

Программирование в Linux



Межпроцессное взаимодействие. Каналы. Очереди.
Общая память

Каналы — односторонние средства общения

Pipe

- Неименованный канал
- Два открытых дескриптора RO [0] и WO [1]
- Пользоваться могут только родственные процессы
- Существует пока открыт хоть один дескриптор

```
#include <unistd.h>

/* On Alpha, IA-64, MIPS, SuperH, and SPARC/SPARC64; see NOTES */
struct fd_pair {
    long fd[2];
};
struct fd_pair pipe();

/* On all other architectures */
int pipe(int pipefd[2]);

#define _GNU_SOURCE          /* See feature_test_macros(7) */
#include <fcntl.h>           /* Obtain O_* constant definitions */
#include <unistd.h>

int pipe2(int pipefd[2], int flags);
```

Fifo

- Именованный канал, очередь
- Объект в файловой системе
- Все, кто знает к нему путь, могут подключиться
- Существует до тех пор, пока не удалят (rm)

```
#include <sys/types.h>
#include <sys/stat.h>

int mkfifo(const char *pathname, mode_t mode);

#include <fcntl.h>          /* Definition of AT_* constants */
#include <sys/stat.h>

int mkfifoat(int dirfd, const char *pathname, mode_t mode);
```

Особые случаи

- Нет данных в канале/буфер заполнился — блокировка
- Именованный канал открыли только с одной стороны — блокировка
- Все WO дескрипторы закрыты — eof для читателя
- Все RO дескрипторы закрыты — SIGPIPE при записи
- Запись менее PIPE_BUF байт атомарна (данные не перемешиваются при записи из многих процессов)

Блокировки отключаются флагом O_NONBLOCK.

При завязке на SIGPIPE стоит выставлять O_CLOEXEC (автозаккрытие при exec)

Специальные операции для каналов

```
#define _GNU_SOURCE          /* See feature_test_macros(7) */
#include <fcntl.h>

ssize_t splice(int fd_in, loff_t *off_in, int fd_out,
               loff_t *off_out, size_t len, unsigned int flags);
```

Переслать len байт из fd_in в fd_out, избегая копирований userspace/kernelspace.
Один из дескрипторов должен быть каналом

```
#include <fcntl.h>

ssize_t tee(int fd_in, int fd_out, size_t len, unsigned int flags);
```

Дублировать len байт из fd_in в fd_out, не извлекая из fd_in.
Оба дескриптора должны быть каналами

```
#include <fcntl.h>
#include <sys/uio.h>

ssize_t vmsplice(int fd, const struct iovec *iov,
                 unsigned long nr_segs, unsigned int flags);
```

fd — один из концов pipe.
Отобразить в канал страницы адресного пространства или
Прочитать из канала в эти страницы

```

89 int main() {
90     struct pipe_ends {
91         int read_end;
92         int write_end;
93     };
94
95     pipe_ends pipe_ends;
96     if (pipe(reinterpret_cast<int(&)[2]>(pipe_ends)) != 0) {
97         perror("pipe");
98         return EXIT_FAILURE;
99     }
100
101     pid_t child = fork();
102
103     if (child > 0) {
104         close(pipe_ends.read_end);
105         WritePipe(pipe_ends.write_end);
106     }
107
108     if (child == 0) {
109         close(pipe_ends.write_end);
110         ReadPipe(pipe_ends.read_end);
111     }
112
113     perror("fork");
114     return EXIT_FAILURE;
115 }

```

```

21 static constexpr auto TerminateHandler = [](int s) {
22     terminated.store(true);
23     termination_signal.store(s);
24 };
25
26 [[noreturn]] void WritePipe(int pipefd) {
27     printf("Start writer: (%d)\n", getpid());
28     signal(SIGPIPE, TerminateHandler);
29     signal(SIGINT, SIG_IGN);
30
31     std::vector<std::vector<Point>> points = {
32         { {1, 2}, {3, 4} },
33         { {5, 6} }
34     };
35
36     std::vector<iovec> data_ptrs;
37     for (const auto& line : points) {
38         data_ptrs.push_back(iovec {
39             .iov_base = reinterpret_cast<void*>(
40                 reinterpret_cast<const void*>(
41                     line.data())),
42             .iov_len = line.size() * sizeof(Point)
43         });
44     }
45
46     while (!terminated.load()) {
47         vmsplICE(pipefd,
48             data_ptrs.data(),
49             data_ptrs.size(),
50             SPLICE_F_MORE);
51         if (terminated.load()) {
52             break;
53         }
54         sleep(1);
55     }
56
57     printf("Writer: (%d) is terminated.\n",
58         getpid());
59     printf("Termination signal is %d\n",
60         termination_signal.load());
61     close(pipefd);
62     exit(EXIT_SUCCESS);
63 }

```

```

65 [[noreturn]] void ReadPipe(int pipefd) {
66     printf("Start reader: (%d)\n", getpid());
67
68     signal(SIGTERM, TerminateHandler);
69     signal(SIGINT, TerminateHandler);
70
71     Point cur;
72     while (!terminated.load()) {
73         auto sz = read(pipefd,
74             reinterpret_cast<void*>(&cur),
75             sizeof(cur));
76         if (terminated.load()) {
77             break;
78         }
79         if (sz != sizeof(cur)) {
80             printf("Reader: (%d).\n",
81                 getpid());
82             printf("read something unexpected",
83                 getpid());
84             break;
85         }
86         printf("Reader: (%d) got Point: (%d, %d)\n",
87             getpid(), cur.x, cur.y);
88         sleep(1);
89     }
90     if (terminated.load()) {
91         printf("Reader: (%d) is manually terminated\n",
92             getpid());
93     } else {
94         printf("Reader: (%d) is spuriously terminated\n",
95             getpid());
96     }
97     close(pipefd);
98     exit(EXIT_SUCCESS);
99 }

```

```

125 int main() {
126     signal(SIGINT, SIG_IGN);
127     signal(SIGTERM, SIG_IGN);
128
129     if (mkfifo(FIFO_FILE_NAME, 0644) != 0) {
130         perror("mkfifo");
131         return EXIT_FAILURE;
132     }
133
134     struct DeferDeletion {
135         ~DeferDeletion() {
136             remove(FIFO_FILE_NAME);
137         }
138     } deleter;
139
140     pid_t reader = clone([], (void*){
141         RunReader();
142     }, reader_stack + STASK_SIZE, SIGCHLD, NULL);
143
144     if (reader < 0) {
145         perror("clone");
146         return EXIT_FAILURE;
147     }
148
149     pid_t writer = clone([], (void*){
150         RunWriter();
151     }, writer_stack + STASK_SIZE, SIGCHLD, NULL);
152
153     if (writer < 0) {
154         perror("clone");
155         kill(reader, SIGKILL);
156         return EXIT_FAILURE;
157     }
158
159     waitpid(reader, NULL, 0);
160     waitpid(writer, NULL, 0);
161     return EXIT_SUCCESS;
162 }

```

```

33 [[noreturn]] void RunWriter() {
34     printf("Start writer (%d)\n", getpid());
35     signal(SIGPIPE, TerminateHandler);
36     signal(SIGINT, SIG_IGN);
37
38     int pipefd = open(FIFO_FILE_NAME, O_WRONLY);
39     if (pipefd < 0) {
40         perror("open");
41         std::exit(EXIT_FAILURE);
42     }
43
44 > std::vector<std::vector<Point>> points = {...
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78 [[noreturn]] void RunReader() {
79     printf("Start reader (%d)\n", getpid());
80
81     signal(SIGTERM, TerminateHandler);
82     signal(SIGINT, TerminateHandler);
83
84     int pipefd = open(FIFO_FILE_NAME, O_RDONLY);
85     if (pipefd < 0) {
86         perror("open");
87         std::exit(EXIT_FAILURE);
88     }
89
90     Point cur;
91 > while (!terminated.load()) {...
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108 > if (terminated.load()) {...
109
110
111
112
113
114
115     close(pipefd);
116     exit(EXIT_SUCCESS);
117 }

```


“popen” — подсоединяемся к stdin/stdout запускаемого процесса

```
16 struct pipe_ends {
17     int read_end;
18     int write_end;
19 };
20
21 pid_t MyPopen(const char* command, pipe_ends* communication_pipe) {
22     pipe_ends stdin_ends;
23     pipe_ends stdout_ends;
24     if (pipe2(reinterpret_cast<int(&)[2]>(stdin_ends), O_CLOEXEC) != 0) {
25         return -1;
26     }
27     if (pipe2(reinterpret_cast<int(&)[2]>(stdout_ends), O_CLOEXEC) != 0) {
28         return -1;
29     }
30     communication_pipe->read_end = stdout_ends.read_end;
31     communication_pipe->write_end = stdin_ends.write_end;
32
33     pid_t child = fork();
34     if (child < 0) {
35         return -1;
36     }
37
38     if (child != 0) {
39         close(stdin_ends.read_end);
40         close(stdout_ends.write_end);
41         return child;
42     }
43
44     dup2(stdin_ends.read_end, STDIN_FILENO); // redefine STDIN/STDOUT
45     dup2(stdout_ends.write_end, STDOUT_FILENO);
46     execl("/bin/sh", "sh", "-c", command, NULL);
47     perror("execl");
48     exit(EXIT_FAILURE);
49 }
```

dup(fd) — клонировать дескриптор
dup2(fd, newfd) — клонировать дескриптор
fd, назначить клону значение newfd

Очередь сообщений — высокоуровневый named pipe

Общение не отдельными байтами, а цельными пачками — сообщениями

POSIX (man mq_overview)

Library interface	System call
<code>mq_close(3)</code>	<code>close(2)</code>
<code>mq_getattr(3)</code>	<code>mq_getsetattr(2)</code>
<code>mq_notify(3)</code>	<code>mq_notify(2)</code>
<code>mq_open(3)</code>	<code>mq_open(2)</code>
<code>mq_receive(3)</code>	<code>mq_timedreceive(2)</code>
<code>mq_send(3)</code>	<code>mq_timedsend(2)</code>
<code>mq_setattr(3)</code>	<code>mq_getsetattr(2)</code>
<code>mq_timedreceive(3)</code>	<code>mq_timedreceive(2)</code>
<code>mq_timedsend(3)</code>	<code>mq_timedsend(2)</code>
<code>mq_unlink(3)</code>	<code>mq_unlink(2)</code>

System V (man svipc)

<code>msgget</code>	— открыть/создать очередь
<code>msgctl</code>	— настроить параметры
<code>msgsnd</code>	— отправить сообщение
<code>msgrcv</code>	— принять сообщение

Существуют до тех пор, пока не удалят, или до выгрузки OS

POSIX mq — очередь с приоритетами


```

19 constexpr auto MQQUEUE_NAME = "/my_queue";
20
21 struct Point
22 {
23     int x;
24     int y;
25 };

```

```

152 int main()
153 {
154     signal(SIGINT, SIG_IGN);
155     signal(SIGTERM, SIG_IGN);
156
157     mq_attr queue_attrs;
158     queue_attrs.mq_maxmsg = 10; // max messages
159     queue_attrs.mq_msgsize = sizeof(Point);
160     mqd_t mq = mq_open(MQQUEUE_NAME,
161                       O_CREAT | O_RDWR,
162                       0644, &queue_attrs);
163     if (mq == mqd_t(-1))
164     {
165         perror("mq_open");
166         return EXIT_FAILURE;
167     }
168     mq_close(mq);
169     struct DeferDeletion
170     {
171         ~DeferDeletion()
172         {
173             mq_unlink(MQQUEUE_NAME);
174         }
175     } deleter;
176

```

```

35 [[noreturn]] void RunWriter()
36 {
37     printf("Start writer (%d)\n", getpid());
38     signal(SIGINT, TerminateHandler);
39     signal(SIGTERM, TerminateHandler);
40
41     int pipefd = mq_open(MQQUEUE_NAME, O_WRONLY);
42     if (pipefd < 0) ...
43     mq_attr attrs;
44     mq_getattr(pipefd, &attrs);
45     printf("Queue attrs: maxmessagesize=%ld, "
46           "max_messages=%ld\n",
47           attrs.mq_msgsize,
48           attrs.mq_maxmsg);
49
50     if (sizeof(Point) != attrs.mq_msgsize) ...
51
52     std::vector<Point> points =
53         {{1, 2}, {3, 4}, {5, 6}};
54
55     while (!terminated.load())
56     {
57         for (const auto &p : points)
58         {
59             mq_send(pipefd,
60                   (const char *)&p,
61                   sizeof(p), 0);
62         }
63         sleep(1);
64     }
65
66     printf("Writer (%d) is terminated."
67           "Termination signal is %d\n",
68           getpid(),
69           termination_signal.load());
70     mq_close(pipefd);
71     exit(EXIT_SUCCESS);
72 }

```

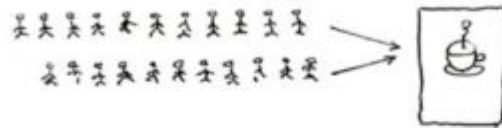
```

82 [[noreturn]] void RunReader()
83 {
84     printf("Start reader (%d)\n", getpid());
85     signal(SIGTERM, TerminateHandler);
86     signal(SIGINT, SIG_IGN);
87
88     mqd_t pipefd = mq_open(MQQUEUE_NAME, O_RDONLY);
89     if (pipefd < 0) ...
90     mq_attr attrs;
91     mq_getattr(pipefd, &attrs);
92     printf("Queue attrs: maxmessagesize=%ld, "
93           "max_messages=%ld\n",
94           attrs.mq_msgsize,
95           attrs.mq_maxmsg);
96
97     if (sizeof(Point) != attrs.mq_msgsize) ...
98
99     Point cur;
100     timespec timeout{1, 0};
101     while (!terminated.load())
102     {
103         auto sz = mq_timedreceive(pipefd,
104                                   (char *)&cur,
105                                   sizeof(cur),
106                                   nullptr, // ignore prio
107                                   &timeout);
108         if (sz < 0 && errno == ETIMEDOUT)
109         {
110             printf("Reader timeout!\n");
111             break;
112         }
113         if (terminated.load()) ...
114         if (sz != sizeof(cur)) ...
115         printf("Reader (%d) got Point {%d, %d}\n",
116               getpid(), cur.x, cur.y);
117         sleep(1);
118     }
119     if (terminated.load()) ...
120     else ...
121     mq_close(pipefd);
122     exit(EXIT_SUCCESS);
123 }

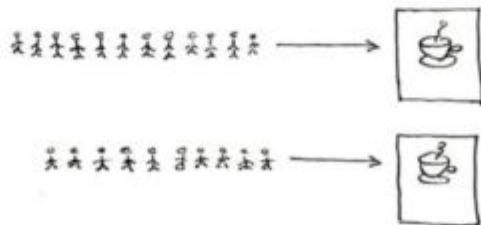
```

Конкурентное vs параллельное исполнение

Concurrent = Two Queues One Coffee Machine



Parallel = Two Queues Two Coffee Machines



© Joe Armstrong 2013

Shared memory

Анонимная

```
mmap(... MAP_ANONYMOUS  
| MAP_SHARED)
```

Существует пока все
использующие процессы не
размапят или не завершатся.

Именованная

POSIX: (man shm_overview)

```
fd = shm_open("/name")  
ftruncate(fd, N)  
ptr = mmap(..., fd);  
...  
shm_unlink("/name");
```

Именованная

System V: (man svipc)

```
shmget — создать  
shmat  — получить отображение  
shmctl — настроить  
shmdt  — размапить отображение
```

Именованная память существует, пока не удалят или до выгрузки OS

Синхронизация конкурентного доступа

Семафор — счетчик с блокировкой. Можно увеличивать и уменьшать. Попытка опустить значение ниже нуля влечет блокировку.

Мьютекс — бинарный семафор.

POSIX семафоры (man sem_overview)

sem_open — открыть/создать именованный семафор

sem_close

sem_unlink — удалить именованный семафор

sem_init — создать анонимный семафор в указанной области памяти

sem_destroy

sem_post — увеличить счетчик на единицу

sem_wait — уменьшить счетчик на единицу

Есть System V семафоры.

Fast userspace mutex (futex)

- Семафоры — объекты ядра. Работа с ними — по два переключения в kernelspace.
- Переключение нужно только для блокировки/уведомления.
- Часто семафор свободен. И переключение в kernelspace избыточно
- Можно перенести проверку в userspace

```
23 static int
24 futex(uint32_t *uaddr, int futex_op, uint32_t val,
25       const struct timespec *timeout, uint32_t *uaddr2, uint32_t val3)
26 {
27     return syscall(SYS_futex, uaddr, futex_op, val,
28                   timeout, uaddr2, val3);
29 }
```

```
33 static void fwait(uint32_t *futexp) {
34     while (1) {
35         /* Is the futex available? */
36         const uint32_t one = 1;
37         if (atomic_compare_exchange_strong(futexp, &one, 0))
38             break; /* Yes */
39
40         /* Futex is not available; wait. */
41         long s = futex(futexp, FUTEX_WAIT, 0, NULL, NULL, 0);
42         if (s == -1 && errno != EAGAIN)
43             errExit("futex-FUTEX_WAIT");
44     }
45 }
```

```
50 static void fpost(uint32_t *futexp) {
51     const uint32_t zero = 0;
52     if (atomic_compare_exchange_strong(futexp, &zero, 1)) {
53         long s = futex(futexp, FUTEX_WAKE, 1, NULL, NULL, 0);
54         if (s == -1)
55             errExit("futex-FUTEX_WAKE");
56     }
57 }
```

IPC namespace, изоляция и просмотр через fs

Все System V объекты IPC изолируются при создании нового IPC namespace. Объекты именуются уникальными idшниками, найти их через vfs затруднительно.

POSIX:

mq_* — изолируются через IPC namespace.
Можно примонтировать vfs, для мониторинга

```
# mkdir /dev/mqueue  
# mount -t mqueue none /dev/mqueue
```

Именованные POSIX семафоры и shared memory через IPC namespace не изолируются.
Нужно изолировать через mount namespace (/dev/shm)