

Java NIO浅析

王烨 · 2016-11-04 18:26

NIO (Non-blocking I/O , 在Java领域, 也称为New I/O) , 是一种同步非阻塞的I/O模型, 也是I/O多路复用的基础, 已经被越来越多地应用到大型应用服务器, 成为解决高并发与大量连接、I/O处理问题的有效方式。

那么NIO的本质是什么样的呢? 它是怎样与事件模型结合来解放线程、提高系统吞吐的呢?

本文会从传统的阻塞I/O和线程池模型面临的问题讲起, 然后对比几种常见I/O模型, 一步步分析NIO怎么利用事件模型处理I/O, 解决线程池瓶颈处理海量连接, 包括利用面向事件的方式编写服务端/客户端程序。最后延展到一些高级主题, 如Reactor与Proactor模型的对比、Selector的唤醒、Buffer的选择等。

注: 本文的代码都是伪代码, 主要是为了示意, 不可用于生产环境。

传统BIO模型分析

让我们先回忆一下传统的服务器端同步阻塞I/O处理 (也就是BIO , Blocking I/O) 的经典编程模型:

```
{
    ExecutorService executor = Executors.newFixedThreadPool(100); //线程池

    ServerSocket serverSocket = new ServerSocket();
    serverSocket.bind(8088);
    while(!Thread.currentThread().isInterrupted()) { //主线程死循环等待新连接到来
        Socket socket = serverSocket.accept();
        executor.submit(new ConnectIOHandler(socket)); //为新的连接创建新的线程
    }

    class ConnectIOHandler extends Thread{
        private Socket socket;
        public ConnectIOHandler(Socket socket){
            this.socket = socket;
        }
        public void run(){
            while(!Thread.currentThread().isInterrupted() && !socket.isClosed()) { //死循环处理读写事件
                String something = socket.read().... //读取数据
                if(something != null){
                    .... //处理数据
                    socket.write().... //写数据
                }
            }
        }
    }
}
```

这是一个经典的每连接每线程的模型，之所以使用多线程，主要原因在于`socket.accept()`、`socket.read()`、`socket.write()`三个主要函数都是同步阻塞的，当一个连接在处理I/O的时候，系统是阻塞的，如果是单线程的话必然就挂死在那里；但CPU是被释放出来的，开启多线程，就可以让CPU去处理更多的事情。其实这也是所有使用多线程的本质：

1. 利用多核。
2. 当I/O阻塞系统，但CPU空闲的时候，可以利用多线程使用CPU资源。

现在的多线程一般都使用线程池，可以让线程的创建和回收成本相对较低。在活动连接数不是特别高（小于单机1000）的情况下，这种模型是比较不错的，可以让每一个连接专注于自己的I/O并且编程模型简单，也不用过多考虑系统的过载、限流等问题。线程池本身就是一个天然的漏斗，可以缓冲一些系统处理不了的连接或请求。

不过，这个模型最本质的问题在于，严重依赖于线程。但线程是很“贵”的资源，主要表现在：

1. 线程的创建和销毁成本很高，在Linux这样的操作系统中，线程本质上就是一个进程。创建和销毁都是重量

级的系统函数。

2. 线程本身占用较大内存，像Java的线程栈，一般至少分配512K ~ 1M的空间，如果系统中的线程数过千，恐怕整个JVM的内存都会被吃掉一半。
3. 线程的切换成本是很高的。操作系统发生线程切换的时候，需要保留线程的上下文，然后执行系统调用。如果线程数过高，可能执行线程切换的时间甚至会大于线程执行的时间，这时候带来的表现往往是系统load偏高、CPU使用率特别高（超过20%以上），导致系统几乎陷入不可用的状态。
4. 容易造成锯齿状的系统负载。因为系统负载是用活动线程数或CPU核心数，一旦线程数量高但外部网络环境不是很稳定，就很容易造成大量请求的结果同时返回，激活大量阻塞线程从而使系统负载压力过大。

所以，当面对十万甚至百万级连接的时候，传统的BIO模型是无能为力的。随着移动端应用的兴起和各种网络游戏的盛行，百万级长连接日趋普遍，此时，必然需要一种更高效的I/O处理模型。

NIO是怎么工作的

很多刚接触NIO的人，第一眼看到的就是Java相对晦涩的API，比如：Channel，Selector，Socket什么的；然后就是一坨上百行的代码来演示NIO的服务端Demo.....瞬间头大有没有？

我们不管这些，抛开现象看本质，先分析下NIO是怎么工作的。

常见I/O模型对比

所有的系统I/O都分为两个阶段：等待就绪和操作。举例来说，读函数，分为等待系统可读和真正的读；同理，写函数分为等待网卡可以写和真正的写。

需要说明的是等待就绪的阻塞是不使用CPU的，是在“空等”；而真正的读写操作的阻塞是使用CPU的，真正在“干活”，而且这个过程非常快，属于memory copy，带宽通常在1GB/s级别以上，可以理解为基本不耗时。

下图是几种常见I/O模型的对比：



以socket.read()为例子：

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。

对于NIO，如果TCP RecvBuffer有数据，就把数据从网卡读到内存，并且返回给用户；反之则直接返回0，永远不会阻塞。

最新的AIO(Async I/O)里面会更进一步：不但等待就绪是非阻塞的，就连数据从网卡到内存的过程也是异步的。

换句话说，BIO里用户最关心“我要读”，NIO里用户最关心“我可以读了”，在AIO模型里用户更需要关注的是“读完了”。

NIO一个重要的特点是：socket主要的读、写、注册和接收函数，在等待就绪阶段都是非阻塞的，真正的I/O操作是同步阻塞的（消耗CPU但性能非常高）。

如何结合事件模型使用NIO同步非阻塞特性

回忆BIO模型，之所以需要多线程，是因为在进行I/O操作的时候，一是没有办法知道到底能不能写、能不能读，只能“傻等”，即使通过各种估算，算出来操作系统没有能力进行读写，也没法在socket.read()和socket.write()函数中返回。这两个函数无法进行有效的中断。所以除了多线程

`socket.read()`和`socket.write()`函数不返回，这两个函数无法进行有效的中断。所以除了多线程性另起炉灶，没有好的办法利用CPU。

NIO的读写函数可以立刻返回，这就给了我们不开线程利用CPU的最好机会：如果一个连接不能读写（`socket.read()`返回0或者`socket.write()`返回0），我们可以把这件事记下来，记录的方式通常是在Selector上注册标记位，然后切换到其它就绪的连接（channel）继续进行读写。

下面具体看下如何利用事件模型单线程处理所有I/O请求：

NIO的主要事件有几个：读就绪、写就绪、有新连接到来。

我们首先需要注册当这几个事件到来的时候所对应的处理器。然后在合适的时机告诉事件选择器：我对这个事件感兴趣。对于写操作，就是写不出去的时候对写事件感兴趣；对于读操作，就是完成连接和系统没有办法承载新读入的数据的时；对于accept，一般是服务器刚启动的时候；而对于connect，一般是connect失败需要重连或者直接异步调用connect的时候。

其次，用一个死循环选择就绪的事件，会执行系统调用（Linux 2.6之前是select、poll，2.6之后是epoll，Windows是IOCP），还会阻塞的等待新事件的到来。新事件到来的时候，会在selector上注册标记位，标示可读、可写或者有连接到来。

注意，select是阻塞的，无论是通过操作系统的通知（epoll）还是不停的轮询(select，poll)，这个函数是阻塞的。所以你可以放心大胆地在一个while(true)里面调用这个函数而不用担心CPU空转。

所以我们的程序大概的模样是：

```

interface ChannelHandler{
    void channelReadable(Channel channel);
    void channelWritable(Channel channel);
}

class Channel{
    Socket socket;
    Event event;//读，写或者连接
}

//IO线程主循环：
class IoThread extends Thread{
    public void run(){
        Channel channel;
        while(channel=Selector.select()){//选择就绪的事件和对应的连接
            if(channel.event==accept){
                registerNewChannelHandler(channel);//如果是新连接，则注册一个新的读写处理器
            }
            if(channel.event==write){
                getChannelHandler(channel).channelWritable(channel);//如果可以写，则执行写事件
            }
            if(channel.event==read){
                getChannelHandler(channel).channelReadable(channel);//如果可以读，则执行读事件
            }
        }
    }
}

Map<Channel, ChannelHandler> handlerMap;//所有channel的对应事件处理器
}

```

这个程序很简短，也是最简单的Reactor模式：注册所有感兴趣的事件处理器，单线程轮询选择就绪事件，执行事件处理器。

优化线程模型

由上面的示例我们大概可以总结出NIO是怎么解决掉线程的瓶颈并处理海量连接的：

NIO由原来的阻塞读写（占用线程）变成了单线程轮询事件，找到可以进行读写的网络描述符进行读写。除了事件的轮询是阻塞的（没有可干的事情必须要阻塞），剩余的I/O操作都是纯CPU操作，没有必要开启多线程。

并且由于线程的节约，连接数大的时候因为线程切换带来的问题也随之解决，进而为处理海量连接提供了可能。

单线程处理I/O的效率确实非常高，没有线程切换，只是拼命的读、写、选择事件。但现在的服

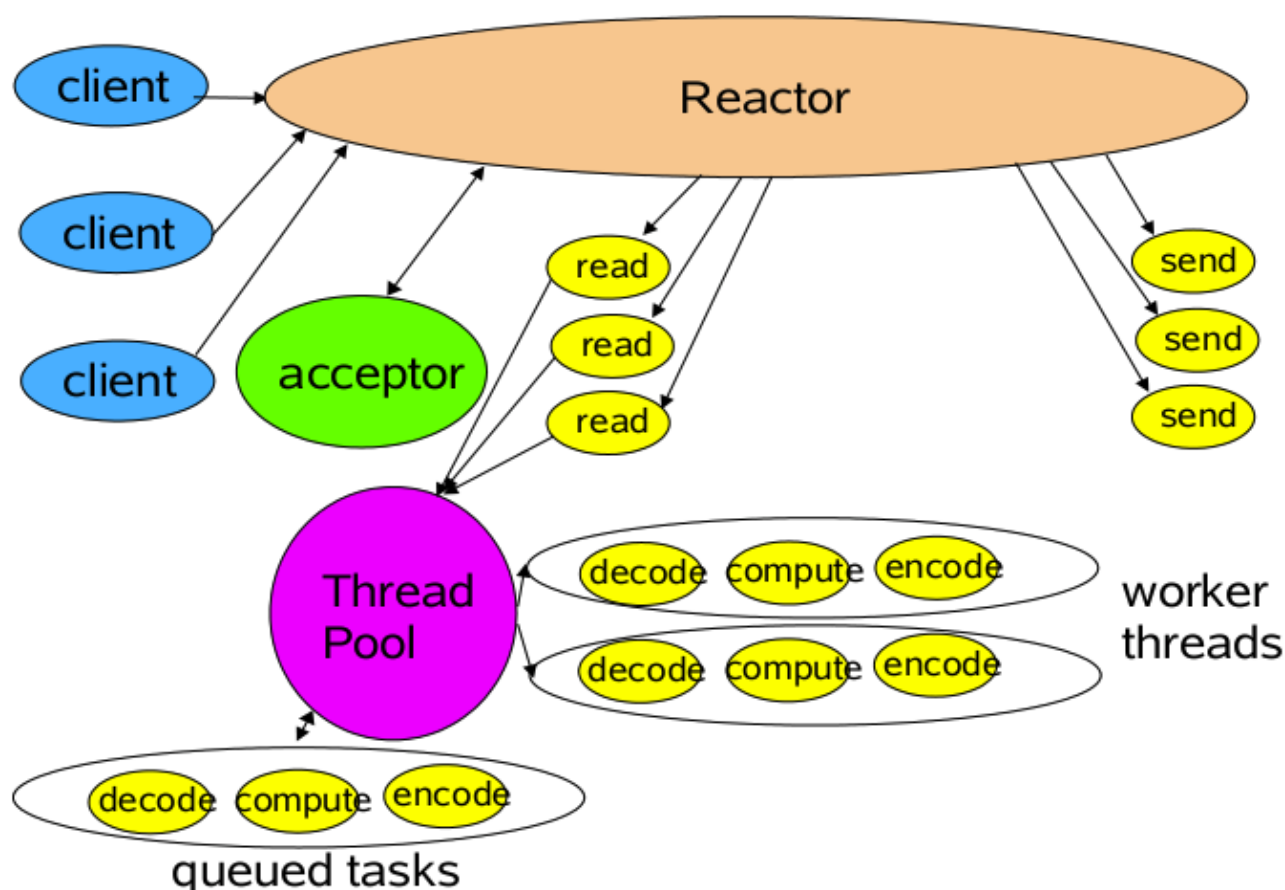
务器，一般都是多核处理器，如果能够利用多核心进行I/O，无疑对效率会有更大的提高。

仔细分析一下我们需要的线程，其实主要包括以下几种：

1. 事件分发器，单线程选择就绪的事件。
2. I/O处理器，包括connect、read、write等，这种纯CPU操作，一般开启CPU核心个线程就可以。
3. 业务线程，在处理完I/O后，业务一般还会有自己的业务逻辑，有的还会有其他的阻塞I/O，如DB操作，RPC等。只要有阻塞，就需要单独的线程。

Java的Selector对于Linux系统来说，有一个致命限制：同一个channel的select不能被并发的调用。因此，如果有多个I/O线程，必须保证：一个socket只能属于一个IoThread，而一个IoThread可以管理多个socket。

另外连接的处理和读写的处理通常可以选择分开，这样对于海量连接的注册和读写就可以分发。虽然read()和write()是比较高效无阻塞的函数，但毕竟会占用CPU，如果面对更高的并发则无能为力。



NIO在客户端的魔力

通过上面的分析，可以看出NIO在服务端对于解放线程，优化I/O和处理海量连接方面，确实有自己的用武之地。那么在客户端上，NIO又有什么使用场景呢？

常见的客户端BIO+连接池模型，可以建立n个连接，然后当某一个连接被I/O占用的时候，可以使用其他连接来提高性能。

但多线程的模型面临和服务端相同的问题：如果指望增加连接数来提高性能，则连接数又受制于线程数、线程很贵、无法建立很多线程，则性能遇到瓶颈。

每连接顺序请求的Redis

对于Redis来说，由于服务端是全局串行的，能够保证同一连接的所有请求与返回顺序一致。这样可以使用单线程+队列，把请求数据缓冲。然后pipeline发送，返回future，然后channel可读时，直接在队列中把future取回来，done()就可以了。

伪代码如下：

```
class RedisClient implements ChannelHandler{
    private BlockingQueue CmdQueue;
    private EventLoop eventLoop;
    private Channel channel;
    class Cmd{
        String cmd;
        Future result;
    }
    public Future get(String key){
        Cmd cmd= new Cmd(key);
        queue.offer(cmd);
        eventLoop.submit(new Runnable(){
            List list = new ArrayList();
            queue.drainTo(list);
            if(channel.isWritable()){
                channel.writeAndFlush(list);
            }
        });
    }
    public void ChannelReadFinish(Channel channel, Buffer Buffer){
        List result = handleBuffer();//处理数据
        //从cmdQueue取出future，并设值，future.done();
    }
    public void ChannelWritable(Channel channel){
        channel.flush();
    }
}
```

这样做，能够充分的利用pipeline来提高I/O能力，同时获取异步处理能力。

多连接短连接的HttpClient

类似于竞对抓取的项目，往往需要建立无数的HTTP短连接，然后抓取，然后销毁，当需要单机抓取上千网站线程数又受制的时候，怎么保证性能呢？

何不尝试NIO，单线程进行连接、写、读操作？如果连接、读、写操作系统没有能力处理，简单的注册一个事件，等待下次循环就好了。

如何存储不同的请求/响应呢？由于http是无状态没有版本的协议，又没有办法使用队列，好像办法不多。比较笨的办法是对于不同的socket，直接存储socket的引用作为map的key。

常见的RPC框架，如Thrift，Dubbo

这种框架内部一般维护了请求的协议和请求号，可以维护一个以请求号为key，结果的result为future的map，结合NIO+长连接，获取非常不错的性能。

NIO高级主题

Proactor与Reactor

一般情况下，I/O 复用机制需要事件分发器（event dispatcher）。事件分发器的作用，即将那些读写事件源分发给各读写事件的处理器，就像送快递的在楼下喊：谁谁谁的快递到了，快来拿吧！开发人员在开始的时候需要在分发器那里注册感兴趣的事件，并提供相应的处理器（event handler），或者是回调函数；事件分发器在适当的时候，会将请求的事件分发给这些handler或者回调函数。

涉及到事件分发器的两种模式称为：Reactor和Proactor。Reactor模式是基于同步I/O的，而Proactor模式是和异步I/O相关的。在Reactor模式中，事件分发器等待某个事件或者可应用或个操作的状态发生（比如文件描述符可读写，或者是socket可读写），事件分发器就把这个事件传给事先注册的事件处理函数或者回调函数，由后者来做实际的读写操作。

而在Proactor模式中，事件处理器（或者代由事件分发器发起）直接发起一个异步读写操作（相当于请求），而实际的工作是由操作系统来完成的。发起时，需要提供的参数包括用于存放读到数据的缓存区、读的数据大小或用于存放外发数据的缓存区，以及这个请求完后的回调函数等信息。事件分发器得知了这个请求，它默默等待这个请求的完成，然后转发完成事件给相应的事件处理器或者回调。举例来说，在Windows上事件处理器投递了一个异步IO操作（称为overlapped技术），事件分发器等IO Complete事件完成。这种异步模式的典型实现是基于操作系统底层异步API的，所以我们可称之为“系统级别”的或者“真正意义上”的异步，因为具体的读写是由操作系统代劳的。

举个例子，将有助于理解Reactor与Proactor二者的差异，以读操作为例（写操作类似）。

在Reactor中实现读

- 注册读就绪事件和相应的事件处理器。
- 事件分发器等待事件。
- 事件到来，激活分发器，分发器调用事件对应的处理器。
- 事件处理器完成实际的读操作，处理读到的数据，注册新的事件，然后返还控制权。

在Proactor中实现读：

- 处理器发起异步读操作（注意：操作系统必须支持异步IO）。在这种情况下，处理器无视IO就绪事件，它关注的是完成事件。
- 事件分发器等待操作完成事件。
- 在分发器等待过程中，操作系统利用并行的内核线程执行实际的读操作，并将结果数据存入用户自定义缓冲区，最后通知事件分发器读操作完成。
- 事件分发器呼唤处理器。
- 事件处理器处理用户自定义缓冲区中的数据，然后启动一个新的异步操作，并将控制权返回事件分发器。

可以看出，两个模式的相同点，都是对某个I/O事件的事件通知（即告诉某个模块，这个I/O操作可以进行或已经完成）。在结构上，两者也有相同点：事件分发器负责提交IO操作（异步）、查询设备是否可操作（同步），然后当条件满足时，就回调handler；不同点在于，异步情况下（Proactor），当回调handler时，表示I/O操作已经完成；同步情况下（Reactor），回调handler时，表示I/O设备可以进行某个操作（can read 或 can write）。

下面，我们将尝试应对为Proactor和Reactor模式建立可移植框架的挑战。在改进方案中，我们将Reactor原来位于事件处理器内的Read/Write操作移至分发器（不妨将这个思路称为“模拟异步”），以此寻求将Reactor多路同步I/O转化为模拟异步I/O。以读操作为例子，改进过程如下：

- 注册读就绪事件和相应的事件处理器。并为分发器提供数据缓冲区地址，需要读取数据量等信息。
- 分发器等待事件（如在select()上等待）。
- 事件到来，激活分发器。分发器执行一个非阻塞读操作（它有完成这个操作所需的全部信息），最后调用对应处理器。
- 事件处理器处理用户自定义缓冲区的数据，注册新的事件（当然同样要给出数据缓冲区地址，需要读取的数据量等信息），最后将控制权返还分发器。

如我们所见，通过对多路I/O模式功能结构的改造，可将Reactor转化为Proactor模式。改造前后，模型实际完成的工作量没有增加，只不过参与者间对工作职责稍加调换。没有工作量的改变，自然不会造成性能的削弱。对如下各步骤的比较，可以证明工作量的恒定：

标准/典型的Reactor：

- 步骤1：等待事件到来（Reactor负责）。
- 步骤2：将读就绪事件分发给用户定义的处理器（Reactor负责）。
- 步骤3：读数据（用户处理器负责）。
- 步骤4：处理数据（用户处理器负责）。

改进实现的模拟Proactor：

- 步骤1：等待事件到来（Proactor负责）。
- 步骤2：得到读就绪事件，执行读数据（现在由Proactor负责）。
- 步骤3：将读完成事件分发给用户处理器（Proactor负责）。
- 步骤4：处理数据（用户处理器负责）。

对于不提供异步I/O API的操作系统来说，这种办法可以隐藏Socket API的交互细节，从而对外暴露一个完整的异步接口。借此，我们就可以进一步构建完全可移植的，平台无关的，有通用对外接口的解决方案。

代码示例如下：

```

interface ChannelHandler{
    void channelReadComplate(Channel channel, byte[] data);
    void channelWritable(Channel channel);
}

class Channel{
    Socket socket;
    Event event;//读，写或者连接
}

//IO线程主循环：
class IoThread extends Thread{
public void run(){
    Channel channel;
    while(channel=Selector.select()){//选择就绪的事件和对应的连接
        if(channel.event==accept){
            registerNewChannelHandler(channel);//如果是新连接，则注册一个新的读写处理器
            Selector.interestedReader();
        }
        if(channel.event==write){
            getChannelHandler(channel).channelWritable(channel);//如果可以写，则执行写事件
        }
        if(channel.event==read){
            byte[] data = channel.read();
            if(channel.read()==0)//没有读到数据，表示本次数据读完了
            {
                getChannelHandler(channel).channelReadComplate(channel, data;//处理读完成事件
            }
            if(过载保护){
                Selector.interestedReader();
            }
        }
    }
}

Map<Channel, ChannelHandler> handlerMap;//所有channel的对应事件处理器
}

```

Selector.wakeup()

主要作用

解除阻塞在Selector.select()/select(long)上的线程，立即返回。

两次成功的select之间多次调用wakeup等价于一次调用。