



链滴

I/O

作者: [AshShawn](#)

原文链接: <https://ld246.com/article/1649339054232>

来源网站: [链滴](#)

许可协议: [署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

<p></p>

<p></p>- 数据到达网卡,网卡通过 DMA 搬运数据到内核 DMA 缓冲区 - 网卡向 cpu 发起硬中断,cpu 回调中断程序,创建 sk_buffer - 网卡中断程序向内核发起软中断,通知内核网络数据到达 - 内核网络协议栈处理数据头,找到对应 socket,拷贝数据到对应的 socket 缓冲区 - 程序调用内核的 read 函数读取 socket 缓冲区, - 如果 socket 缓冲区没有数据,则程序阻塞(BIO) - 当 socket 接受缓冲区有数据时,将内核空间的数据拷贝至用户空间,系统调用 read 返回, - 应用程序通过 <code>系统调用</code> 从 <code>用户态</code> 转为 <code>内核态</code> 的开销以及系统调用 <code>返回</code> 时从 <code>内核态</code> 转为 <code>用户</code> 的开销。 - 网络数据从 <code>内核空间</code> 通过 <code>CPU拷贝</code> 到 <code>用户空间</code> 的开销。 - 内核线程 <code>ksoftirqd</code> 响应 <code>软中断</code> 的开销,主要为数据包解析及数据拷贝到 socket 缓冲区。 - <code>CPU</code> 响应 <code>硬中断</code> 的开销,主要为 DMA 缓冲区数据拷贝到 sk_buffer。 - <code>DMA拷贝</code> 网络数据包到 <code>内存</code> 中的开销。 <p></p>- send 函数用户态切换成内核态,发送完毕从内核态切换成用户态开销 - 网卡发送完毕后出发 cpu 硬中断,以及在硬中断中出发软中断 - 内存拷贝开销 - 发送数据拷贝至 sk_buffer - sk_buffer 的副本拷贝 - 发送数据过大时,会进行分片拷贝成多个小的 sk_buffer

<h2 id="3--阻塞-非阻塞-同步-异步">3. 阻塞、非阻塞、同步、异步</h2>
<p>同步与异步关注的是通信机制,简单来说就是调用方与服务方之间的协同机制</p>
<blockquote>
<p>同步: 调用方 A 发起调用后,服务方 B 在处理完调用方法之前不会返回,即 AB 是串行执行的</p>
</blockquote>
<blockquote>
<p>异步: 调用方 A 发起调用后,服务方 B 会立刻返回,并在处理完调用方法后通过回调等方式通知 A , 即 AB 步调可以不一致</p>
</blockquote>
<p>不同于同步与异步, 阻塞与非阻塞关注的是一方的状态</p>
<blockquote>
<p>阻塞: 阻塞当前线程(进入循环或者挂起线程),使线程不能处理其他事</p>
</blockquote>
<blockquote>
<p>非阻塞: 不阻塞当前线程,线程仍然可以处理其他事</p>
</blockquote>
<p></p>
<p>整个流程分为两个阶段:</p>

数据准备阶段: 数据从网卡拷贝到内存再到 socket 数据缓冲区
数据拷贝阶段: 内核空间(socket 数据缓冲区)拷贝到用户空间

<h3 id="1-阻塞">1.阻塞</h3>
<p></p>
<p>第一阶段和第二阶段都会 <code>阻塞等待</code></p>
<h3 id="2-非阻塞">2.非阻塞</h3>
<p></p>
<p>第一阶段 <code>不会等待</code>, 但是在第二阶段还是会 <code>等待</code>。</p>
<h3 id="3-同步">3.同步</h3>
<p></p>
<p><code>同步模式</code> 在数据准备好后, 是由 <code>用户线程</code> 的 <code>内核</code> 来执行 <code>第二阶段</code>。所以应用程序会在第二阶段发生 <code>阻塞</code>, 直到数据从 <code>内核空间</code> 拷贝到 <code>用户空间</code>, 系统调用才会返回。<p>
<p>Linux 下的 <code>epoll</code> 和 Mac 下的 <code>kqueue</code> 都属于 <code>同步 IO</code>。</p>
<h3 id="4-异步">4.异步</h3>
<p></p>
<p><code>异步模式</code> 下是由 <code>内核</code> 来执行第二阶段的数据拷贝操作, 当 <code>内核</code> 执行完第二阶段, 会通知用户线程 IO 操作已经完成, 并将数据回调给用户线程 所以在 <code>异步模式</code> 下 <code>数据准备阶段</code> 和 <code>数据拷贝阶段</c

de> 均是由 `内核` 来完成, 不会对应用程序造成任何阻塞。</p>

4.五种 IO 模型</h2>

4.1 阻塞 IO(BIO)</h3>

</p>

阻塞读</h4>

调用系统 `read` 函数,用户线程从用户态切换为内核态,查看 socket 缓冲区是否有数据</p>

socket 缓冲区有数据, 将数据从内核空间拷贝到用户空间,系统 IO 调用返回

socket 缓冲区无数据,用户线程让出 CPU,进入 `阻塞状态`。当数据到达 `socket` 接收缓冲区后, 内核唤醒 `阻塞状态` 中的用户线程进入 `就绪状态`, 随后经过 CPU 的调度获取到 `CPU quota` 进入 `运行状态`, 将内核空间的数据拷贝到用户空间, 随后系统调用返回。

阻塞写</h4>

调用系统 `send` 函数,用户线程从用户态切换为内核态,将发送数据从用户空间拷贝到内核空间中的 `Socket` 发送缓冲区中。</p>

`Socket` 发送缓冲区能够容纳下发送数据时, 用户线程会将全部的发送数据写入 `Socket` 缓冲区,执行发送流程后返回

`Socket` 发送缓冲区空间不够, 无法容纳下全部发送数据时, 用户线程让出 CPU 进入 `阻塞状态`, 直到 `Socket` 发送缓冲区能够容纳下全部发送数据时, 内核唤醒用户线程, 执行后续发送流程。

阻塞 IO 模型</h3>

每个请求单独线程处理

并发量增大则服务端瞬间会创建大量线程,造成资源浪费

创建好线程后,若长期处于空闲时间,则线程们一直阻塞

线程切换开销大

适用场景</h3>

链接少

并发低

重型操作

4.2 非阻塞 IO(NIO)</h3>

</p>

用尽量少的线程去处理更多的连接</p>

非阻塞读</h4>

调用系统 `read` 函数,用户线程从用户态切换为内核态,查看 socket 缓冲区是否有数据</p>

socket 缓冲区有数据, 将数据从内核空间拷贝到用户空间,系统 IO 调用返回

socket 缓冲区无数据,系统调用立马返回,线程 `不会阻塞`,也 `不交出cpu`,而是会继续 `轮训` 直到 `Socket` 接收缓冲区中有数据为止

非阻塞写</h4>

<p>调用系统 `send` 函数,用户线程从用户态切换为内核态,将发送数据从用户空间到内核空间中的 `Socket` 发送缓冲区中。</p>

<code>Socket</code> 发送缓冲区能够容纳下发送数据时, 用户线程会将全部的发送数据写入 `Socket` 缓冲区,执行发送流程后返回

<code>Socket</code> 发送缓冲区空间不够, 无法容纳下全部发送数据时, <code>能写多少</code>, 写不下了, 就立即返回。并将写入到发送缓冲区的字节数返回给应用程序, 方便线程不断的 <code>轮训</code> 尝试将 <code>剩下的数据</code> 写入发送缓冲区中。

<h5 id="非阻塞IO模型">非阻塞 IO 模型</h5>

<p></p>

<p>利用一个线程或者很少的线程, 去 <code>不断地轮询</code> 每个 <code>Socket</code> 的接收缓冲区是否有数据到达, 如果没有数据, <code>不必阻塞</code> 线, 而是接着去 <code>轮询</code> 下一个 <code>Socket</code> 接收缓冲区, 直到轮询到数后, 处理连接上的读写, 或者交给业务线程池去处理, 轮询线程则 <code>继续轮询</code> 其他的 <code>Socket</code> 接收缓冲区。</p>

<h5 id="适用场景">适用场景</h5>

<p>在 <code>非阻塞IO模型</code> 下, 需要用户线程去 <code>不断地</code> 发起 <code>系统调用</code> 去轮训 <code>Socket</code> 接收缓冲区, 这就需要用户线程不断地从 <code>用户态</code> 切换到 <code>内核态</code>, <code>内核态</code> 切换到 <code>用户态</code>。</p>

<p>单纯的 <code>非阻塞IO</code> 模型还是无法适用于高并发的场景</p>

<h3 id="4-3-IO多路复用">4.3 IO 多路复用</h3>

<p></p>

多路: 我们的核心需求是要用尽可能少的线程来处理尽可能多的连接, 这的 <code>多路</code> 指的就是我们需要处理的众多连接。

复用: 核心需求要求我们使用 <code>尽可能少的线程</code>, <code>尽可能少的系统开销</code> 去处理 <code>尽可能多</code> 的连接 (<code>多路</code>) 那么这里的 <code>复用</code> 指的就是用 <code>有限的资源</code>, 比如用一个线程或者定数量的线程去处理众多连接上的读写事件。换句话说, 在 <code>阻塞IO模型</code> 中一个连接就需要分配一个独立的线程去专门处理这个连接上的读写, 到了 <code>IO多路复用模型</code>, 多个连接可以 <code>复用</code> 这一个独立的线程去处理这多个连接上的读写。

<p>非阻塞 IO 中其实实现了多路复用,只不过是用户态去不断轮询,会导致用户态与内核态频繁切换,费资源</p>

<p>我们可以把频繁的轮询操作交给操作系统内核来替我们完成, 这样就避免了在 <code>用户空间</code> 频繁的去使用系统调用来轮询所带来的性能开销。</p>

<h5 id="select">select</h5>

<p></p>

<p><code>select</code> 系统调用将 <code>轮询</code> 的操作交给了 <code>内核</code> 来帮助我们完成, 从而避免了在 <code>用户空间</code> 不断的发起轮询所带来的的系统性能开。</p>

首先用户线程在发起 <code>select</code> 系统调用时会 <code>阻塞</code> 在 <code>select</code> 系统调用上。此时, 用户线程从 <code>用户态</code> 切换到了 <code>内核态</code> 完成了一次 <code>上下文切换</code>

用户线程将需要监听的 <code>Socket</code> 对应的文件描述符 <code>fd</code> 数组

过 `select` 系统调用传递给内核。此时，用户线程将 `用户空间` 中文件描述符 `fd` 数组 `拷贝` 到 `内核空间`。

- 当用户线程调用完 `select` 后开始进入 `阻塞状态`，`内核` 开始轮询遍历 `fd` 数组，查看 `fd` 对应的 `Socket` 接收缓冲区中是否有数据到来。如果有数据到来，则将 `fd` 对应 `BitMap` 的值设置为 `1`。如果没有数据到来，则保持值为 `0`。
- 内核遍历一遍 `fd` 数组后，如果发现有些 `fd` 上有 IO 数据到，则将修改后的 `fd` 数组返回给用户线程。此时，会将 `fd` 数组从 `内核空间` 拷贝到 `用户空间`。
- 当内核将修改后的 `fd` 数组返回给用户线程后，用户线程解除 `阻塞`，由用户线程开始遍历 `fd` 数组然后找出 `fd` 数组中值为 `1` 的 `Socket` 文件描述符。最后对这些 `Socket` 发起系统调用读取数据。
- 由于内核在遍历的过程中已经修改了 `fd` 数组，所以在用户线程遍历完 `fd` 数组后获取到 `IO就绪` 的 `Socket` 后，就需要 `重置` `fd` 数组，并重新调用 `select` 传入重置后的 `fd` 数组，让内核发起新一轮遍历轮询

这里的**文件描述符数组** 其实是一个 `BitMap`，`BitMap` 下标为 `文件描述符fd`，下标对应的值为：`1` 表示该 `fd` 上有读写事件，`0` 表示该 `fd` 上没有读写事件。

性能开销

- 发生 2 次上下文 `切换`**
- 发生 2 次文件描述符集合的 `拷贝`**
- 虽然由原来在 `用户空间` 发起轮询 `优化成了` 在 `内核空间` 发起轮询但 `select` 不会告诉用户线程到底是哪些 `Socket` 上发生了 `IO就绪` 事件，只是对 `IO就绪` 的 `Socket` 作了标记，用户线程依然要 `遍历` 文件描述符集合去查找具体 `IO就绪` 的 `Socket`。时间复杂度依然为 `O(n)`。
- `内核` 会对原始的 `文件描述符集合` 进行修改。导致每次在 `用户空间` 重新发起 `select` 调用时，都需要对 `文件描述符集合` 进行 `重置`。
- `BitMap` 结构的文件描述符集合，长度为固定的 `1024`，所以能监听 `0~1023` 的文件描述符。
- `select` 系统调用 不是线程安全的。
- 并发量增大,开销线性增长,只适合 1000 个左右的并发

poll

poll 与 select 原理之一,主要变化点为修改了 1024 个文件描述符的限制,这样就没有了最大描述数量的限制（当然还会受到系统文件描述符限制）

epoll

select 和 poll 现有问题:

- 每次都需要全量传入 fd 集合,导致大量 fd 在用户空间与内核空间频繁复制
- 内核不会通知具体 `IO就绪` 的 `socket`，只是在这些 `IO就绪` 的 `socket` 上打好标记，所以当 `select` 系统调用返回时，在 `用户空间` 还是需要 `完整遍历` 一遍 `socket` 文件描述集合来获取具体 `IO就绪` 的 `socket`。
- 在 `内核空间` 中也是通过遍历的方式来得到 `IO就绪` 的 `socket`。
- 红黑树

<p></p>

<p>epoll 内部使用红黑树来保存所有监听的 socket,添加和查找复杂度为 $O(\log n)$,节点中每个数字为 socket 的句柄</p>

<ol start="2">

就绪队列

<p></p>

<p>当 socket 从网络中获取到数据后,会发生通知给 epoll, epoll 会将当前 socket 添加到就绪队中,并且唤醒等待中的进程(也就是调用 <code>epoll_wait</code> 的进程)。</p>

<p>当进程被唤醒后,就会从就绪队列中,把就绪的 socket 复制到用户提供的数组中。</p>

<p>当 <code>epoll_wait</code> 返回后,用户就可以从 <code>events</code> 数组中获取到就绪的 socket,并可对其进行读写操作。</p>

<ol start="3">

常用方法

<p>1) 调用 epoll_create()建立一个 epoll 对象(在 epoll 文件系统中为这个句柄对象分配资源)</p>

<p>2) 调用 epoll_ctl 向 epoll 对象中添加这 100 万个连接的套接字</p>

<p>3) 调用 epoll_wait 收集发生的事件的连接</p>

<p>epoll 增加优化:</p>

<code>epoll</code> 在内核中通过 <code>红黑树</code> 管理海量的连接,所以在调用 <code>epoll_wait</code> 获取 <code>IO就绪</code> 的 socket 时,不需要传入监听的 socket 文件描述符。从而避免了海量的文件描述符集合在 <code>用户空间</code> 和 <code>内核空间</code> 中来回复制。

<blockquote>

<p><code>select, poll</code> 每次调用时都需要传递全量的文件描述符集合,导致大量频繁的页操作。</p>

</blockquote>

<code>epoll</code> 仅会通知 <code>IO就绪</code> 的 socket。避免了在用户空间遍历开销。

<blockquote>

<p><code>select, poll</code> 只会在 <code>IO就绪</code> 的 socket 上打好标记,依然全量返回,所以在用户空间还需要用户程序在一次遍历全量集合找出具体 <code>IO就绪</code> 的 socket。</p>

</blockquote>

<code>epoll</code> 通过在 <code>socket</code> 的等待队列上注册回调函数 <code>epoll_callback</code> 通知用户程序 <code>IO就绪</code> 的 socket。避免了在内核中轮询的开销。

<blockquote>

<p>大部分情况下 <code>socket</code> 上并不总是 <code>IO活跃</code> 的,在面对海量连接的情况下,<code>select, poll</code> 采用内核轮询的方式获取 <code>IO活跃</code> 的 socket,无疑是性能低下的核心原因。</p>

</blockquote>

<h3 id="4-4-信号驱动IO">4.4 信号驱动 IO</h3>

<p></p>

<h3 id="4-5-异步IO-AIO-">4.5 异步 IO(AIO)</h3>

<p></p>

<h3 id="4-6-总结">4.6 总结</h3>

<p></p>

<p></p>

<p></p>

<p></p>

<p></p>

<p></p>

<p></p>

<p>java 中的 NIO 其实是多路复用 IO 模型,与 linux 中的 NIO(非阻塞 IO)存在异同点</p>

<p>Blocking IO 主要是里面的方法操作步骤都是同步的,而 Non-Blocking IO 则不是。它是通过册事件来触发对应的 handler 来执行。比起 thread per connection 的方式,它能更好的利用资源因为回调的机制对于异步的通信方式来说减少了轮询或者其它强制同步机制的开销,效率算是比较理的</p>

<p>NIO 中的非阻塞主要体现在</p>

read 和 write 都是立刻返回的,不像 bio 中可能会阻塞,因为 nio 中 read 和 write 针对的都是 buffer, bio 中则是针对流,效率较低

事件线程是阻塞的,即调用了 epoll_wait,但是事件处理线程是非阻塞的

<h2 id="5-用户空间的IO线程模型">5.用户空间的 IO 线程模型</h2>

<h3 id="5-1-Reactor">5.1 Reactor</h3>

<p><code>Reactor</code> 是利用 <code>NIO</code> 对 <code>IO线程</code> 进行不同的分工: </p>

使用前边我们提到的 <code>IO多路复用模型</code> 比如 <code>select,poll,epoll,kqueue</code>,进行 IO 事件的注册和监听。

将监听到 <code>就绪的IO事件</code> 分发 <code>dispatch</code> 到各个具体的处理 <code>Handler</code> 中进行相应的 <code>IO事件处理</code>。

<p>通过 <code>IO多路复用技术</code> 就可以不断的监听 <code>IO事件</code>,不断的发 <code>dispatch</code>,就像一个 <code>反应堆</code> 一样,看起来像不断的产生 <code>IO事件</code>,因此我们称这种模式为 <code>Reactor</code> 模型。</p>

<h5 id="单Reactor单线程">单 Reactor 单线程</h5>

<p></p>

- 单 `Reactor` 意味着只有一个 `epoll` 对象，用来监听所有的件，比如 `连接事件`，`读写事件`。
- `单线程` 意味着只有一个线程来执行 `epoll_wait` 获取 `O就绪` 的 `Socket`，然后对这些就绪的 `Socket` 执行读，以及后边的业务处理也依然是这个线程。

<p></p> - `单reactor` 是只有一个 `epoll` 对象来监听所有的 `IO件`，一个线程来调用 `epoll_wait` 获取 `IO就绪` 的 `Socket` - 但是当 `IO就绪事件` 产生时，这些 `IO事件` 对应处理的业务 `Handler`，我们是通过线程池来执行。这样相比 `单Reactor单线程` 型提高了执行效率，充分发挥了多核 CPU 的优势。 <p></p> - 我们由原来的 `单Reactor` 变为了 `多Reactor`。<code>主Reactor</code> 用来优先 `专门` 做优先级最高的事情，也就是处理 `连接事件`，对应的处理 `Handler` 就是图中的 `acceptor`。 - 当创建好连接，建立好对应的 `socket` 后，在 `acceptor` 中要监听的 `read事件` 注册到 `从Reactor` 中，由 `从Reactor` 来监听 `socket` 上的 `读写` 事件。 <p>等待后续更新.....</p> <p>等待后续更新....</p> <p>netty 中 IO 模型为 Reactor 模型,三种类型都涉及,单主要使用的是 `主从Reactor多线程型`</p> <p></p> - `Reactor` 在 `netty` 中是以 `group` 的形式现的，`netty` 中将 `Reactor` 分为两组，一组是 `MainReactorGroup` 也就是我们在编码中常常看到的 `EventLoopGroup bossGroup` 另一组是 `SubReactorGroup` 也就是我们在编码中常常看到的 `EventLoopGroup workerGroup`。 - `MainReactorGroup` 中通常只有一个 `Reactor`，专门负责最重要的事情，也就是监听连接 `accept` 事件。当有连接事件产生时，在对应的处理 `handler acceptor` 中创建初始化相应的 `NioSocketChannel`（表一个 `Socket连接`）。然后以 `负载均衡` 的方式在 `Sub`

`ReactorGroup` 中选取一个 `Reactor`，注册上去，监听 `Read事件`。

- `MainReactorGroup` 中只有一个 `Reactor` 的原因是，通常我们服务端程序只会 `绑定监听` 一个端口，如果要 `绑定监听` 多个端口，就会配置多个 `Reactor`。
- `SubReactorGroup` 中有多个 `Reactor`，具体 `Reactor` 的个数可以由系统参数 `-D io.netty.eventLoopThreads` 指定。默认的 `Reactor` 的个数为 `CPU核数 * 2`。`SubReactorGroup` 中的 `Reactor` 主要负责监听 `读写事件`，每一个 `Reactor` 负责监听一组 `socket连接`。将全量的连接 `分摊` 在一个 `Reactor` 中。
- 一个 `Reactor` 分配一个 `IO线程`，这个 `IO线程` 负责从 `Reactor` 中获取 `IO就绪事件`，执行 `IO调用` 取 `IO数据`，执行 `PipeLine`。

`Socket连接` 在创建后就被 `固定的分配` 给一个 `Reactor`，所以一个 `Socket连接` 也只会被一个固定的 `IO线程` 执行，每个 `Socket连接` 分配一个独立的 `PipeLine` 实例，用来排这个 `Socket连接` 上的 `IO处理逻辑`。这种 `无锁串行化` 的设计的目的是为了防止多线程并发执行同一个 `socket 连接上的 IO逻辑处理`，防止出现 `线程安全问题`。同时使系统吞吐量达到最大化

由于每个 `Reactor` 中只有一个 `IO线程`，这个 `IO线程` 既要执行 `IO活跃Socket连接` 对应的 `PipeLine` 中的 `ChannelHandler`，又要从 `Reactor` 中获取 `IO就绪事件`，执行 `IO调用`。所以 `PipeLine` 中 `ChannelHandler` 中执行的逻辑不能耗时太长，尽量将耗时的业务逻辑处理放入单独的业务线程池中处理，则会影响其他连接的 `IO读写`，从而进一步影响整个服务程序的 `IO吞吐`。

- 当 `IO请求` 在业务线程中完成相应的业务逻辑处理后，在业务线程中利用持有的 `ChannelHandlerContext` 引用将响应数据在 `PipeLine` 中反向传播，最终写回给客户端。

文章参考: https://ld246.com/forward?goto=https%3A%2F%2Fmp.weixin.qq.com/s/Ylf_ClsjIWJf8nTUQuwMoA