

ep\_insert()

a. 创建新的epitem对象, 并初始化成员:  
rdlink, pwqlist, flink指向自己,  
ep指向 eventpoll对象,  
ffd保存被监视的fd及file结构, event保存事件类型,

b. init\_poll\_funcptr(&epq.pt, ep\_ptable\_queue\_proc)  
设置 ep\_pqueue 成员 epi= epitem对象, pt 成员回调函数  
\_qproc= ep\_ptable\_queue\_proc  
ep\_queue的作用是能够在回调函数中取得对应的epitem对象

c. ep\_item\_poll()  
调用被监视fd的poll函数, 获取当前fd的事件状态位, 以便后续检查是否有用户关心的事件发生.  
如果fd是TCP套接字, poll 函数是 sock\_poll()且调用过程是 sock\_poll()-->tcp\_poll()-->sock\_poll\_wait()-->poll\_wait()-->b中的回调函数

d. list\_add\_tail\_rcu(&epi->flink, &tfile->f\_ep\_links)  
被监视 fd的file成员的f\_ep\_links链接到 epitem对象的flink

e. ep\_rbtrees\_insert(ep, epi)  
将fd的管理者 epitem对象插入到eventpoll对象的红黑树中, 以便 eventpoll 可以快速访问, 查找,插入,删除的时间复杂度均为  $O(\lg N)$ .  
如此, 每当有一个新的fd被监听时, 红黑树中都会增加一个结点

c1. ep\_ptable\_queue\_proc()  
取得epitem对象,  
创建新的 eppoll\_entry 对象

c1a. init\_waitqueue\_func\_entry()  
设置eppoll\_entry 对象的等待队列的回调函数 = ep\_poll\_callback,

c1b.  
设置eppoll\_entry 对象的等待队列的对头 = 被监听的fd的 等待队列的对头  
设置eppoll\_entry 对象的base=epitem对象

c1c.  
将eppoll\_entry 对象的wait添加到被监听的fd的等待队列上

c1d.  
将eppoll\_entry 对象的llink添加到epitem对象的pwqlist

c1e.  
epitem 对象的 nwait++, 显示 epitem 等待队列有一项

struct eventpoll

spinlock_t lock
struct mutex mtx
wait_queue_head_t wq
wait_queue_head_t poll_wait
struct list_head rdllist
struct rb_root rbr
struct epitem *ovflist
struct wakeup_source *ws
struct file *file
int visited
struct list_head visited_list_link

struct epitem

struct rb_node rbn
struct list_head rdllink
struct epitem *next
struct epoll_filefd ffd
int nwait
struct list_head pwqlist
struct eventpoll *ep
struct list_head flink
struct wakeup_source __rcu *ws
struct epoll_event event

1. 关联 epitem和 fd 上事件发生时callback 函数  
2. 成员 wait 挂载到fd的设备等待队列上

struct eppoll\_entry

struct list_head llink
struct epitem *base
wait_queue_head_t *whead
wait_queue_t wait

被监视fd的设备等待队列对头

wait\_queue\_head\_t

spinlock_t lock
struct list_head task_list

wait\_queue\_t

unsigned int flags
void *private
wait_queue_func_t func
struct list_head task_list

被监视fd对应的file结构

struct file

...
void *private_data
struct list_head f_ep_links
...

struct file \*[]

[0]
[1]
[listenfd]

struct epoll\_filefd

struct file *file
int fd

wait\_queue\_t

unsigned int flags	= 0
void *private	= current
wait_queue_func_t func	= default_wake_function
struct list_head task_list	

struct socket

socket_wq *wq
struct file *file
struct sock *sk
...

socket fd 对应的等待队列

在创建 socket时, 在函数 sock\_alloc\_file() 中, private\_data 指向 socket 对象

ep\_poll\_callback() 中将就绪的 epitem 链到 eventpoll的就绪链表

epoll\_wait() 中将当前进程链到 eventpoll的等待队列

被监视fd对应的file结构

socket fd 对应的等待队列