

nio 和 epoll

计算机网络io的底层

select和epoll

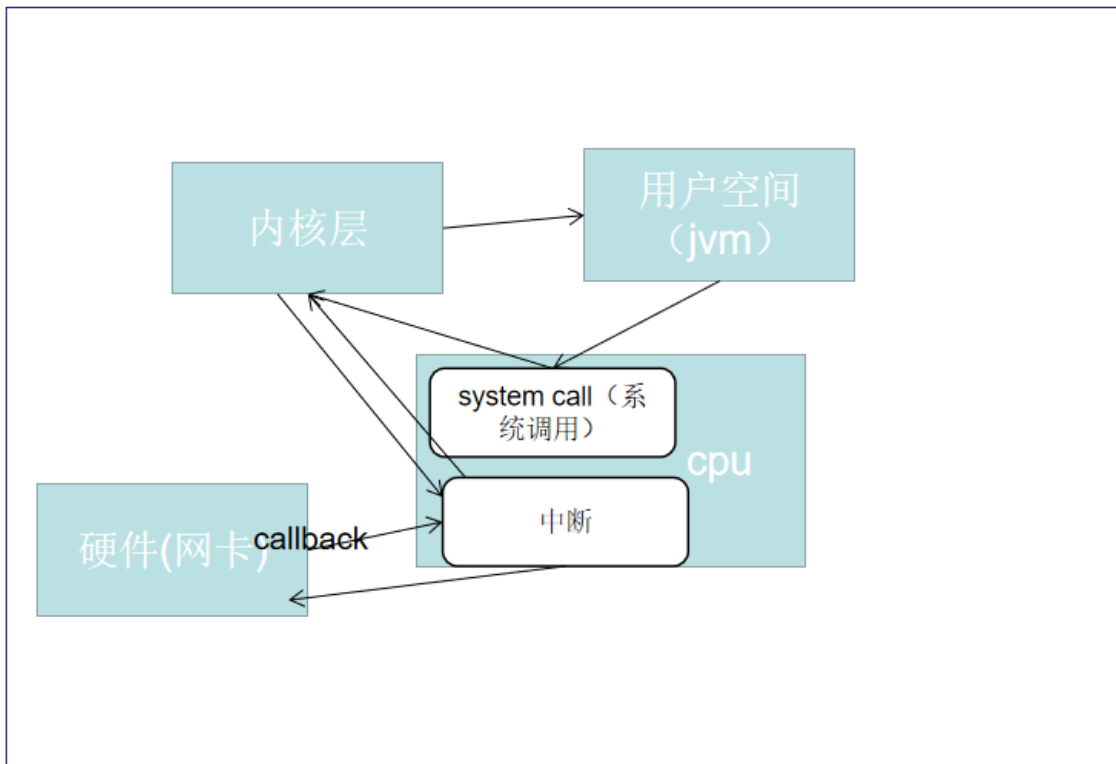
nio 多路复用



nio是java中同步非阻塞io，epoll 是linux 底层io 多路复用的一种技术，nio实现同步非阻塞是需要epoll底层支持，这篇文章我们就来搞懂epoll和nio的底层原理

计算机网络io的底层

在了解epoll的工作原理前我们需要先了解在linux系统中，一次网络io是怎么去完成的。我们来看下面这张图：



图中有两个需要理解的要点：

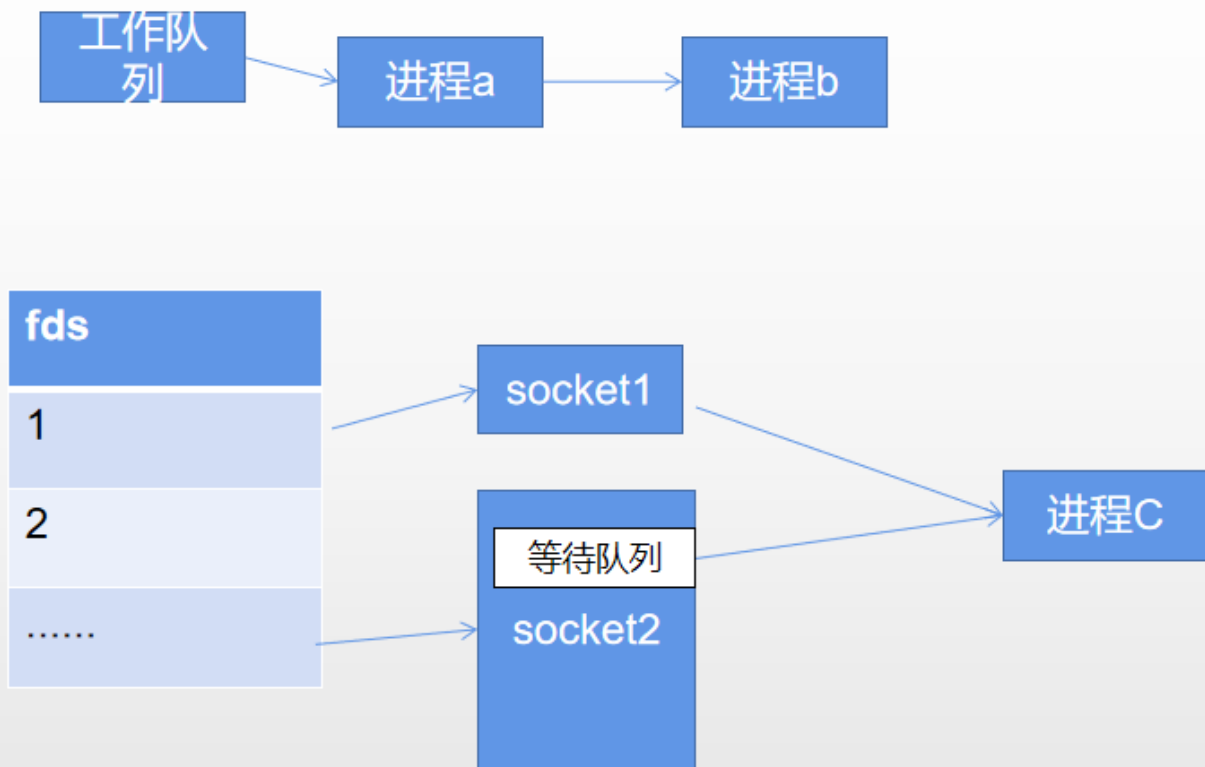
- 用户空间的程序（我们自己写的代码）是无法直接访问硬件层面的数据的，需要通过系统调用去将数据读取到内核空间再转到用户空间；
- 计算机执行程序时，会有优先级；一般而言，由硬件产生的信号需要cpu立马做出回应（不然数据可能就丢失），所以它的优先级很高。cpu理应中断掉正在执行的程序，去做出响应；当cpu完成对硬件的响应后，再重新执行用户程序。这种机制叫做中断；

当我们建立一个socket，用户程序通过系统调用向内核空间读取网卡的数据；网卡收到网络数据后，通过回调向cpu发出一个中断信号，操作系统便能得知有新数据到来，再通过中断程序去处理数据，然后将数据加载到内核空间，我们的程序再通过系统调用读取到网卡的数据。

上述的流程只是网络io的一个整体流程，我们想要明白为什么会有io多路复用这种机制出现还要进一步了解服务器建立socket的过程。当服务器建立一个socket，操作系统会创建一个由文件系统管理的socket对象。这个socket对象包含了发送缓冲区、接收缓冲区、等待队列等成员，等待队列存放所有需要等待该socket事件的进程。其伪代码如下：

```
1 //创建socket
2 int s = socket(AF_INET, SOCK_STREAM, 0);
3 //绑定
4 bind(s, ...)
5 //监听
6 listen(s, ...)
7 //接受客户端连接
8 int c = accept(s, ...)
9 //等待接收客户端数据
10 recv(c, ...);
11 //将数据打印出来
12 printf(...)
```

recv是个阻塞方法，当程序运行到recv时，它会一直等待，直到接收到数据才往下执行。当recv阻塞，是不会消耗cpu资源的，而我们网络io往往会有很多个socket，那么有没有什么办法，用一个进程去等待多个recv，也就是io多路复用？如下图所示，当进程c执行到recv方法等待网络数据时会被放入等待队列，cpu不会去执行进程c，也就是进程c不会占用cpu资源，当网卡有数据到来的时候网卡就会通过中断来将进程A唤醒，并挂到工作队列中让cpu运行它。



图中的fd是文件描述符，一个socket就是一个文件，socket句柄就是一个文件描述符。

那这个怎么和传统的BIO联系起来呢，我们看bio建立socket的代码：

```
1  {
2      ExecutorService executor = Executors.newFixedThreadPool(100); //线程池
3
4      ServerSocket serverSocket = new ServerSocket();
5      serverSocket.bind(8088);
6      while(!Thread.currentThread().isInterrupted()){ //主线程死循环等待新连接到来
7          Socket socket = serverSocket.accept();
8          executor.submit(new ConnectIOHandler(socket)); //为新的连接创建新的线程
9      }
10
11     class ConnectIOHandler extends Thread{
12         private Socket socket;
13         public ConnectIOHandler(Socket socket){
14             this.socket = socket;
15         }
16         public void run(){
17             while(!Thread.currentThread().isInterrupted() && !socket.isClosed()){ //死循环处理读写事件
18                 String something = socket.read().... //读取数据
19                 if(something != null){
20                     .... //处理数据
21                     socket.write().... //写数据
22                 }
23             }
24         }
25     }
```

传统的BIO里面socket.read()，如果TCP RecvBuffer里没有数据，函数会一直阻塞，直到收到数据，返回读到的数据。也就是说图中的进程如果进入了等待队列（recv方法阻塞），socket.read()也会跟着阻塞，直到有数据过来，它相当于直接调用操作系统recv方法。

select和epoll

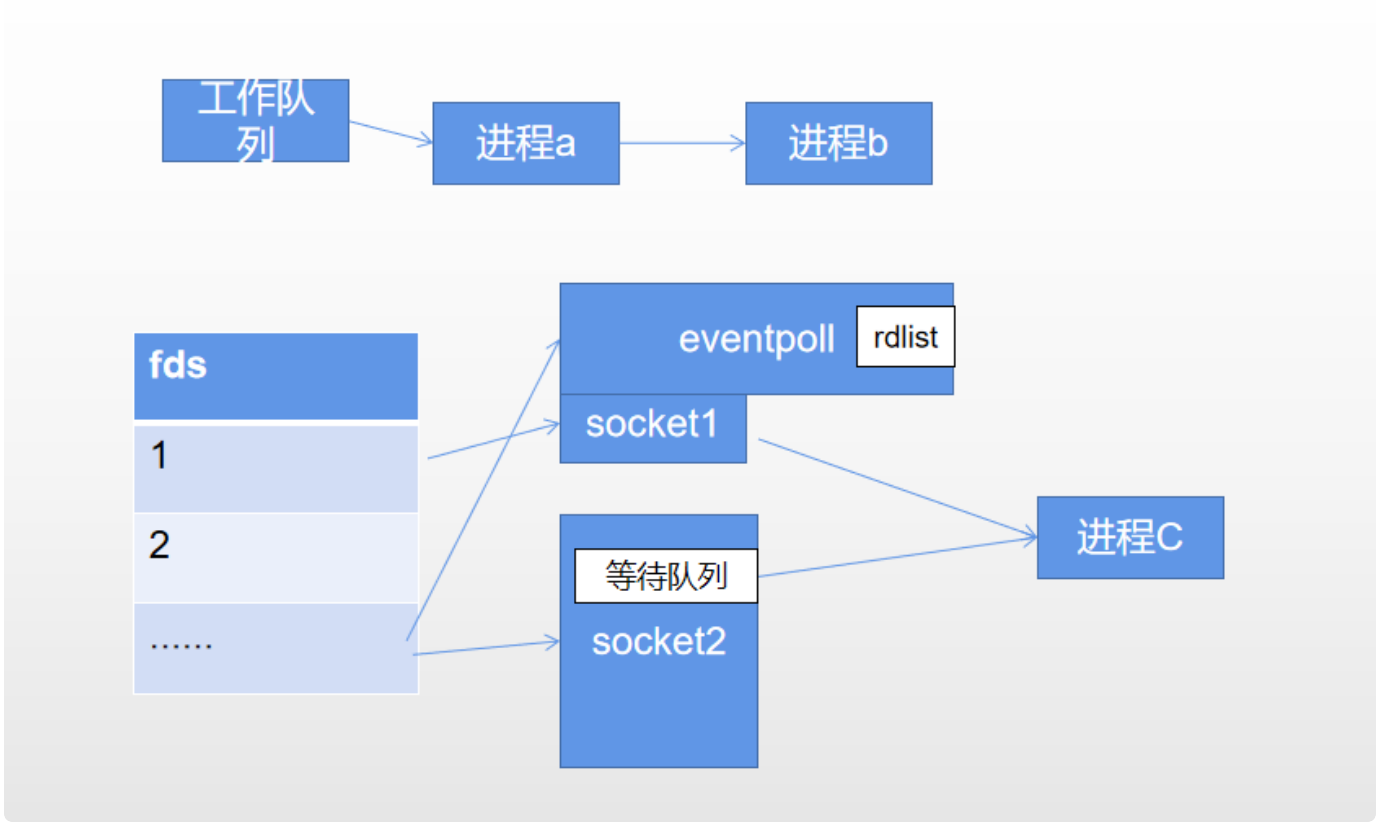
图中的机制还存在一个问题，前文提到过如果socket中有数据过来，网卡会给一个中断给到cpu，那么中断程序又是怎么去找到对应的fd的，这里面的工作流程是怎么样的？我们来看select是怎么做的：

数组fds存放着所有需要监视的socket，然后调用select，如果fds中的所有socket都没有数据，select会阻塞；当有一个socket接收到数据，发起一个中断，select返回，唤醒进程，将进程从所

有socket的等待队列中移除，遍历fds，通过FD_ISSET判断具体哪个socket收到数据，然后做出处理。该方法实现了一个进程监听多个socket，而不是傻傻的等recv返回（就像bio那样），但是这种方式有致命的缺点：

- 每次调用select都需要将进程加入到所有监视socket的等待队列，每次唤醒都需要从每个队列中移除。这里涉及了两次遍历，而且每次都要将整个fds列表传递给内核，有一定的开销。
- 进程被唤醒后，程序并不知道哪些socket收到数据，还需要遍历一次

为了解决上述缺点，于是有了epoll。在epoll中，内核维护一个“就绪列表”，引用收到数据的socket，避免了进程被唤醒后需要遍历才能找到socket，于是原来的模型就被改成了这样：



eventpoll 是 epoll 创建的对象，它内部维护了“就绪列表”；内核会将eventpoll添加到socket的等待队列中，当socket收到数据后，中断程序会操作eventpoll对象，而不是直接操作进程。当有socket就绪，发起中断，中断会将对应的socket引用存入rdlist中，并唤醒进程，进程只需要找rdlist就能快速处理就绪的socket；另：eventpoll通过红黑树保存监视的socket，而rdlist也并非直接引用socket，而是通过epitem间接引用，红黑树的节点也是epitem对象，图中简化了模型。

nio 多路复用

既然操作系统提供给我们这么好的机制，自然不能白白浪费，我们应用程序也想一个线程监听多个socket，谁有数据就起一个线程让它接收数据干活，干完活就扔到线程池里面去；为了更好的设计一套

多路复用的java框架，我们需要知道epoll的api:

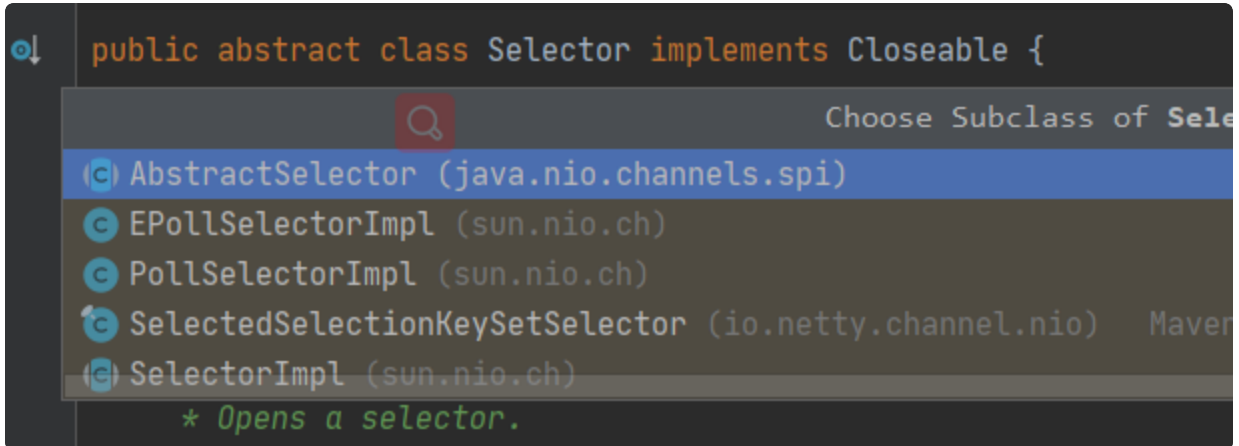
- `epoll_create`: 当某个进程调用`epoll_create`方法时，内核会创建一个eventpoll对象
- `epoll_ctl`: `epoll` 注册并监听事件的函数；
- `epoll_wait`: 等待文件描述符epfd上的事件。

我们来看NIO是如何进行IO的：

```
1  public static void main(String[] args) throws IOException, InterruptedExc
    eption {
2
3      ServerSocketChannel serverSocket = ServerSocketChannel.open();
4      serverSocket.socket().bind(new InetSocketAddress(8080));
5      serverSocket.configureBlocking(false);
6      Selector selector = Selector.open();
7      serverSocket.register(selector, SelectionKey.OP_ACCEPT);
8      System.out.println("服务启动成功");
9
10     while (true) {
11         selector.select();
12         Set<SelectionKey> selectionKeys = selector.selectedKeys();
13         Iterator<SelectionKey> iterator = selectionKeys.iterator();
14         // 遍历SelectionKey对事件进行处理
15         while (iterator.hasNext()) {
16             SelectionKey key = iterator.next();
17             iterator.remove();
18             // 如果是OP_ACCEPT事件, 则进行连接获取和事件注册
19             if (key.isAcceptable()) {
20                 ServerSocketChannel server = (ServerSocketChannel) ke
y.channel();
21                 SocketChannel socketChannel = server.accept();
22                 socketChannel.configureBlocking(false);
23                 // 这里只注册了读事件, 如果需要给客户端发送数据可以注册写事件
24                 socketChannel.register(selector, SelectionKey.OP_REA
D);
25                 System.out.println("客户端连接成功");
26             }
27             // 如果是OP_READ事件, 则进行读取和打印
28             if (key.isReadable()) {
29                 SocketChannel socketChannel = (SocketChannel) key.chan
nel();
30                 ByteBuffer byteBuffer = ByteBuffer.allocate(128);
31                 int read = socketChannel.read(byteBuffer);
32                 // 如果有数据, 把数据打印出来
33                 if (read > 0) {
34                     System.out.println("接收到消息: " + new String(byteB
uffer.array()));
35                 } else if (read == -1) {
36                     // 如果客户端断开连接, 关闭Socket
37                     System.out.println("客户端断开连接");
38                     socketChannel.close();
39                 }
40             }
```

```
41         }
42     }
43 }
44 }
45 }
```

然后查看一下`Selector`类linux系统下的的实现



它包括了 `PollSelectorImpl` , `EPollSelectorImpl` , `SelectedSelectionKeySetSelector` 实现类, 我们知道io多路复用在linux 中有Select, Poll Epoll 三种实现, 自然就对应了三个类, 我们重点关注 `EPollSelectorImpl` 的源码, 在NIO使用示例中 `Selector` 是通过 `Selector.open()` 创建的, 我们重点看看如何创建的这个 `Selector` :


```
1 public class EPollSelectorProvider
2     extends SelectorProviderImpl
3 {
4     public AbstractSelector openSelector() throws IOException {
5         return new EPollSelectorImpl(this);
6     }
7
8     public Channel inheritedChannel() throws IOException {
9         return InheritedChannel.getChannel();
10    }
11 }
12
13
14 EPollSelectorImpl(SelectorProvider sp) throws IOException {
15     super(sp);
16
17     this.epfd = EPoll.create();
18     this.pollArrayAddress = EPoll.allocatePollArray(NUM_EPOLLEVENTS);
19
20     try {
21         this.eventfd = new EventFD();
22         IOUtil.configureBlocking(IOUtil.newFD(eventfd.efd()), false);
23     } catch (IOException ioe) {
24         EPoll.freePollArray(pollArrayAddress);
25         FileDispatcherImpl.closeIntFD(epfd);
26         throw ioe;
27     }
28
29     // register the eventfd object for wakeups
30     EPollctl(epfd, EPOLL_CTL_ADD, eventfd.efd(), EPOLLIN);
31 }
```

我们看到在 `EPollSelectorImpl` 构造器中调用了 `EPoll.create()` 这个便是epoll的api。我们也可以顺便看看`EPoll`类的源码：

```
1 class EPoll {
2     private EPoll() { }
3     static native int create() throws IOException;
4
5     static native int ctl(int epfd, int opcode, int fd, int events);
6
7     static native int wait(int epfd, long pollAddress, int numfds, int time
    out)
8         throws IOException;
9 }
```

通过idea 可以看到 `sun.nio.ch.EPollSelectorImpl#doSelect` 调用了 `EPoll.wait` :

```
@Override
protected int doSelect(Consumer<SelectionKey> action, long timeout)
    throws IOException
{
    assert Thread.holdsLock( obj: this);

    // epoll_wait timeout is int
    int to = (int) Math.min(timeout, Integer.MAX_VALUE);
    boolean blocking = (to != 0);
    boolean timedPoll = (to > 0);

    int numEntries;
    processUpdateQueue();
    processDeregisterQueue();
    try {
        begin(blocking);

        do {
            long startTime = timedPoll ? System.nanoTime() : 0;
            long comp = BlockerImpl.isRead-only(g);
            try {
                numEntries = EPoll.wait(epfd, pollArrayAddress, NUM_EPOLLEVENTS, to);
            } finally {
                Blocker.end(comp);
            }
        } while (numEntries == 0);
    }
}
```

`EPoll.ctl` 在多处调用:


```

/**
 * Selects a set of keys whose corresponding channels are ready for I/O
 * operations.
 *
 * <p> This method performs a blocking <a href="#selop">selection
 * operation</a>. It returns only after at least one channel is selected,
 * this selector's {@link #wakeup wakeup} method is invoked, or the current
 * thread is interrupted, whichever comes first. </p>
 *
 * @return The number of keys, possibly zero, whose ready-operation sets
 *         now indicate readiness for at least one category of operations
 *         for which the channel was not previously detected to be ready
 *
 * @throws IOException
 *         If an I/O error occurs
 *
 * @throws ClosedSelectorException
 *         If this selector is closed
 */
public abstract int select() throws IOException;

```

大概意思就是选择一个准备好的通道。

通过上述的源码，我们发现NIO其实就是对linux多路复用的一个封装。离上一次写公众号已经很久了，懒了懒了。