, но программа все равно зависла. Когда мы произвели форк, то файловый дескриптор на запись скопировался из родительского процесса в процесс - потомок. Операционная система не знает, что никто больше не пишет в конец этого пайпа.

```

1 #include <unistd.h>
2 #include <stdio.h>
3 #include <sys/wait.h>
4
5 int main() {
6     pid_t child;
7     int fds[2];
8     pipe(fds);
9
10    if (child = fork()) {
11        // parent process
12        for (int i = 0; i < 5; ++i) {
13            write(fds[1], "hello", 5);
14            sleep(1);
15        }
16        close(fds[1]); // closing the descriptor after usage
17        wait(NULL);
18    } else {
19        // child process
20        close(fds[1]);
21        while (1) {
22            char buf[10] = {0};
23            if (read(fds[0], buf, sizeof(buf)) <= 0) {
24                return 0;
25            }
26            puts(buf);
27        }
28    }
29 }

```

Снова запустим программу и убедимся, что теперь все хорошо. Теперь сделаем то же самое, но в другую сторону. Пусть родительский процесс вечно пишет, а дочерний процесс читает только несколько строк.

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <sys/wait.h>
4
5 int main() {
6     pid_t child;
7     int fds[2];
8     pipe(fds);
9
10    if (child = fork()) {
11        // parent process
12        for (int i = 0; i < 1000; ++i) {
13            write(fds[1], "hello", 5);
14            sleep(1);
15        }
16        close(fds[1]);
17        wait(NULL);
18    } else {
19        // child process
20        close(fds[1]);
21        for (int i = 0; i < 5; ++i) {
22            char buf[10] = {0};
23            if (read(fds[0], buf, sizeof(buf)) <= 0) {
24                return 0;
25            }
26            puts(buf);
27        }
28    }
29 }
```

Родительский процесс не завершается. Разумеется, проблема точно такая же.

1.3 Конечный размер буфера пайпа

Важно понимать, что буфер неименованного канала (пайпа) имеет конечный размер. Из-за конечного размера буфера могут возникать ситуации, когда запись в пайп блокируется, если буфер уже заполнен и нет процесса, который бы читал из него данные. Например, если процесс, записывающий в пайп, генерирует данные быстрее, чем процесс, читающий из пайпа, может их обработать, то рано или поздно буфер пайпа заполнится, и системный вызов `write()` будет заблокирован до тех пор, пока не освободится место в буфере. Попробуем в родительском процессе закрыть читающий процесс:

```
1 int main() {
2     pid_t child;
3     int fds[2];
4     pipe(fds);
5
6     if (child = fork()) {
7         // parent process
8         close(fds[0]);
9         for (int i = 0; i < 100000; ++i) {
10            write(fds[1], "hello", 5);
11            sleep(1);
12        }
13        close(fds[1]);
14        wait(NULL);
15    }
16 }
```

Итог: лишние файловые дескрипторы на пайпы могут мешать нам организовать конвейер так, как нам бы хотелось. Будем всегда их тщательно закрывать, когда создаем дочерние процессы. Теперь попробуем запустить два дочерних процесса:

```

1 #include <unistd.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <sys/wait.h>
5
6 int main() {
7     pid_t child;
8     int fds[2];
9     pipe(fds);
10
11     if (child = fork()) {
12         // parent process
13         close(fds[0]);
14         dup2(fds[1], STDOUT_FILENO);
15         execlp("yes", "yes", NULL);
16         perror("yes");
17         exit(EXIT_FAILURE);
18     } else {
19         // child process
20         close(fds[1]);
21         dup2(fds[0], STDOUT_FILENO);
22         execlp("head", "head", NULL);
23         exit(EXIT_FAILURE);
24     }
25 }
```

Новая конструкция сработала, потому что head завершился, в yes прислали сегрпоре. Однако в конструкции не хватает того, что у программы yes и head остаются лишние файловые дескрипторы, которые хоть и не мешают, но все равно было бы неплохо их закрыть.

```

1 #include <unistd.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <sys/wait.h>
5
6 int main() {
7     pid_t child;
8     int fds[2];
9     pipe(fds);
10
11     if (child = fork()) {
12         // parent process
13         close(fds[0]);
14         dup2(fds[1], STDOUT_FILENO);
15         close(fds[1]);
16         execlp("yes", "yes", NULL);
17         perror("yes");
18         exit(EXIT_FAILURE);
19     } else {
20         // child process
21         close(fds[1]);
22         dup2(fds[0], STDOUT_FILENO);
23         close(fds[0]);
24         execlp("head", "head", NULL);
25         exit(EXIT_FAILURE);
26     }
27 }
```

Не всегда удобно считать количество дочерних процессов. Рассмотрим системные вызовы `waitid` и `waitpid`. Если мы хотим дождаться окончания всех дочерних процессов, то мы можем запустить `waitpid` со специальным значением `pid -1`, которое означает "ждать любого дочернего процесса".

```

1  } else if (!fork()) {
2      // child process 2
3      close(fds[1]);
4      dup2(fds[0], STDIN_FILENO);
5      close(fds[0]);
6      execlp("head", "head", NULL);
7      perror("head");
8      exit(EXIT_FAILURE);
9  }
10 close(fds[1]);
11 close(fds[0]);
12 while (waitpid(-1, NULL, 0) != -1);
13 }
```

Если бы мы хотели дождаться завершения только всех завершившихся дочерних процессов (например, мы регулярно их запускаем и не хотим, чтобы они висели в статусе зомби):

```

1  } else if (!fork()) {
2      // child process 2
3      close(fds[1]);
4      dup2(fds[0], STDIN_FILENO);
5      close(fds[0]);
6      execlp("head", "head", NULL);
7      perror("head");
8      exit(EXIT_FAILURE);
9  }
10 close(fds[1]);
11 close(fds[0]);
12 while (waitpid(-1, NULL, WNOHANG) > 0);
13 }
```

До сих пор мы использовали данные конвейеры для передачи текстовых данных и запуска имеющихся команд, которые с ними работают. На самом деле эти механизмы нам необходимы не для реализации Shell, а для распараллеливания работы в нескольких процессах. Чаще всего будем запускать несколько дочерних процессов и раздавать им задания через конвейер. Попробуем это сделать.

```

1  #include <unistd.h>
2  #include <stdio.h>
3  #include <stdlib.h>
4  #include <sys/wait.h>
5
6  int main() {
7      pid_t child;
8      pipe(fds);
9
10     for (int i = 0; i < 4; ++i) {
11         if (!fork()) {
12             // STDIN_FILENO=0, STDOUT_FILENO=1
13             while (read(fds[0], &task, sizeof(task)) > 0 {
14                 printf("pid %d: sleeping for %d seconds\n", getpid(), task);
15                 sleep(task);
16             }
17             return EXIT_SUCCESS;
18         }
19     }
20     for (int i = 1; i < 10; ++i) {
21         int task = i % 3 + 1;
22         write(fds[1], &task, sizeof(task));

```

```

23     }
24     close(fds[1]);
25     while (waitpid(-1, NULL, 0) != -1);
26 }

```

1.4 Обработка процессов-зомби

Некоторые задачи (если мы будем им давать задание спать дольше) могут завершиться сильно раньше чем мы начнем дожидаться завершения всех потомков. Это нас не особенно смущает, но если мы будем запускать новые процессы, создавать ротацию среди них, то эти зомби могут накапливаться. Мы могли бы периодически в родительском процессе зачищать всех зомби, запуская `waitpid` с опцией `WNOHANG`. Тогда, если у нас где-то в процессе работы родительского процесса возникли зомби, то мы их всех поawaitим.

```

1     for (int i = 1; i < 10; ++i) {
2         int task = i % 3 + 1;
3         write(fds[1], &task, sizeof(task));
4         while (waitpid(-1, NULL, WNOHANG) != -1);
5     }
6     close(fds[1]);
7     while (waitpid(-1, NULL, 0) != -1);
8 }

```

Есть еще один способ узнавать о завершении дочерних процессов. Оказывается имеется специальный сигнал `SIGCHLD` который приходит нам по завершении дочернего процесса. Мы можем повесить на него обработчик.

```

1 #include <unistd.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <string.h>
5 #include <sys/wait.h>
6
7 void handler(int sig) {
8     WRITE(STDOUT_FILENO, "reaping child\n", strlen("reaping child\n"));
9     wait(NULL);
10 }
11
12 int main() {
13     pid_t child;
14     pipe(fds);
15
16     signal(SIGCHLD, handler);
17
18     for (int i = 0; i < 4; ++i) {
19         if (!fork()) {
20             // STDIN_FILENO=0, STDOUT_FILENO=1
21             while (read(fds[0], &task, sizeof(task)) > 0 {
22                 printf("pid %d: sleeping for %d seconds\n", getpid(), task);
23                 sleep(task);
24             }
25             return EXIT_SUCCESS;
26         }
27
28         for (int i = 1; i < 10; ++i) {
29             int task = i % 3 + 1;
30             write(fds[1], &task, sizeof(task));
31         }
32         close(fds[1]);
33         while (waitpid(-1, NULL, 0) != -1);
34     }

```

Есть еще один способ завершения потомков, если нас не интересует, как они завершились, то мы можем сделать еще так:

```
signal(SIGCHLD, SIG_IGN);
```

Потомки завершаются и не переходят в статус зомби, поскольку мы явно сообщаем ОС, что их завершение нас не интересует. Соответственно нам не нужен обработчик. Если дочерние процессы завершаются раньше родительского то они просто удаляются из таблицы процессов.

1.5 Превращение сигналов в файловые дескрипторы

Есть 1 трюк и 1 механизм, которые превращают сигналы в файловые дескрипторы. Классический трюк на эту тему называется *self-pipe trick*. Вы создаете себе отдельный канал ровно для того, чтобы обрабатывать какой-нибудь сигнал, и в обработчике сигнала просто пишете в pipe 1 байт данных, тогда вы можете дожидаться прихода этого сигнала в основном коде программы, просто читая из этого пайпа. Второй механизм, чтобы превратить сигнал в файловый дескриптор, называется `signalfd`. Это нестандартный, Linux-специфичный механизм, но он крайне удобен, он работает таким способом: вы передаете ему маску сигналов, которые вы хотите обрабатывать путем чтения из файловых дескрипторов, и после этого блокируете эти сигналы, чтобы они не доставлялись обычным способом, а только через `signalfd`, после этого читаете из этого файлового дескриптора специального вида структуры и пользуетесь информацией о пришедших вам сигналах, не беспокоясь, что они приходят асинхронно.