



链滴

字节跳动在 Go 网络库上的实践

作者: [panjf2000](#)

原文链接: <https://ld246.com/article/1589257501102>

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>两种方式各有优缺点，netpoll 采用前者策略，水平触发时效性更好，容错率高，主动 I/O 可以集内存使用和管理，提供 nocopy 操作并减少 GC。事实上一些热门开源网络库也是采用方式一的设计如 easygo、evio、gnet 等。</p>

<p>但使用 LT 也带来另一个问题，即底层主动 I/O 和上层代码并发操作 buffer，引入额外的并发开。比如：I/O 读数据写 buffer 和上层代码读 buffer 存在并发读写，反之亦然。为了保证数据正确性同时不引入锁竞争，现有的开源网络库通常采取 同步处理 buffer(easygo, evio) 或者将 buffer 再 copy 一份提供给上层代码(gnet) 等方式，均不适合业务处理或存在 copy 开销。</p>

<p>另一方面，常见的 bytes、bufio、ringbuffer 等 buffer 库，均存在 growth 需要 copy 原数组据，以及只能扩容无法缩容，占用大量内存等问题。因此我们希望引入一种新的 Buffer 形式，一举决上述两方面的问题。</p>

<p>NoCopy Buffer 基于链表数组实现，如下图所示，我们将 []byte 数组抽象为 block，并以链表接的形式将 block 组合为 NoCopy Buffer，同时引入了引用计数、nocopy API 和对象池。</p> <p></p> <p>NoCopy Buffer 相比常见的 bytes、bufio、ringbuffer 等有以下优势：</p> 读写并行无锁，支持 nocopy 地流式读写 读写分别操作头尾指针，相互不干扰。 高效扩缩容 扩容阶段，直接在尾指针后添加新的 block 即可，无需 copy 原数组。 缩容阶段，头指针会直接释放使用完毕的 block 节点，完成缩容。每个 block 都有独立的引用计，当释放的 block 不再有引用时，主动回收 block 节点。 灵活切片和拼接 buffer (链表特性) 支持任意读取分段(nocopy)，上层代码可以 nocopy 地并行处理数据流分段，无需关心生命周，通过引用计数 GC。 支持任意拼接(nocopy)，写 buffer 支持通过 block 拼接到尾指针后的形式，无需 copy，保证据只写一次。 NoCopy Buffer 池化，减少 GC 将每个 []byte 数组视为 block 节点，构建对象池维护空闲 block，由此复用 block，减少内存占和 GC。 <p>基于该 NoCopy Buffer，我们实现了 NoCopy Thrift，使得编解码过程内存零分配零拷贝。</p> <p>RPC 调用通常采用短连接或者长连接池的形式，一次调用绑定一个连接，那么当上下游规模很大情况下，网络中存在的连接数以 MxN 的速度扩张，带来巨大的调度压力和计算开销，给服务治理造困难。因此，我们希望引入一种 "在单一长连接上并行处理调用" 的形式，来减少网络中的连接数，种方案即称为 "连接多路复用"。</p> <p>当前业界也存在一些开源的连接多路复用方案，掣肘于代码层面的束缚，这些方案均需要 copy buffer 来实现数据分包和合并，导致实际性能并不理想。而上述 NoCopy Buffer 基于其灵活切片和拼 原文链接：[字节跳动在 Go 网络库上的实践](#)

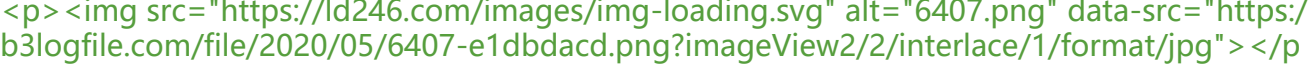
的特性，很好的支持了 nocopy 的数据分包和合并，使得实现高性能连接多路复用方案成为可能。 </p></p><p>基于 netpoll 的连接多路复用设计如下图所示，我们将 NoCopy Buffer(及其分片) 抽象为虚拟连接，使得上层代码保持同 net.Conn 相同的调用体验。与此同时，在底层代码上通过协议分包将真实接上的数据灵活的分配到虚拟连接上；或通过协议编码合并发送虚拟连接数据。 </p><p> </p><p>连接多路复用方案包含以下核心要素： </p>虚拟连接实质上是 NoCopy Buffer，目的是替换真正的连接，规避内存 copy。 上层的业务逻辑/编解码 均在虚拟连接上完成，上层逻辑可以异步独立并行执行。 Shared map引入分片锁来减少锁力度。 在调用端使用 sequence id 来标记请求，并使用分片锁存储 id 对应的回调。 在接收响应数据后，根据 sequence id 来找到对应回调并执行。 协议分包和编码如何识别完整的请求响应数据包是连接多路复用方案可行的关键，因此需要引入协议。 这里采用 thrift header protocol 协议，通过消息头判断数据包完整性，通过 sequence id 标记请求和响应的对应关系。 <h2 id="ZeroCopy">ZeroCopy</h2><p>这里所说的 ZeroCopy，指的是 Linux 所提供的 ZeroCopy 的能力。上一章中我们说了业务层零拷贝，而众所周知，当我们调用 sendmsg 系统调用发包的时候，实际上仍然是会产生一次数据的拷贝，并且在大包场景下这个拷贝的消耗非常明显。以 100M 为例，perf 可以看到如下结果： </p><p> </p><p>这还仅仅是普通 tcp 发包的占用，在我们的场景下，大部分服务都会接入 Service Mesh，所以一次发包中，一共会有 3 次拷贝：业务进程到内核、内核到 sidecar、sidecar 再到内核。这使得有包需求的业务，拷贝所导致的 CPU 占用会特别明显，如下图： </p><p> </p><p>为了解决这个问题，我们选择了使用 Linux 提供的 ZeroCopy API（在 4.14 以后支持 send；5.4 以后支持 receive）。但是这引入了一个额外的工程问题：ZeroCopy send API 和原先调用方式不兼容，无法很好地共存。这里简单介绍一下 ZeroCopy send 的工作方式：业务进程调用 sendmsg 之后，endmsg 会记录下 iovec 的地址并立即返回，这时候业务进程不能释放这段内存，需要通过 epoll 等内核回调一个信号表明某段 iovec 已经发送成功之后才能释放。由于我们并不希望更改业务方的使用法，需要对上层提供同步收发的接口，所以很难基于现有的 API 同时提供 ZeroCopy 和非 ZeroCopy 的抽象；而由于 ZeroCopy 在小包场景下是有性能损耗的，所以也不能将这个作为默认的选项。 </p><p>于是，字节跳动框架组和字节跳动内核组合作，由内核组提供了同步的接口：当调用 sendmsg 时候，内核会监听并拦截内核原先给业务的回调，并且在回调完成后才会让 sendmsg 返回。这使得我们无需更改原有模型，可以很方便地接入 ZeroCopy send。同时，字节跳动内核组还实现了基于 unix domain socket 的 ZeroCopy，可以使得业务进程与 Mesh sidecar 之间的通信也达到零拷贝。 </p><p>在使用了 ZeroCopy send 后，perf 可以看到内核不再有 copy 的占用： </p></div><div data-bbox="655 936 930 951" data-label="Page-Footer"><p>原文链接：字节跳动在 Go 网络库上的实践</p></div>



从 CPU 占用数值上看，大包场景下 ZeroCopy 能够比非 ZeroCopy 节省一半的 CPU。

Go 调度导致的延迟问题分享

在我们实践过程中，发现我们新写的 netpoll 虽然在 avg 延迟上表现胜于 Go 原生的 net 库，是在 p99 和 max 延迟上要普遍略高于 Go 原生的 net 库，并且尖刺也会更加明显，如下图（Go 1.1，蓝色为 netpoll + 多路复用，绿色为 netpoll + 长连接，黄色为 net 库 + 长连接）：



我们尝试了很多种办法去优化，但是收效甚微。最终，我们定位出这个延迟并非是由于 netpoll 身的开销导致的，而是由于 go 的调度导致的，比如说：

由于在 netpoll 中，SubReactor 本身也是一个 goroutine，受调度影响，不能保证 EpollWait 调之后马上执行，所以这一块会有延迟；

同时，由于用来处理 I/O 事件的 SubReactor 和用来处理连接监听的 MainReactor 本身也是 go routine，所以实际上很难保证在多核情况之下，这些 Reactor 能并行执行；甚至在最极端情况之下，能这些 Reactor 会挂在同一个 P 下，最终变成了串行执行，无法充分利用多核优势；

由于 EpollWait 回调之后，SubReactor 内是串行处理 I/O 事件的，导致排在最后的事件可能会长尾问题；

在连接多路复用场景下，由于每个连接绑定了一个 SubReactor，故延迟完全取决于这个 SubReactor 的调度，导致尖刺会更加明显。

由于 Go 在 runtime 中对于 net 库有做特殊优化，所以 net 库不会有以上情况；同时 net 库是 routine-per-connection 的模型，所以能确保请求能并行执行而不会相互影响。

对于以上这个问题，我们目前解决的思路有两个：

修改 Go runtime 源码，在 Go runtime 中注册一个回调，每次调度时调用 EpollWait，把获取的 fd 传递给回调执行；

与字节跳动内核组合作，支持同时批量读/写多个连接，解决串行问题。另外，经过我们的测试，o 1.14 能够使得延迟略有降低同时更加平稳，但是所能达到的极限 QPS 更低。希望我们的思路能够业界同样遇到此问题的同学提供一些参考。

后记

希望以上的分享能够对社区有所帮助。同时，我们也在加速建设 netpoll 以及基于 netpoll 的新架 KiteX。欢迎各位感兴趣的同学加入我们，共同建设 Go 语言生态！

参考资料

http://man7.org/linux/man-pages/man7/epoll.7.html

https://golang.org/src/runtime/proc.o

https://github.com/panjf2000/gnet

https://github.com/tidwall/evio

原文

字节跳动在 Go 网络库的实践

 </p>