

10 非阻塞 IO

这章来讲另一种 IO 模型非阻塞 IO(NIO)，也就是同步非阻塞 IO。上一章说过，IO 操作的两个阶段先查询再读写，而非阻塞 IO 在查询阶段的处理和阻塞 IO 不同。应用程序需要进行 IO 操作前，先发起查询，驱动程序根据数据情况返回查询结果，如果返回查询结果 NG，应用程序就不执行读写操作了。如果应用程序非要读写的话，就继续去查询，直到驱动程序返回数据准备完成，才会做下一步的读写操作。

10.1 Linux 中的 NIO

非阻塞 IO 的处理方式是轮询。Linux 中提供了应用程序的轮询机制和相应的驱动程序系统调用。

10.1.1 应用程序中的轮询方法

应用程序中提供了三种轮询的方法：select、poll、epoll。实际上他们也是多路复用 IO 的解决方法，这里就不展开说了。

1) select

select 有良好的跨平台支持性，但是他的单个进程能够监视的文件描述符的数量存在最大限制(Linux 中一般为 1024)。

函数原型：

```
int select(int maxfdp, fd_set *readfds, fd_set *writefds, fd_set *errorfds, struct timeval *timeout);
```

参数说明：

maxfdp：是集合中所有文件描述符的范围，即等于所有文件描述符的最大值加 1。

readfds：struct fd_set 结构体可以理解为是文件描述符(file descriptor)的集合，也就是文件句柄，他的每一位就代表一个描述符，readfds 用于监视指定描述符集的读变化，只要其中有一位可读，select 就会返回一个大于 0 的值。可以已使用这些宏来操作 fd_set 变量：

```
/* 把 fdset 所有位置 0, 清楚 fdset 和所有文件句柄的关系 */
FD_ZERO(fd_set *fdset)
/* 把 fdset 某个位置 1, 把 fdset 和文件句柄 fd 关联 */
FD_SET(int fd, fd_set *fdset)
/* 把 fdset 某个位置 0, 取消 fdset 和文件句柄的关联 */
FD_CLR(int fd, fd_set *fdset)
/* 判断某个文件句柄是否为 1, 即是否可操作 */
FD_ISSET(int fd, fdset *fdset)
```

writefds：用于监视文件是否可写。

errorfds：用于监视文件异常。

timeout：struct timeval 用来代表时间值，有两个成员，一个是秒数，另一个是毫秒数。定义如下：

```
1. struct timeval{
2.     long tv_sec;    /*秒 */
3.     long tv_usec;   /*微秒 */
4. }
```

这个参数的值有三种情况

如果传入 **NULL**，即不传入时间结构，就是一个阻塞函数，直到某个文件描述符发生变化才会返回。

如果把秒和微妙都设为 **0**，就是一个非阻塞函数，会立刻返，文件无变化返回 **0**，有变化返回正值。

如果值大于 **0**，则意为超时时间，**select** 若在 **timeout** 时间内没有检测到文件描述符变化，则会直接返回 **0**，有变化则返回正值。

使用示例：

```
1. void main(void)
2. {
3.     /* ret 获取返回值, fd 获取文件句柄 */
4.     int ret, fd;
5.     /* 定义一个监视文件读变化的描述符合集 */
6.     fd_set readfds;
7.     /* 定义一个超时时间结构体 */
8.     struct timeval timeout;
9.
10.    /* 获取文件句柄, O_NONBLOCK 表示非阻塞访问 */
11.    fd = open("dev_xxx", O_RDWR | O_NONBLOCK);
12.
13.    /* 初始化描述符合集 */
14.    FD_ZERO(&readfds);
15.    /* 把文件句柄 fd 指向的文件添加到描述符 */
16.    FD_SET(fd, &readfds);
17.
18.    /* 超时时间初始化为 1.5 秒 */
19.    timeout.tv_sec = 1;
20.    timeout.tv_usec = 500000;
21.
22.    /* 调用 select, 注意第一个参数为 fd+1 */
23.    ret = select(fd + 1, &readfds, NULL, NULL, &timeout);
24.
25.    switch (ret)
26.    {
27.        case 0:
28.        {
29.            /* 超时 */
30.            break;
31.        }
```

```

32.         case -1:
33.         {
34.             /* 出错 */
35.             break;
36.         }
37.         default:
38.         {
39.             /* 监视的文件可操作 */
40.             /* 判断可操作的文件是不是文件句柄 fd 指向的文件 */
41.             if(FD_ISSET(fd, &readfds))
42.             {
43.                 /* 操作文件 */
44.             }
45.             break;
46.         }
47.     }
48. }

```

在 23 行调用 `select` 函数之前，做了很多准备工作，主要是 `select` 函数输入参数的初始化。

注意 11 行 `open` 函数输入参数中的 `O_NONBLOCK` 属性，如果需要非阻塞的访问文件，则需要添加这个属性。

41 行，在 `ret` 返回大于 0 时，使用宏定义 `FD_ISSET` 判断可操作的句柄是不是我们需要的句柄，在只等待一个文件的情况下，可以不做这个判断。

2) poll

`poll` 本质上和 `select` 没有区别，但是他的最大连接数没有限制。

函数原型：

```
int poll (struct pollfd *fds, unsigned int nfds, int timeout);
```

参数说明：

fds: `struct pollfd` 结构体是文件句柄和事件的组合，定义如下：

```

1. struct pollfd {
2.     int fd;
3.     short events;
4.     short revents;
5. };

```

`fd` 是文件句柄，`events` 是对于这个文件需要监视的事件类型，`revents` 是内核返回的事件类型。事件类型有：

```

POLLIN    //有数据可读
POLLPRI   //有紧急数据可读
POLLOUT   //数据可写
POLLERR   //指定文件描述符发生错误
POLLHUP   //指定文件描述符挂起
POLLNVAL  //无效请求

```

POLLRDNORM //有数据可读

nfds: poll 监视的文件句柄数量，也就是 fds 的数组长度。

timeout: 超时时间，单位为毫秒。

使用示例：

```
1. void main(void)
2. {
3.     /* ret 获取返回值, fd 获取文件句柄 */
4.     int ret, fd;
5.     /* 定义 struct pollfd 结构体变量 */
6.     struct pollfd fds[1];
7.
8.     /* 非阻塞访问文件 */
9.     fd = open(filename, O_RDWR | O_NONBLOCK);
10.
11.    /* 初始化 struct pollfd 结构体变量 */
12.    fds[0].fd = fd;
13.    fds[0].events = POLLIN;
14.
15.    /* 调用 poll */
16.    ret = poll(fds, sizeof(fds), 1500);
17.    if(ret == 0)
18.    {
19.        /* 超时 */
20.    }
21.    else if (ret < 0)
22.    {
23.        /* 错误 */
24.    }
25.    else
26.    {
27.        /* 操作数据 */
28.    }
29. }
```

3) epoll

可以理解为 event poll，设计用于大并发时的 IO 查询，常用于网络编程，暂不介绍。

10.1.2 驱动程序中的 poll 函数

应用程序中调用 select、poll、epoll 时，系统调用就会执行驱动程序中 file_operations 的 poll 函数。也就是我们需要实现的函数。圆原型如下：

unsigned int (*poll) (struct file *filp, struct poll_table_struct *wait)

参数说明：

filp: 应用程序传递过来的值, 应用程序 open 之后获得的目标文件句柄。

wait: 应用程序传递过来的值, 代表应用程序线程。我们需要在 poll 函数中调用 poll_wait 将应用程序添线程 wait 添加到 poll_table 等待队列中, poll_wait 函数原型如下:

```
void poll_wait(struct file * filp, wait_queue_head_t * wait_address, poll_table *p)
```

wait 作为参数 p 传递给 poll_wait 函数。

返回值: 返回值和 struct pollfd 结构体中的事件类型相同。

10.2 实验

这章的实验目标和上一张相同, 使用 ps_key1 控制 ps_led1, 并使 cpu 占用率相较于第八章降低。

10.2.1 原理图

led 部分和章节 1.3.1 相同。

key 部分和章节 6.1 相同。

10.2.2 设备树

和章节 6.2 相同。

10.2.3 驱动代码

使用 petalinux 新建名为“ax-nio-driv”的驱动程序, 并执行 petalinux-config -c rootfs 命令选上新增的驱动程序。

在 ax-nio-driv.c 文件中输入下面的代码:

```
1.  /** ===== **
2.  *Author : ALINX Electronic Technology (Shanghai) Co., Ltd.
3.  *Website: http://www.alinx.com
4.  *Address: Room 202, building 18,
5.           No.518 xinbrick Road,
6.           Songjiang District, Shanghai
7.  *Created: 2020-3-2
8.  *Version: 1.0
9.  ** ===== **/
10.
11. #include <linux/module.h>
12. #include <linux/kernel.h>
13. #include <linux/init.h>
14. #include <linux/ide.h>
15. #include <linux/types.h>
```

```

16. #include <linux/errno.h>
17. #include <linux/cdev.h>
18. #include <linux/of.h>
19. #include <linux/of_address.h>
20. #include <linux/of_gpio.h>
21. #include <linux/device.h>
22. #include <linux/delay.h>
23. #include <linux/init.h>
24. #include <linux/gpio.h>
25. #include <linux/semaphore.h>
26. #include <linux/timer.h>
27. #include <linux/of_irq.h>
28. #include <linux/irq.h>
29. #include <linux/wait.h>
30. #include <linux/poll.h>
31. #include <asm/uaccess.h>
32. #include <asm/mach/map.h>
33. #include <asm/io.h>
34.
35. /* 设备节点名称 */
36. #define DEVICE_NAME      "nio_led"
37. /* 设备号个数 */
38. #define DEVID_COUNT      1
39. /* 驱动个数 */
40. #define DRIVE_COUNT      1
41. /* 主设备号 */
42. #define MAJOR_U
43. /* 次设备号 */
44. #define MINOR_U          0
45.
46. /* 把驱动代码中会用到的数据打包进设备结构体 */
47. struct alinx_char_dev {
48. /** 字符设备框架 */
49.     dev_t          devid;          //设备号
50.     struct cdev     cdev;          //字符设备
51.     struct class    *class;        //类
52.     struct device   *device;       //设备
53.     struct device_node *nd;        //设备树的设备节点
54. /** gpio */
55.     int             alinx_key_gpio; //gpio 号
56. /** 并发处理 */
57.     atomic_t        key_sts;       //记录按键状态，为 1 时被按下
58. /** 中断 */
59.     unsigned int    irq;           //中断号

```

```

60. /** 定时器 */
61.     struct timer_list timer;           //定时器
62. /** 等待队列 */
63.     wait_queue_head_t wait_q_h;       //等待队列头
64. };
65. /* 声明设备结构体 */
66. static struct alinx_char_dev alinx_char = {
67.     .cdev = {
68.         .owner = THIS_MODULE,
69.     },
70. };
71.
72. /** 回掉 */
73. /* 中断服务函数 */
74. static irqreturn_t key_handler(int irq, void *dev)
75. {
76.     /* 按键按下或抬起时会进入中断 */
77.     /* 开启 50 毫秒的定时器用作防抖动 */
78.     mod_timer(&alinx_char.timer, jiffies + msecs_to_jiffies(50));
79.     return IRQ_RETVAL(IRQ_HANDLED);
80. }
81.
82. /* 定时器服务函数 */
83. void timer_function(unsigned long arg)
84. {
85.     /* value 用于获取按键值 */
86.     unsigned char value;
87.     /* 获取按键值 */
88.     value = gpio_get_value(alinx_char.alinx_key_gpio);
89.     if(value == 0)
90.     {
91.         /* 按键按下，状态置 1 */
92.         atomic_set(&alinx_char.key_sts, 1);
93. /** 等待队列 */
94.         /* 唤醒进程 */
95.         wake_up_interruptible(&alinx_char.wait_q_h);
96.     }
97.     else
98.     {
99.         /* 按键抬起 */
100.    }
101. }
102.
103. /** 系统调用实现 */

```

```

104. /* open 函数实现, 对应到 Linux 系统调用函数的 open 函数 */
105. static int char_drv_open(struct inode *inode_p, struct file *file_p)
106. {
107.     printk("gpio_test module open\n");
108.     return 0;
109. }
110.
111. /* read 函数实现, 对应到 Linux 系统调用函数的 read 函数 */
112. static ssize_t char_drv_read(struct file *file_p, char __user *buf, size_t len, loff_t *loff
    f_t_p)
113. {
114.     unsigned int keysts = 0;
115.     int ret;
116.
117.     /* 读取 key 的状态 */
118.     keysts = atomic_read(&alinx_char.key_sts);
119.     /* 判断文件打开方式 */
120.     if(file_p->f_flags & O_NONBLOCK)
121.     {
122.         /* 如果是非阻塞访问, 说明以满足读取条件 */
123.     }
124.     /* 判断当前按键状态 */
125.     else if(!keysts)
126.     {
127.         /* 按键未被按下(数据未准备好) */
128.         /* 以当前进程创建并初始化为队列项 */
129.         DECLARE_WAITQUEUE(queue_mem, current);
130.         /* 把当前进程的队列项添加到队列头 */
131.         add_wait_queue(&alinx_char.wait_q_h, &queue_mem);
132.         /* 设置当前进程成为可被信号打断的状态 */
133.         __set_current_state(TASK_INTERRUPTIBLE);
134.         /* 切换进程, 是当前进程休眠 */
135.         schedule();
136.
137.         /* 被唤醒, 修改当前进程状态为 RUNNING */
138.         set_current_state(TASK_RUNNING);
139.         /* 把当前进程的队列项从队列头中删除 */
140.         remove_wait_queue(&alinx_char.wait_q_h, &queue_mem);
141.
142.         /* 判断是否是被信号唤醒 */
143.         if(signal_pending(current))
144.         {
145.             /* 如果是直接返回错误 */
146.             return -ERESTARTSYS;

```



```

147.     }
148.     else
149.     {
150.         /* 被按键唤醒 */
151.     }
152. }
153. else
154. {
155.     /* 按键被按下(数据准备好了) */
156. }
157.
158. /* 读取 key 的状态 */
159. keysts = atomic_read(&alinx_char.key_sts);
160. /* 返回按键状态值 */
161. ret = copy_to_user(buf, &keysts, sizeof(keysts));
162. /* 清除按键状态 */
163. atomic_set(&alinx_char.key_sts, 0);
164. return 0;
165. }
166.
167. /* poll 函数实现 */
168. unsigned int char_drv_poll(struct file *filp, struct poll_table_struct *wait)
169. {
170.     unsigned int ret = 0;
171.
172.     /* 将应用程序添加到等待队列中 */
173.     poll_wait(filp, &alinx_char.wait_q_h, wait);
174.
175.     /* 判断 key 的状态 */
176.     if(atomic_read(&alinx_char.key_sts))
177.     {
178.         /* key 准备好了, 返回数据可读 */
179.         ret = POLLIN;
180.     }
181.     else
182.     {
183.
184.     }
185.
186.     return ret;
187. }
188.
189. /* release 函数实现, 对应到 Linux 系统调用函数的 close 函数 */
190. static int char_drv_release(struct inode *inode_p, struct file *file_p)

```

```
191. {
192.     printk("gpio_test module release\n");
193.     return 0;
194. }
195.
196. /* file_operations 结构体声明, 是上面 open、write 实现函数与系统调用函数对应的关键 */
197. static struct file_operations ax_char_fops = {
198.     .owner    = THIS_MODULE,
199.     .open     = char_drv_open,
200.     .read     = char_drv_read,
201.     .poll     = char_drv_poll,
202.     .release  = char_drv_release,
203. };
204.
205. /* 模块加载时会调用的函数 */
206. static int __init char_drv_init(void)
207. {
208.     /* 用于接受返回值 */
209.     u32 ret = 0;
210.
211.     /** 并发处理 **/
212.     /* 初始化原子变量 */
213.     atomic_set(&alinx_char.key_sts, 0);
214.
215.     /** gpio 框架 **/
216.     /* 获取设备节点 */
217.     alinx_char.nd = of_find_node_by_path("/alinxkey");
218.     if(alinx_char.nd == NULL)
219.     {
220.         printk("alinx_char node not find\r\n");
221.         return -EINVAL;
222.     }
223.     else
224.     {
225.         printk("alinx_char node find\r\n");
226.     }
227.
228.     /* 获取节点中 gpio 标号 */
229.     alinx_char.alinx_key_gpio = of_get_named_gpio(alinx_char.nd, "alinxkey-gpios", 0);
230.     if(alinx_char.alinx_key_gpio < 0)
231.     {
232.         printk("can not get alinxkey-gpios");
233.         return -EINVAL;
234.     }
```

```
235.     printk("alinxkey-gpio num = %d\r\n", alinx_char.alinx_key_gpio);
236.
237.     /* 申请 gpio 标号对应的引脚 */
238.     ret = gpio_request(alinx_char.alinx_key_gpio, "alinxkey");
239.     if(ret != 0)
240.     {
241.         printk("can not request gpio\r\n");
242.         return -EINVAL;
243.     }
244.
245.     /* 把这个 io 设置为输入 */
246.     ret = gpio_direction_input(alinx_char.alinx_key_gpio);
247.     if(ret < 0)
248.     {
249.         printk("can not set gpio\r\n");
250.         return -EINVAL;
251.     }
252.
253. /** 中断 **/
254.     /* 获取中断号 */
255.     alinx_char.irq = gpio_to_irq(alinx_char.alinx_key_gpio);
256.     /* 申请中断 */
257.     ret = request_irq(alinx_char.irq,
258.                      key_handler,
259.                      IRQF_TRIGGER_FALLING | IRQF_TRIGGER_RISING,
260.                      "alinxkey",
261.                      NULL);
262.     if(ret < 0)
263.     {
264.         printk("irq %d request failed\r\n", alinx_char.irq);
265.         return -EFAULT;
266.     }
267.
268. /** 定时器 **/
269.     alinx_char.timer.function = timer_function;
270.     init_timer(&alinx_char.timer);
271.
272. /** 等待队列 **/
273.     init_waitqueue_head(&alinx_char.wait_q_h);
274.
275. /** 字符设备框架 **/
276.     /* 注册设备号 */
277.     alloc_chrdev_region(&alinx_char.devid, MINOR_U, DEVID_COUNT, DEVICE_NAME);
278.
```

```

279.     /* 初始化字符设备结构体 */
280.     cdev_init(&alinx_char.cdev, &ax_char_fops);
281.
282.     /* 注册字符设备 */
283.     cdev_add(&alinx_char.cdev, alinx_char.devid, DRIVE_COUNT);
284.
285.     /* 创建类 */
286.     alinx_char.class = class_create(THIS_MODULE, DEVICE_NAME);
287.     if(IS_ERR(alinx_char.class))
288.     {
289.         return PTR_ERR(alinx_char.class);
290.     }
291.
292.     /* 创建设备节点 */
293.     alinx_char.device = device_create(alinx_char.class, NULL,
294.                                       alinx_char.devid, NULL,
295.                                       DEVICE_NAME);
296.     if (IS_ERR(alinx_char.device))
297.     {
298.         return PTR_ERR(alinx_char.device);
299.     }
300.
301.     return 0;
302. }
303.
304. /* 卸载模块 */
305. static void __exit char_drv_exit(void)
306. {
307.     /** gpio **/
308.     /* 释放 gpio */
309.     gpio_free(alinx_char.alinx_key_gpio);
310.
311.     /** 中断 **/
312.     /* 释放中断 */
313.     free_irq(alinx_char.irq, NULL);
314.
315.     /** 定时器 **/
316.     /* 删除定时器 */
317.     del_timer_sync(&alinx_char.timer);
318.
319.     /** 字符设备框架 **/
320.     /* 注销字符设备 */
321.     cdev_del(&alinx_char.cdev);
322.

```

```

323.     /* 注销设备号 */
324.     unregister_chrdev_region(alinx_char.devid, DEVID_COUNT);
325.
326.     /* 删除设备节点 */
327.     device_destroy(alinx_char.class, alinx_char.devid);
328.
329.     /* 删除类 */
330.     class_destroy(alinx_char.class);
331.
332.     printk("timer_led_dev_exit_ok\n");
333. }
334.
335. /* 标记加载、卸载函数 */
336. module_init(char_drv_init);
337. module_exit(char_drv_exit);
338.
339. /* 驱动描述信息 */
340. MODULE_AUTHOR("Alinx");
341. MODULE_ALIAS("alinx char");
342. MODULE_DESCRIPTION("NIO LED driver");
343. MODULE_VERSION("v1.0");
344. MODULE_LICENSE("GPL");

```

驱动代码在上一章的代码基础上，增加了 poll 实现，并稍微修改了 read 函数。

201 行在 file_operations 结构体中添加 poll 函数实现。

168 行实现 poll 函数，调用一下 poll_wait 函数，之后哦按段数据状态，如果数据准备好，就返回 POLLIN 状态标识。

120 行在 read 函数中稍作修改，先判断文件打开方式，如果是非阻塞的方式访问，就不去做队列相关的操作了，直接返回数据给用户，否则按照阻塞访问处理。

10.2.4 测试代码

测试代码在 6.4 节的基础上修改，新建 QT 工程名为“ax_nioled_test”，新建 main.c，输入下列代码：

```

1. #include "stdio.h"
2. #include "unistd.h"
3. #include "sys/types.h"
4. #include "sys/stat.h"
5. #include "fcntl.h"
6. #include "stdlib.h"
7. #include "string.h"
8. #include "poll.h"
9. #include "sys/select.h"

```

```
10. #include "sys/time.h"
11. #include "linux/ioctl.h"
12.
13. int main(int argc, char *argv[])
14. {
15.
16.     /* ret 获取返回值, fd 获取文件句柄 */
17.     int ret, fd, fd_l;
18.     /* 定义一个监视文件读变化的描述符合集 */
19.     fd_set readfds;
20.     /* 定义一个超时时间结构体 */
21.     struct timeval timeout;
22.
23.     char *filename, led_value = 0;
24.     unsigned int key_value;
25.
26.     if(argc != 2)
27.     {
28.         printf("Error Usage\r\n");
29.         return -1;
30.     }
31.
32.     filename = argv[1];
33.     /* 获取文件句柄, O_NONBLOCK 表示非阻塞访问 */
34.     fd = open(filename, O_RDWR | O_NONBLOCK);
35.     if(fd < 0)
36.     {
37.         printf("can not open file %s\r\n", filename);
38.         return -1;
39.     }
40.
41.     while(1)
42.     {
43.         /* 初始化描述符合集 */
44.         FD_ZERO(&readfds);
45.         /* 把文件句柄 fd 指向的文件添加到描述符 */
46.         FD_SET(fd, &readfds);
47.
48.         /* 超时时间初始化为 1.5 秒 */
49.         timeout.tv_sec = 1;
50.         timeout.tv_usec = 500000;
51.
52.         /* 调用 select, 注意第一个参数为 fd+1 */
53.         ret = select(fd + 1, &readfds, NULL, NULL, &timeout);
```



```

98.
99.             ret = close(fd_1);
100.            if(ret < 0)
101.            {
102.                printf("file /dev/gpio_leds close failed\r\n");
103.                break;
104.            }
105.        }
106.    }
107.    break;
108. }
109. }
110. }
111. close(fd);
112. return ret;
113. }

```

73 行的 `read` 函数开始，之后的代码与 6.4 节是一样的，通过判断 `key` 的状态，来改变 `led` 的状态。

在调用 `read` 之前，先调用 `select` 函数来检测数据状态，用法和 10.1.1 节中的用法相同，就不重复说明了。

10.2.5 运行测试

测试方式和现象和上一章一样，步骤如下：

```

mount -t nfs -o nolock 192.168.1.107:/home/alinx/work /mnt
cd /mnt
mkdir /tmp/qt
mount qt_lib.img /tmp/qt
cd /tmp/qt
source ./qt_env_set.sh
cd /mnt
insmod ./ax-concled-dev.ko
insmod ./ax-nio-dev.ko
cd ./build-ax_nioled_test-ZYNQ-Debug
./ax-key-test /dev/nio_led&
top

```

此外，可以尝试一下，把测试程中的超时时间改成 0 或者 `NULL` 来贯彻现象。