

Linux epoll模型详解及源码分析

原创 KiteRunner24 最后发布于2018-06-03 15:50:50 阅读数 12537 ☆ 收藏

一、epoll简介

epoll是当前在Linux下开发大规模并发网络程序的热门选择，epoll在Linux2.6内核中正式引入，和select相似，都是**IO多路复用 (IO multiplexing)**

按照man手册的说法，**epoll是为处理大批量句柄而做了改进的poll。**

Linux下有以下几个经典的服务器模型：

1、PPC模型和TPC模型

PPC (Process Per Connection) 模型和TPC (Thread Per Connection) 模型的设计思想类似，就是给每一个到来的连接都分配一个独立的进程或者对于这两种模型，其需要耗费较大的时间和空间资源。当管理连接数较多时，进程或线程的切换开销较大。因此，这类模型能接受的最大连接数都不会几百个左右。

2、select模型

对于select模型，其主要有以下几个特点：

- **最大并发数限制**：由于一个进程所打开的fd（文件描述符）是有限制的，由FD_SETSIZE设置，默认值是1024/2048，因此，select模型的最大并发数就受限了。
- **效率问题**：每次进行select调用都会线性扫描全部的fd集合。这样，效率就会呈现线性下降。
- **内核/用户空间内存拷贝问题**：select在解决将fd消息传递给用户空间时采用了**内存拷贝**的方式。这样，其处理效率不高。

3、poll模型

对于poll模型，其虽然解决了select最大并发数的限制，但依然没有解决掉select的效率问题和内存拷贝问题。

4、epoll模型

对比于其他模型，epoll做了如下改进：

支持一个进程打开较大数目的文件描述符（fd）

select模型对一个进程所打开的文件描述符是有一定限制的，其由FD_SETSIZE设置，默认为1024/2048。这对于那些需要支持上万个连接数目的高并发场景来说，显然太少了，这个时候，可以选择两种方案：一是可以选择修改FD_SETSIZE宏然后重新编译内核，不过这样做也会带来网络效率的下降；二是可以选择线程模型（传统的Apache方案），不过虽然Linux中创建线程的代价比较小，但仍然是不可忽视的，加上进程间数据同步远不及线程间同步的高效，所以线程模型并不是完美的方案。

但是，epoll则没有对描述符数目的限制，**它所支持的文件描述符上限是整个系统最大可以打开的文件数目**，例如，在1GB内存的机器上，这个限制大概在10万左右。

IO效率不会随文件描述符（fd）的增加而线性下降

传统的select/poll的一个致命弱点就是当你拥有一个很大的socket集合时，不过任一时间只有部分socket是活跃的，select/poll每次调用都会线性扫描整个集合，这将导致IO处理效率呈现线性下降。

但是，epoll不存在这个问题，它只会对活跃的socket进行操作，这是因为在内核实现中，epoll是根据每个fd上面的callback函数实现的。因此，只有当fd上有数据可读或可写时，才会主动去调用callback函数，其他idle状态socket则不会。在这一点上，epoll实现了一个伪AIO，其内部推动力在内核。

在一些benchmark中，如果所有的socket基本上都是活跃的，如高速LAN环境，epoll并不比select/poll效率高，相反，过多使用epoll_ctl，其效率反而下降。但是，一旦使用idle connections模拟WAN环境，epoll的效率就远在select/poll之上了。

使用mmap加速内核与用户空间的消息传递

无论是select，poll还是epoll，它们都需要内核把fd消息通知给用户空间。因此，如何避免不必要的内存拷贝就很重要了。对于该问题，epoll通过在内核与用户空间之间使用mmap同一块内存来实现。

内核微调

这一点其实不算epoll的优点了，而是整个Linux平台的优点，Linux赋予开发者微调内核的能力。比如，内核TCP/IP协议栈使用内存池管理sk_buff结构以在运行期间动态调整这个内存池大小（skb_head_pool）来提高性能，该参数可以通过使用 `echo xxxx > /proc/sys/net/core/hot_list_length` 来尝试使用最新的NAPI网卡驱动架构来处理数据包数量巨大但数据包本身很小的特殊场景。

二、epoll API

epoll只有 `epoll_create`、`epoll_ctl` 和 `epoll_wait` 这三个系统调用。其定义如下：

```
1 #include <sys/epoll.h>
2
3 int epoll_create(int size);
4
5 int epoll_ctl(int epfd, int op, int fd, struct epoll_event *event);
6
7 int epoll_wait(int epfd, struct epoll_event *events, int maxevents, int timeout);
```

1、epoll_create

```
1 #include <sys/epoll.h>
2
3 int epoll_create(int size);
```

可以调用`epoll_create`方法创建一个epoll的句柄。

需要注意的是，当创建好epoll句柄后，它就会占用一个fd值。在使用完epoll后，必须调用close函数进行关闭，否则可能导致fd被耗尽。

2、epoll_ctl

```
1 #include <sys/epoll.h>
2
3 int epoll_ctl(int epfd, int op, int fd, struct epoll_event *event);
```

epoll的事件注册函数，它不同于select是在监听事件时告诉内核要监听什么类型的事件，而是通过`epoll_ctl`注册要监听的事件类型。

第一个参数`epfd`：`epoll_create`函数的返回值。

第二个参数`events`：表示动作类型。有三个宏来表示：

- * `EPOLL_CTL_ADD`：注册新的fd到`epfd`中；
- * `EPOLL_CTL_MOD`：修改已经注册的fd的监听事件；
- * `EPOLL_CTL_DEL`：从`epfd`中删除一个fd。

第三个参数`fd`：需要监听的fd。

第四个参数`event`：告诉内核需要监听什么事件。

`struct epoll_event`结构如下所示：

```
1 // 保存触发事件的某个文件描述符相关的数据
2 typedef union epoll_data {
3     void *ptr;
4     int fd;
5     __uint32_t u32;
6     __uint64_t u64;
7 } epoll_data_t;
8
9 // 感兴趣的事件和被触发的事件
10 struct epoll_event {
11     __uint32_t events; // Epoll events
12     epoll_data_t data; // User data variable
13 };
```

如上所示，对于Epoll Events，其可以是以下几个宏的集合：

- `EPOLLIN`：表示对应的文件描述符可读（包括对端Socket）；
- `EPOLLOUT`：表示对应的文件描述符可写；

- EPOLLPRI：表示对应的文件描述符有紧急数据可读（带外数据）；
- EPOLLERR：表示对应的文件描述符发生错误；
- EPOLLHUP：表示对应的文件描述符被挂断；
- EPOLLET：将EPOLL设为边缘触发（Edge Triggered），这是相对于水平触发（Level Triggered）而言的。
- EPOLLONESHOT：只监听一次事件，当监听完这次事件之后，如果还需要继续监听这个socket，需要再次

3、epoll_wait

```
1 #include <sys/epoll.h>
2
3 int epoll_wait(int epfd, struct epoll_event *events, int maxevents, int timeout);
```

收集在epoll监控的事件中已经发生的事件。参数events是分配好的epoll_event结构体数组，epoll将会把发生的事件赋值到events数组中（events不针，内核只负责把数据赋值到这个event数组中，不会去帮助我们在用户态分配内存）。maxevents告诉内核这个events数组有多大，这个maxevents创建epoll_create时的size。参数timeout是超时时间（毫秒）。如果函数调用成功，则返回对应IO上已准备好的文件描述符数目，如果返回0则表示已

三、epoll工作模式

1. LT模式（Level Triggered，水平触发）

该模式是epoll的缺省工作模式，其同时支持阻塞和非阻塞socket。内核会告诉开发者一个文件描述符是否就绪，如果开发者不采取任何操作，内核仍

2. ET模式（Edge Triggered，边缘触发）

该模式是一种高速处理模式，当且仅当状态发生变化时才会获得通知。在该模式下，其假定开发者在接收到一次通知后，会完整地处理该事件，因此内这一事件。注意，缓冲区中还有未处理的数据不能说是状态变化，因此，在ET模式下，开发者如果只读取了一部分数据，其将再也得不到通知了。正确开发者自己确认读完了所有的字节（一直调用read/write直到出错EAGAIN为止）。

Nginx默认采用的就是ET（边缘触发）。

四、epoll高效性探讨

epoll的高效性主要体现在以下三个方面：

（1）select/poll每次调用都要传递所要监控的所有fd给select/poll系统调用，这意味着每次调用select/poll时都要将fd列表从用户空间拷贝到内核，因此，这会造成性能低效。对于epoll_wait，每次调用epoll_wait时，其不需要将fd列表传递给内核，epoll_ctl不需要每次都拷贝所有的fd列表，只需要作。因此，在调用epoll_create函数之后，内核已经在内核开始准备数据结构用于存放需要监控的fd了。其后，每次epoll_ctl只是对这个数据结构进行操作即可。

（2）内核使用slab机制，为epoll提供了快速的数据结构。在内核里，一切都是文件。因此，epoll向内核注册了一个文件系统，用于存储所有被监控epoll_create时，就会在这个虚拟的epoll文件系统中创建一个file节点。epoll在被内核初始化时，同时会分配出epoll自己的内核告诉cache区，用于有希望监控的fd。这些fd会以红黑树的形式保存在内核cache里，以支持快速查找、插入和删除。这个内核高速cache，就是建立连续的物理内存页，然后slab层，简单的说，就是物理上分配好想要的size的内存对象，每次使用时都使用空闲的已分配好的对象。

（3）当调用epoll_ctl往epfd注册百万个fd时，epoll_wait仍然能够快速返回，并有效地将发生的事件fd返回给用户。原因在于，当我们调用epoll_ctl除了帮我们在epoll文件系统新建file节点，同时在内核cache创建红黑树用于存储以后由epoll_ctl传入的fd外，还会再建立一个list链表，用于存储准备就绪的fd。当调用epoll_wait时，仅仅观察这个list链表中是否有数据即可。如果list链表中有数据，则返回这个链表中的所有元素；如果list链表中没有数据，则返回到timeout超时返回。所以，epoll_wait非常高效，而且，通常情况下，即使我们需要监控百万计的fd，但大多数情况下，一次也只返回少量准备就绪的fd。因此，每次调用epoll_wait，其仅需要从内核态复制少量的fd到用户空间而已。那么，这个准备就绪的list链表是怎么维护的呢？过程如下：当我们执行epoll_ctl除了把fd放入到epoll文件系统里file对象对应的红黑树之外，还会给内核中断处理程序注册一个回调函数，其告诉内核，如果这个fd的中断到了，就把就绪的list链表中。

如此，一棵红黑树、一张准备就绪的fd链表以及少量的内核cache，就帮我们解决了高并发下fd的处理问题。

总结一下：

- 执行epoll_create时，创建了红黑树和就绪list链表；
- 执行epoll_ctl时，如果增加fd，则检查在红黑树中是否存在，存在则立即返回，不存在则添加到红黑树中，然后向内核注册回调函数，用于当中断准备就绪的list链表中插入数据。
- 执行epoll_wait时立即返回准备就绪链表里的数据即可。

五、epoll源码分析

eventpoll_init过程:

```
1 static int __init eventpoll_init(void)
2 {
3     int error;
4
5     init_MUTEX(&epsem);
6
7     /* Initialize the structure used to perform safe poll wait head wake ups */
8     ep_poll_safewake_init(&psw);
9
10    /* Allocates slab cache used to allocate "struct epitem" items */
11    epi_cache = kmem_cache_create("eventpoll_epi", sizeof(struct epitem),
12        0, SLAB_HWCACHE_ALIGN|EPI_SLAB_DEBUG|SLAB_PANIC,
13        NULL, NULL);
14
15    /* Allocates slab cache used to allocate "struct eppoll_entry" */
16    pwq_cache = kmem_cache_create("eventpoll_pwq",
17        sizeof(struct eppoll_entry), 0,
18        EPI_SLAB_DEBUG|SLAB_PANIC, NULL, NULL);
19
20    /*
21     * Register the virtual file system that will be the source of inodes
22     * for the eventpoll files
23     */
24    error = register_filesystem(&eventpoll_fs_type);
25    if (error)
26        goto epanic;
27
28    /* Mount the above commented virtual file system */
29    eventpoll_mnt = kern_mount(&eventpoll_fs_type);
30    error = PTR_ERR(eventpoll_mnt);
31    if (IS_ERR(eventpoll_mnt))
32        goto epanic;
33
34    DNPRINTK(3, (KERN_INFO "[%p] eventpoll: successfully initialized.\n",
35        current));
36    return 0;
37
38    epanic:
39    panic("eventpoll_init() failed\n");
40 }
```

其中, epoll用slab分配器kmem_cache_create分配内存用于存放 struct epitem 和 struct eppoll_entry。

当向系统中添加一个fd时, 就会创建一个epitem结构体, 这是内核管理epoll的基本数据结构:

```
1 /*
2  * Each file descriptor added to the eventpoll interface will
3  * have an entry of this type linked to the hash.
4  */
5 struct epitem {
6     /* RB-Tree node used to link this structure to the eventpoll rb-tree */
7     struct rb_node rbn;    // 用于主结构管理的红黑树
8
9     /* List header used to link this structure to the eventpoll ready list */
10    struct list_head rdllink; // 事件就绪队列
11
12    /* The file descriptor information this item refers to */
13    struct epoll_filefd ffd; // 用于主结构中的链表
14
15    /* Number of active wait queue attached to poll operations */
16    int nwait; // 事件个数
17
18    /* List containing poll wait queues */
19    struct list_head pwqlist; // 双向链表, 保存着被监控文件的等待队列
20
21    /* The "container" of this item */
22 }
```

```

23 struct eventpoll *ep; // 该项属于哪个主结构体
24
25 /* The structure that describe the interested events and the source fd */
26 struct epoll_event event; // 注册的感兴趣的时间
27
28 /*
29 * Used to keep track of the usage count of the structure. This avoids
30 * that the structure will desappear from underneath our processing.
31 */
32 atomic_t usecnt;
33
34 /* List header used to link this item to the "struct file" items list */
35 struct list_head flink;
36
37 /* List header used to link the item to the transfer list */
38 struct list_head txlink;
39
40 /*
41 * This is used during the collection/transfer of events to userspace
42 * to pin items empty events set.
43 */
44 unsigned int revents;
};

```

对于每个epfd，其对应的数据结构为：

```

1 /*
2 * This structure is stored inside the "private_data" member of the file
3 * structure and rapresent the main data sructure for the eventpoll
4 * interface.
5 */
6 struct eventpoll {
7     /* Protect the this structure access */
8     rwlock_t lock;
9
10    /*
11     * This semaphore is used to ensure that files are not removed
12     * while epoll is using them. This is read-held during the event
13     * collection loop and it is write-held during the file cleanup
14     * path, the epoll file exit code and the ctl operations.
15     */
16    struct rw_semaphore sem;
17
18    /* Wait queue used by sys_epoll_wait() */
19    wait_queue_head_t wq;
20
21    /* Wait queue used by file->poll() */
22    wait_queue_head_t poll_wait;
23
24    /* List of ready file descriptors */
25    struct list_head rdllist; // 准备就绪的事件链表
26
27    /* RB-Tree root used to store monitored fd structs */
28    struct rb_root rbr; // 用于管理所有fd的红黑树（根节点）
29 };

```

eventpoll在epoll_create时创建：

```

1 /*
2 * It opens an eventpoll file descriptor by suggesting a storage of "size"
3 * file descriptors. The size parameter is just an hint about how to size
4 * data structures. It won't prevent the user to store more than "size"
5 * file descriptors inside the epoll interface. It is the kernel part of
6 * the userspace epoll_create(2).
7 */
8 asmlinkage long sys_epoll_create(int size)
9 {
10

```

```

11     int error, fd;
12     struct eventpoll *ep;
13     struct inode *inode;
14     struct file *file;
15
16     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_create(%d)\n",
17                 current, size));
18
19     /*
20     * Sanity check on the size parameter, and create the internal data
21     * structure ( "struct eventpoll" ).
22     */
23     error = -EINVAL;
24     if (size <= 0 || (error = ep_alloc(&ep)) != 0) // ep_alloc为eventpoll分配内存并初始化
25         goto eexit_1;
26
27     /*
28     * Creates all the items needed to setup an eventpoll file. That is,
29     * a file structure, and inode and a free file descriptor.
30     */
31     error = ep_getfd(&fd, &inode, &file, ep); // 创建于eventpoll相关的数据结构, 包括file、inode和fd等信息
32     if (error)
33         goto eexit_2;
34
35     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_create(%d) = %d\n",
36                 current, size, fd));
37
38     return fd;
39
40 eexit_2:
41     ep_free(ep);
42     kfree(ep);
43 eexit_1:
44     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_create(%d) = %d\n",
45                 current, size, error));
46     return error;
47 }

```

如上, 内核中维护了一棵红黑树, 大致结构如下:



下面是epoll_ctl函数过程:

```

1  /*
2  * The following function implements the controller interface for
3  * the eventpoll file that enables the insertion/removal/change of
4  * file descriptors inside the interest set. It represents
5  * the kernel part of the user space epoll_ctl(2).
6  */
7  asmLinkage long
8  sys_epoll_ctl(int epfd, int op, int fd, struct epoll_event __user *event)
9  {
10     int error;
11     struct file *file, *tfile;
12     struct eventpoll *ep;
13     struct epitem *epi;
14     struct epoll_event epds;
15
16     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_ctl(%d, %d, %d, %p)\n",
17                 current, epfd, op, fd, event));
18
19     error = -EFAULT;
20     if (ep_op_hash_event(op) &&
21         copy_from_user(&epds, event, sizeof(struct epoll_event)))
22         goto eexit_1;
23
24

```

```

25  /* Get the "struct file *" for the eventpoll file */
26  error = -EBADF;
27  file = fget(epfd); // 获取epfd对应的文件
28  if (!file)
29      goto eexit_1;
30
31  /* Get the "struct file *" for the target file */
32  tfile = fget(fd); // 获取fd对应的文件
33  if (!tfile)
34      goto eexit_2;
35
36  /* The target file descriptor must support poll */
37  error = -EPERM;
38  if (!tfile->f_op || !tfile->f_op->poll)
39      goto eexit_3;
40
41  /*
42  * We have to check that the file structure underneath the file descriptor
43  * the user passed to us _is_ an eventpoll file. And also we do not permit
44  * adding an epoll file descriptor inside itself.
45  */
46  error = -EINVAL;
47  if (file == tfile || !is_file_epoll(file))
48      goto eexit_3;
49
50
51  /*
52  * At this point it is safe to assume that the "private_data" contains
53  * our own data structure.
54  */
55  ep = file->private_data;
56
57  down_write(&ep->sem);
58
59  /* Try to lookup the file inside our hash table */
60  epi = ep_find(ep, tfile, fd); // 在哈希表中查询，防止重复添加
61
62  error = -EINVAL;
63  switch (op) {
64  case EPOLL_CTL_ADD: // 添加节点，调用ep_insert函数
65      if (!epi) {
66          epds.events |= POLLERR | POLLHUP;
67
68          error = ep_insert(ep, &epds, tfile, fd);
69      } else
70          error = -EEXIST;
71      break;
72  case EPOLL_CTL_DEL: // 删除节点，调用ep_remove函数
73      if (epi)
74          error = ep_remove(ep, epi);
75      else
76          error = -ENOENT;
77      break;
78  case EPOLL_CTL_MOD: // 修改节点，调用ep_modify函数
79      if (epi) {
80          epds.events |= POLLERR | POLLHUP;
81          error = ep_modify(ep, epi, &epds);
82      } else
83          error = -ENOENT;
84      break;
85  }
86
87  /*
88  * The function ep_find() increments the usage count of the structure
89  * so, if this is not NULL, we need to release it.
90  */
91  if (epi)
92      ep_release_epitem(epi);
93
94  up_write(&ep->sem);
95

```

```

96
97 eexit_3:
98     fput(tfile);
99 eexit_2:
100    fput(file);
101 eexit_1:
102    DNPRTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_ctl(%d, %d, %d, %p) = %d\n",
103              current, epfd, op, fd, event, error));

    return error;
}

```

对于ep_insert函数，基本代码如下：

```

1  static int ep_insert(struct eventpoll *ep, struct epoll_event *event,
2                      struct file *tfile, int fd)
3  {
4      int error, revents, pwake = 0;
5      unsigned long flags;
6      struct epitem *epi;
7      struct ep_queue epq;
8
9      error = -ENOMEM;
10     // 分配一个epitem结构体来保存每个加入的fd
11     if (!(epi = kmem_cache_alloc(epi_cache, SLAB_KERNEL)))
12         goto eexit_1;
13
14     /* Item initialization follow here ... */
15     // 初始化结构体
16     ep_rb_initnode(&epi->rbn);
17     INIT_LIST_HEAD(&epi->rdllink);
18     INIT_LIST_HEAD(&epi->flink);
19     INIT_LIST_HEAD(&epi->txlink);
20     INIT_LIST_HEAD(&epi->pwqlist);
21     epi->ep = ep;
22     ep_set_ffd(&epi->ffd, tfile, fd);
23     epi->event = *event;
24     atomic_set(&epi->usecnt, 1);
25     epi->nwait = 0;
26
27     /* Initialize the poll table using the queue callback */
28     epq.epi = epi;
29     // 安装poll回调函数
30     init_poll_funcptr(&epq.pt, ep_ptable_queue_proc);
31
32     /*
33     * Attach the item to the poll hooks and get current event bits.
34     * We can safely use the file* here because its usage count has
35     * been increased by the caller of this function.
36     */
37     // 将当前item添加至poll hook中，然后获取当前event位
38     revents = tfile->f_op->poll(tfile, &epq.pt);
39
40     /*
41     * We have to check if something went wrong during the poll wait queue
42     * install process. Namely an allocation for a wait queue failed due
43     * high memory pressure.
44     */
45     if (epi->nwait < 0)
46         goto eexit_2;
47
48     /* Add the current item to the list of active epoll hook for this file */
49     spin_lock(&tfile->f_ep_lock);
50     list_add_tail(&epi->flink, &tfile->f_ep_links);
51     spin_unlock(&tfile->f_ep_lock);
52
53     /* We have to drop the new item inside our item list to keep track of it */
54     write_lock_irqsave(&ep->lock, flags);
55
56

```



```

57  /* Add the current item to the rb-tree */
58  ep_rbtrees_insert(ep, epi);
59
60  /* If the file is already "ready" we drop it inside the ready list */
61  if ((revents & event->events) && !ep_is_linked(&epi->rdllink)) {
62      list_add_tail(&epi->rdllink, &ep->rdllist);
63
64      /* Notify waiting tasks that events are available */
65      if (waitqueue_active(&ep->wq))
66          wake_up(&ep->wq);
67      if (waitqueue_active(&ep->poll_wait))
68          pwake++;
69  }
70
71  write_unlock_irqrestore(&ep->lock, flags);
72
73  /* We have to call this outside the lock */
74  if (pwake)
75      ep_poll_safewake(&psw, &ep->poll_wait);
76
77  DNPRINTK(3, (KERN_INFO "[%p] eventpoll: ep_insert(%p, %p, %d)\n",
78              current, ep, tfile, fd));
79
80  return 0;
81
82 eexit_2:
83  ep_unregister_pollwait(ep, epi);
84
85  /*
86   * We need to do this because an event could have been arrived on some
87   * allocated wait queue.
88   */
89  write_lock_irqsave(&ep->lock, flags);
90  if (ep_is_linked(&epi->rdllink))
91      ep_list_del(&epi->rdllink);
92  write_unlock_irqrestore(&ep->lock, flags);
93
94  kmem_cache_free(epi_cache, epi);
95 eexit_1:
96  return error;
97 }

```

其中，`init_poll_funcptr` 和 `tfile->f_op->poll` 将 `ep_ptable_queue_proc` 注册到 `epq.pt` 中的 `qproc` 中。

`ep_ptable_queue_proc` 函数设置了等待队列的 `ep_poll_callback` 回调函数。在设备硬件数据到来时，硬件中断函数唤醒该等待队列上等待的进程时，会调用 `ep_poll_callback`。

`ep_poll_callback` 函数主要的功能是将被监视文件的等待事件就绪时，将文件对应的 `epitem` 实例添加到就绪队列中，当用户调用 `epoll_wait` 时，内核会的事件报告给用户。

`epoll_wait` 的实现如下：

```

1  /*
2   * Implement the event wait interface for the eventpoll file. It is the kernel
3   * part of the user space epoll_wait(2).
4   */
5  asmLinkage long sys_epoll_wait(int epfd, struct epoll_event __user *events,
6                                int maxevents, int timeout)
7  {
8      int error;
9      struct file *file;
10     struct eventpoll *ep;
11
12     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_wait(%d, %p, %d, %d)\n",
13                 current, epfd, events, maxevents, timeout));
14
15     /* The maximum number of event must be greater than zero */
16     if (maxevents <= 0 || maxevents > MAX_EVENTS) // 检查maxevents参数
17         return -EINVAL;

```

```

18
19  /* Verify that the area passed by the user is writeable */
20  // 检查用户空间传入的events指向的内存是否可写
21  if (!access_ok(VERIFY_WRITE, events, maxevents * sizeof(struct epoll_event))) {
22      error = -EFAULT;
23      goto eexit_1;
24  }
25
26  /* Get the "struct file *" for the eventpoll file */
27  error = -EBADF;
28  file = fget(epfd); // 获取epfd对应的eventpoll文件的file实例, file结构是在epoll_create中创建的
29  if (!file)
30      goto eexit_1;
31
32  /*
33   * We have to check that the file structure underneath the fd
34   * the user passed to us _is_ an eventpoll file.
35   */
36  error = -EINVAL;
37  if (!is_file_epoll(file))
38      goto eexit_2;
39
40  /*
41   * At this point it is safe to assume that the "private_data" contains
42   * our own data structure.
43   */
44  ep = file->private_data;
45
46  /* Time to fish for events ... */
47  // 核心处理函数
48  error = ep_poll(ep, events, maxevents, timeout);
49
50
51 eexit_2:
52     fput(file);
53 eexit_1:
54     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_wait(%d, %p, %d, %d) = %d\n",
55                 current, epfd, events, maxevents, timeout, error));
56
57     return error;
58 }

```

其中，调用ep_poll函数，具体流程如下：

```

1  static int ep_poll(struct eventpoll *ep, struct epoll_event __user *events,
2                    int maxevents, long timeout)
3  {
4      int res, eavail;
5      unsigned long flags;
6      long jtimeout;
7      wait_queue_t wait;
8
9      /*
10       * Calculate the timeout by checking for the "infinite" value (-1)
11       * and the overflow condition. The passed timeout is in milliseconds,
12       * that why (t * HZ) / 1000.
13       */
14     jtimeout = (timeout < 0 || timeout >= EP_MAX_MSTIMEO) ?
15         MAX_SCHEDULE_TIMEOUT : (timeout * HZ + 999) / 1000;
16
17     retry:
18         write_lock_irqsave(&ep->lock, flags);
19
20         res = 0;
21         if (list_empty(&ep->rddllist)) {
22             /*
23              * We don't have any available event to return to the caller.
24              * We need to sleep here, and we will be wake up by
25              * ep_poll_callback() when events will become available.

```

```

27     */
28     init_waitqueue_entry(&wait, current);
29     add_wait_queue(&ep->wq, &wait);
30
31     for (;;) {
32         /*
33          * We don't want to sleep if the ep_poll_callback() sends us
34          * a wakeup in between. That's why we set the task state
35          * to TASK_INTERRUPTIBLE before doing the checks.
36          */
37         set_current_state(TASK_INTERRUPTIBLE);
38         if (!list_empty(&ep->rdllist) || !jtimeout)
39             break;
40         if (signal_pending(current)) {
41             res = -EINTR;
42             break;
43         }
44
45         write_unlock_irqrestore(&ep->lock, flags);
46         jtimeout = schedule_timeout(jtimeout);
47         write_lock_irqsave(&ep->lock, flags);
48     }
49     remove_wait_queue(&ep->wq, &wait);
50
51     set_current_state(TASK_RUNNING);
52 }
53
54 /* Is it worth to try to dig for events ? */
55 eavail = !list_empty(&ep->rdllist);
56
57 write_unlock_irqrestore(&ep->lock, flags);
58
59 /*
60  * Try to transfer events to user space. In case we get 0 events and
61  * there's still timeout left over, we go trying again in search of
62  * more luck.
63  */
64
65 if (!res && eavail &&
66     !(res = ep_events_transfer(ep, events, maxevents)) && jtimeout)
67     goto retry;
68
69 return res;
70 }

```

ep_send_events函数用于向用户空间发送就绪事件。ep_send_events函数将用户传入的内存简单封装到ep_send_events_data结构中，然后调用ep_scan_ready_list将就绪队列中的事件传入用户空间的内存。

六、参考

Epoll详解及源码分析——CSDN博客

👍 点赞 13 ☆ 收藏 🔄 分享 ...



KiteRunner24

发布了32 篇原创文章 · 获赞 48 · 访问量 10万+

私信



排名前10的erp软件是什么

国内erp软件品牌排行

10/10/5



想对作者说点什么...



👤 LeBron James 👤 9个月前 good

