

作者：董昊（要转载的同学帮忙把名字和博客链接http://donghao.org/uii/带上，多谢了！）

## epoll原理简介

通过上面的分析，poll运行效率的两个瓶颈已经找出，现在的问题是怎么改进。首先，每次poll都要把1000个fd拷入内核，太不科学了，内核干嘛不自己保存已经拷入的fd呢？答对了，epoll就是自己保存拷入的fd，它的API就已经说明了这一点——不是 epoll\_wait的时候才传入fd，而是通过epoll\_ctl把所有fd传入内核再一起"wait"，这就省掉了不必要的重复拷贝。其次，在 epoll\_wait时，也不是把current轮流的加入fd对应的设备等待队列，而是在设备等待队列醒来时调用一个回调函数（当然，这就需要“唤醒回调”机制），把产生事件的fd归入一个链表，然后返回这个链表上的fd。

## epoll剖析

epoll是个module，所以先看看module的入口eventpoll\_init

```
[fs/eventpoll.c-->eventpoll_init()]
1582 static int __init eventpoll_init(void)
1583 {
1584     int error;
1585
1586     init_MUTEX(&epsem);
1587
1588     /* Initialize the structure used to perform safe poll wait head wake ups */
1589     ep_poll_safewake_init(&psw);
1590
1591     /* Allocates slab cache used to allocate "struct epitem" items */
1592     epi_cache = kmalloc_cache_create("eventpoll_epi", sizeof(struct epitem),
1593                                     0, SLAB_HWCACHE_ALIGN|EPI_SLAB_DEBUG|SLAB_PANIC,
1594                                     NULL, NULL);
1595
1596     /* Allocates slab cache used to allocate "struct eppoll_entry" */
1597     pwq_cache = kmalloc_cache_create("eventpoll_pwq",
1598                                     sizeof(struct eppoll_entry), 0,
1599                                     EPI_SLAB_DEBUG|SLAB_PANIC, NULL, NULL);
1600
1601     /*
1602      * Register the virtual file system that will be the source of inodes
1603      * for the eventpoll files
1604      */
1605     error = register_filesystem(&eventpoll_fs_type);
1606     if (error)
1607         goto epanic;
1608
1609     /* Mount the above commented virtual file system */
1610     eventpoll_mnt = kern_mount(&eventpoll_fs_type);
1611     error = PTR_ERR(eventpoll_mnt);
1612     if (IS_ERR(eventpoll_mnt))
1613         goto epanic;
1614
1615     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: successfully initialized.\n",
1616                 current));
1617
1618     return 0;
```

```
1619 epanic:  
1620     panic("eventpoll_init() failed\n");  
1621 }
```

很有趣，这个module在初始化时注册了一个新的文件系统，叫"eventpollfs"（在eventpoll\_fs\_type结构里），然后挂载此文件系统。另外创建两个内核cache（在内核编程中，如果需要频繁分配小块内存，应该创建kmem\_cache来做“内存池”），分别用于存放struct epitem和eppoll\_entry。如果以后要开发新的文件系统，可以参考这段代码。

现在想想epoll\_create为什么会返回一个新的fd？因为它就是在这个叫做"eventpollfs"的文件系统里创建了一个新文件！如下：

```
[fs/eventpoll.c-->sys_epoll_create()]  
476 asmlinkage long sys_epoll_create(int size)  
477 {  
478     int error, fd;  
479     struct inode *inode;  
480     struct file *file;  
481  
482     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_create(%d)\n",  
483                 current, size));  
484  
485     /* Sanity check on the size parameter */  
486     error = -EINVAL;  
487     if (size <= 0)  
488         goto eexit_1;  
489  
490     /*  
491      * Creates all the items needed to setup an eventpoll file. That is,  
492      * a file structure, and inode and a free file descriptor.  
493      */  
494     error = ep_getfd(&fd, &inode, &file);  
495     if (error)  
496         goto eexit_1;  
497  
498     /* Setup the file internal data structure ( "struct eventpoll" ) */  
499     error = ep_file_init(file);  
500     if (error)  
501         goto eexit_2;
```

函数很简单，其中ep\_getfd看上去是“get”，其实在第一次调用epoll\_create时，它是要创建新inode、新的file、新的fd。而ep\_file\_init则要创建一个struct eventpoll结构，并把它放入file->private\_data，注意，这个private\_data后面还要用到的。

看到这里，也许有人要问了，为什么epoll的开发者不做一个内核的超级大map把用户要创建的epoll句柄存起来，在epoll\_create时返回一个指针？那似乎很直观呀。但是，仔细看看，linux的系统调用有多少是返回指针的？你会发现几乎没有！（特此强调，malloc不是系统调用，malloc调用的brk才是）因为linux做为unix的最杰出的继承人，它遵循了unix的一个巨大优点——一切皆文件，输入输出是文件、socket也是文件，一切皆文件意味着使用这个操作系统的程序可以非常简单，因为一切都是文件操作而已！（unix还不是完全做到，plan 9才算）。而且使用文件系统有个好处：epoll\_create返回的是一个fd，而不是该死的指针，指针如果指错了，你简直没办法判断，而fd则可以通过current->files->fd\_array[]找到其真伪。

epoll\_create好了，该epoll\_ctl了，我们略去判断性的代码：

```
[fs/eventpoll.c-->sys_epoll_ctl()]  
524 asmlinkage long  
525 sys_epoll_ctl(int epfd, int op, int fd, struct epoll_event __user *event)  
526 {
```

```

527 int error;
528 struct file *file, *tfile;
529 struct eventpoll *ep;
530 struct epitom *epi;
531 struct epoll_event epds;
.....
575 epi = ep_find(ep, tfile, fd);
576
577 error = -EINVAL;
578 switch (op) {
579 case EPOLL_CTL_ADD:
580     if (!epi) {
581         epds.events |= POLLERR | POLLHUP;
582
583         error = ep_insert(ep, &epds, tfile, fd);
584     } else
585         error = -EEXIST;
586     break;
587 case EPOLL_CTL_DEL:
588     if (epi)
589         error = ep_remove(ep, epi);
590     else
591         error = -ENOENT;
592     break;
593 case EPOLL_CTL_MOD:
594     if (epi) {
595         epds.events |= POLLERR | POLLHUP;
596         error = ep_modify(ep, epi, &epds);
597     } else
598         error = -ENOENT;
599     break;
600 }

```

原来就是在在一个大的结构（现在先不管是什么大结构）里先ep\_find，如果找到了struct epitom而用户操作是ADD，那么返回-EEXIST；如果是DEL，则ep\_remove。如果找不到struct epitom而用户操作是ADD，就ep\_insert创建并插入一个。很直白。那这个“大结构”是什么呢？看ep\_find的调用方式，ep参数应该是指向这个“大结构”的指针，再看ep = file->private\_data，我们才明白，原来这个“大结构”就是那个在epoll\_create时创建的struct eventpoll，具体再看看ep\_find的实现，发现原来是struct eventpoll的rbr成员（struct rb\_root），原来这是一个红黑树的根！而红黑树上挂的都是struct epitom。

现在清楚了，**一个新创建的epoll文件带有一个struct eventpoll结构，这个结构上再挂一个红黑树，而这个红黑树就是每次epoll\_ctl时fd存放的地方！**

现在数据结构都已经清楚了，我们来看最核心的：

```

[fs/eventpoll.c-->sys_epoll_wait()]
627 asmlinkage long sys_epoll_wait(int epfd, struct epoll_event __user *events,
628                                     int maxevents, int timeout)
629 {
630     int error;
631     struct file *file;
632     struct eventpoll *ep;
633
634     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_wait(%d, %p, %d, %d)\n",
635                 current, epfd, events, maxevents, timeout));
636
637     /* The maximum number of event must be greater than zero */

```

```

638     if (maxevents <= 0)
639         return -EINVAL;
640
641     /* Verify that the area passed by the user is writeable */
642     if ((error = verify_area(VERIFY_WRITE, events, maxevents * sizeof(struct
epoll_event)))) {
643         goto eexit_1;
644
645     /* Get the "struct file *" for the eventpoll file */
646     error = -EBADF;
647     file = fget(epfd);
648     if (!file)
649         goto eexit_1;
650
651     /*
652      * We have to check that the file structure underneath the fd
653      * the user passed to us _is_ an eventpoll file.
654      */
655     error = -EINVAL;
656     if (!IS_FILE_EPOLL(file))
657         goto eexit_2;
658
659     /*
660      * At this point it is safe to assume that the "private_data" contains
661      * our own data structure.
662      */
663     ep = file->private_data;
664
665     /* Time to fish for events ... */
666     error = ep_poll(ep, events, maxevents, timeout);
667
668 eexit_2:
669     fput(file);
670 eexit_1:
671     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: sys_epoll_wait(%d, %p, %d, %d) =
%d\n",
672                 current, epfd, events, maxevents, timeout, error));
673
674     return error;
675 }

```

故伎重演，从file->private\_data中拿到struct eventpoll，再调用ep\_poll

[fs/eventpoll.c-->sys\_epoll\_wait()->ep\_poll()]

```

1468 static int ep_poll(struct eventpoll *ep, struct epoll_event __user *events,
1469                      int maxevents, long timeout)
1470 {
1471     int res, eavail;
1472     unsigned long flags;
1473     long jtimeout;
1474     wait_queue_t wait;
1475
1476     /*
1477      * Calculate the timeout by checking for the "infinite" value ( -1 )
1478      * and the overflow condition. The passed timeout is in milliseconds,

```

```

1479     * that why (t * HZ) / 1000.
1480     */
1481     jtimeout = timeout == -1 || timeout > (MAX_SCHEDULE_TIMEOUT - 1000) / HZ ?
1482         MAX_SCHEDULE_TIMEOUT: (timeout * HZ + 999) / 1000;
1483
1484 retry:
1485     write_lock_irqsave(&ep->lock, flags);
1486
1487     res = 0;
1488     if (list_empty(&ep->rdllist)) {
1489         /*
1490         * We don't have any available event to return to the caller.
1491         * We need to sleep here, and we will be wake up by
1492         * ep_poll_callback() when events will become available.
1493         */
1494         init_waitqueue_entry(&wait, current);
1495         add_wait_queue(&ep->wq, &wait);
1496
1497         for (;;) {
1498             /*
1499             * We don't want to sleep if the ep_poll_callback() sends us
1500             * a wakeup in between. That's why we set the task state
1501             * to TASK_INTERRUPTIBLE before doing the checks.
1502             */
1503             set_current_state(TASK_INTERRUPTIBLE);
1504             if (!list_empty(&ep->rdllist) || !jtimeout)
1505                 break;
1506             if (signal_pending(current)) {
1507                 res = -EINTR;
1508                 break;
1509             }
1510
1511             write_unlock_irqrestore(&ep->lock, flags);
1512             jtimeout = schedule_timeout(jtimeout);
1513             write_lock_irqsave(&ep->lock, flags);
1514         }
1515         remove_wait_queue(&ep->wq, &wait);
1516
1517         set_current_state(TASK_RUNNING);
1518     }

```

....  
又是一个大循环，不过这个大循环比poll的那个好，因为仔细一看——它居然除了睡觉和判断ep->rdllist是否为空以外，啥也没做！

什么也没做当然效率高了，但到底是谁来让ep->rdllist不为空呢？

答案是ep\_insert时设下的回调函数：

```
[fs/eventpoll.c-->sys_epoll_ctl()-->ep_insert()]
923 static int ep_insert(struct eventpoll *ep, struct epoll_event *event,
924                      struct file *tfile, int fd)
925 {
926     int error, revents, pwake = 0;
927     unsigned long flags;
928     struct epitem *epi;
929     struct ep_pqueue epq;
```

```

930
931     error = -ENOMEM;
932     if (!(epi = EPI_MEM_ALLOC())))
933         goto eexit_1;
934
935     /* Item initialization follow here ... */
936     EP_RB_INITNODE(&epi->rbn);
937     INIT_LIST_HEAD(&epi->rdllink);
938     INIT_LIST_HEAD(&epi->fllink);
939     INIT_LIST_HEAD(&epi->txlink);
940     INIT_LIST_HEAD(&epi->pwqlist);
941     epi->ep = ep;
942     EP_SET_FFD(&epi->ffd, tfile, fd);
943     epi->event = *event;
944     atomic_set(&epi->usecnt, 1);
945     epi->nwait = 0;
946
947     /* Initialize the poll table using the queue callback */
948     epq.epi = epi;
949     init_poll_funcptr(&epq.pt, ep_ptable_queue_proc);
950
951     /*
952      * Attach the item to the poll hooks and get current event bits.
953      * We can safely use the file* here because its usage count has
954      * been increased by the caller of this function.
955      */
956     revents = tfile->f_op->poll(tfile, &epq.pt);

```

我们注意949行，其实也就是

```
&(epq.pt)->qproc = ep_ptable_queue_proc;
```

紧接着 tfile->f\_op->poll(tfile, &epq.pt) 其实就是调用被监控文件 ( epoll里叫“target file” ) 的 poll 方法，而这个 poll 其实就是调用 poll\_wait ( 还记得 poll\_wait 吗？每个支持 poll 的设备驱动程序都要调用的 )，最后就是调用 ep\_ptable\_queue\_proc。这是比较难解的一个调用关系，因为不是语言级的直接调用。

ep\_insert 还把 struct epitem 放到 struct file 里的 f\_ep\_links 连表里，以方便查找，struct epitem 里的 fllink 就是担负这个使命的。

```
[fs/eventpoll.c-->ep_ptable_queue_proc()]
883 static void ep_ptable_queue_proc(struct file *file, wait_queue_head_t *whead,
884                                     poll_table *pt)
885 {
886     struct epitem *epi = EP_ITEM_FROM_EPQUEUE(pt);
887     struct epoll_entry *pwq;
888
889     if (epi->nwait >= 0 && (pwq = PWQ_MEM_ALLOC())) {
890         init_waitqueue_func_entry(&pwq->wait, ep_poll_callback);
891         pwq->whead = whead;
892         pwq->base = epi;
893         add_wait_queue(whead, &pwq->wait);
894         list_add_tail(&pwq->llink, &epi->pwqlist);
895         epi->nwait++;
896     } else {
897         /* We have to signal that an error occurred */
898         epi->nwait = -1;
899     }
}
```

```
900 }
```

上面的代码就是ep\_insert中要做的最重要的事：创建struct epoll\_entry，设置其唤醒回调函数为ep\_poll\_callback，然后加入设备等待队列（注意这里的whead就是上一章所说的每个设备驱动都要带的等待队列）。只有这样，当设备就绪，唤醒等待队列上的等待着时，ep\_poll\_callback就会被调用。**每次调用poll系统调用，操作系统都要把current（当前进程）挂到fd对应的所有设备的等待队列上，可以想象，fd多到上千的时候，这样“挂”法很费事；而每次调用epoll\_wait则没有这么罗嗦，epoll只在epoll\_ctl时把current挂一遍（这第一遍是免不了的）并给每个fd一个命令“好了就调回调函数”，如果设备有事件了，通过回调函数，会把fd放入rdllist，而每次调用epoll\_wait就只是收集rdllist里的fd就可以了——epoll巧妙的利用回调函数，实现了更高效的事件驱动模型。**

现在我们猜也能猜出来ep\_poll\_callback会干什么了——肯定是把红黑树上的收到event的epitem（代表每个fd）插入ep->rdllist中，这样，当epoll\_wait返回时，rdllist里就都是就绪的fd了！

[fs/eventpoll.c-->ep\_poll\_callback()]

```
1206 static int ep_poll_callback(wait_queue_t *wait, unsigned mode, int sync, void *key)
1207 {
1208     int pwake = 0;
1209     unsigned long flags;
1210     struct epitem *epi = EP_ITEM_FROM_WAIT(wait);
1211     struct eventpoll *ep = epi->ep;
1212
1213     DNPRINTK(3, (KERN_INFO "[%p] eventpoll: poll_callback(%p) epi=%p
ep=%p\n",
1214             current, epi->file, epi, ep));
1215
1216     write_lock_irqsave(&ep->lock, flags);
1217
1218     /*
1219      * If the event mask does not contain any poll(2) event, we consider the
1220      * descriptor to be disabled. This condition is likely the effect of the
1221      * EPOLLONESHOT bit that disables the descriptor when an event is received,
1222      * until the next EPOLL_CTL_MOD will be issued.
1223      */
1224     if (!(epi->event.events & ~EP_PRIVATE_BITS))
1225         goto is_disabled;
1226
1227     /* If this file is already in the ready list we exit soon */
1228     if (EP_IS_LINKED(&epi->rdllink))
1229         goto is_linked;
1230
1231     list_add_tail(&epi->rdllink, &ep->rdllist);
1232
1233 is_linked:
1234     /*
1235      * Wake up ( if active ) both the eventpoll wait list and the ->poll()
1236      * wait list.
1237      */
1238     if (waitqueue_active(&ep->wq))
1239         wake_up(&ep->wq);
1240     if (waitqueue_active(&ep->poll_wait))
1241         pwake++;
1242
1243 is_disabled:
1244     write_unlock_irqrestore(&ep->lock, flags);
1245
```

```

1246 /* We have to call this outside the lock */
1247 if (pwake)
1248     ep_poll_safewake(&psw, &ep->poll_wait);
1249
1250 return 1;
1251 }

```

真正重要的只有1231行的只一句，就是把struct epitem放到struct eventpoll的rdllist中去。现在我们可以画出epoll的核心数据结构图了：

