
深入剖析 **Go** 语言运行时：**IO** 轮询器

孔俊

2022-11-21

Contents

1	netpoller 概述	2
2	底层机制	3
2.1	非阻塞 IO	3
2.2	IO 多路复用	5
3	数据结构	5
3.1	pollDesc	5
3.2	pollCache	7
4	netpoller 基础函数	8
4.1	初始化轮询器	9
4.2	注册文件描述符	11
4.3	检测文件描述符是否由 poller 管理	11
4.4	删除文件描述符	11
4.5	中断 netpoller	12
4.6	文件描述符状态	12
5	goroutine 挂起与恢复	15
5.1	netpollblock	15
5.2	netpollunblock	19
6	“拦截”阻塞 IO	22
6.1	封装文件描述符	22
6.2	初始化文件描述符	24
6.3	封装系统调用	26
7	实例	31

1 netpoller 概述

考虑一个基于 goroutine-per-connection 模型的 TCP echo server:

```
1 import (
2     "fmt"
3     "io"
4     "log"
5     "net"
6 )
7
8 func worker(conn net.Conn) {
9     defer conn.Close()
10    b := make([]byte, 512)
11    for {
12        size, err := conn.Read(b)
13        if err == io.EOF {
14            break
15        }
16        if err != nil {
17            log.Fatal(err)
18        }
19        size, err = conn.Write(b[0:size])
20        if err != nil {
21            log.Fatal(err)
22        }
23    }
24 }
25
26 func main() {
27     listener, err := net.Listen("tcp", "127.0.0.1:8080")
28     if err != nil {
29         log.Fatal(err)
30     }
31     for {
32         conn, err := listener.Accept()
33         if err != nil {
34             log.Fatal(err)
35         }
36         go worker(conn)
37     }
38 }
```

从用户侧看，系统该调用阻塞 goroutine，Go scheduler 调度其他 goroutine。问题在于，goroutine 复用在线程上，如果 IO 系统调用（如 `read(2)`/`write(2)`）阻塞，直接阻塞 goroutine 所在线程，Go scheduler 将没有机会调度 goroutine！

为了实现 goroutine-per-connection 网络编程模型，必须提供“IO 系统调用阻塞 goroutine 而非 OS 线程”的抽象，这意味着必须“拦截”IO 系统调用对线程的阻塞。

Go runtime 通过非阻塞 IO 和 IO 多路复用机制实现了一个 IO 轮询器（netpoller），提供了“网络 IO

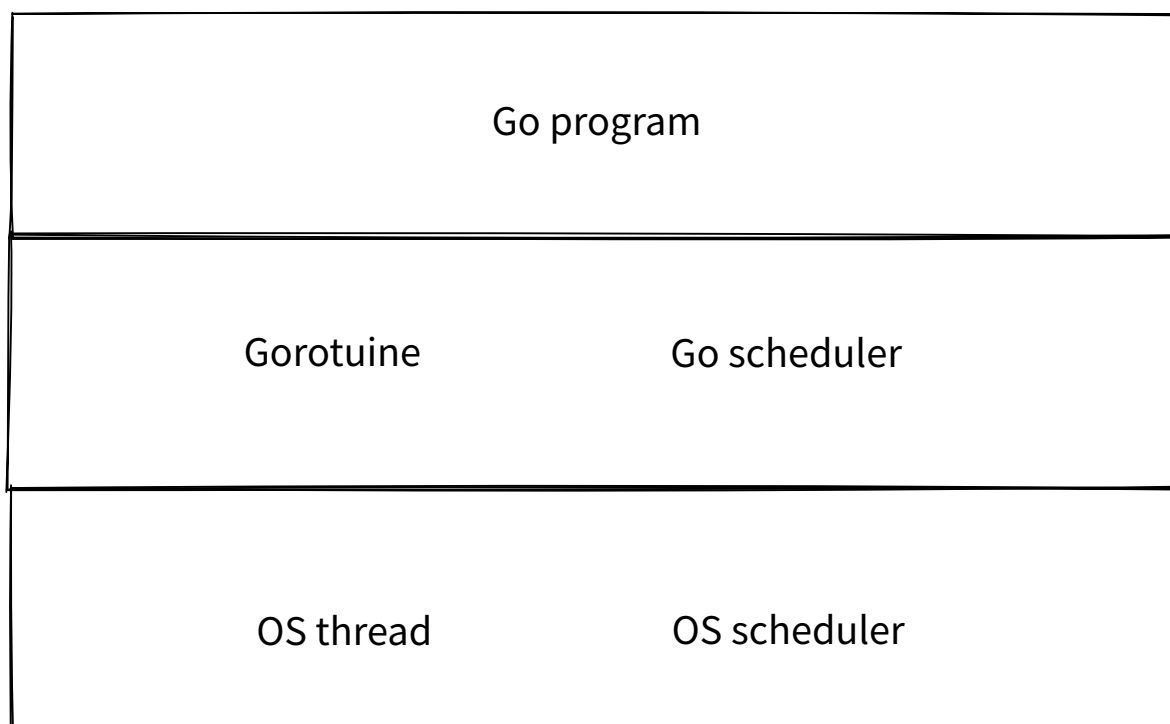


Figure1: Goroutine and thread

阻塞 goroutine” 的抽象：

- 尽量将文件描述符设置为非阻塞的，关联当前 goroutine 到该文件描述符，注册其上的 IO 事件到轮询器中。
- IO 系统调用失败（返回EAGAIN）时，阻塞该 goroutine，切换到其他 goroutine。
- 在特定时机轮询 IO 事件，获取就绪的文件描述符，调度对应的 goroutine。

2 底层机制

2.1 非阻塞 IO

为了避免 IO 系统调用阻塞 OS 线程，必须使用非阻塞 IO。

非阻塞 IO 和阻塞 IO 的差别在于，非阻塞 IO 给了用户处理 IO 未就绪的机会：当文件描述符上不可进行 IO 操作时，非阻塞 IO 系统调用返回错误，而阻塞 IO 直接阻塞线程直到 IO 完成。

非阻塞 IO 不等于高性能，对于非阻塞 IO 和阻塞 IO，从 IO 未就绪到 IO 就绪并传输数据的时间都是相同的。如果直接如图 2 轮询文件描述符，非阻塞 IO 需要调用更多次系统调用，性能反而更差。

非阻塞 IO 高性能的关键在于避免等待 IO 就绪。非阻塞 IO 系统调用返回错误后，用户知道 IO 未就绪，去执行别的操作，IO 就绪后再接着执行先前的逻辑。

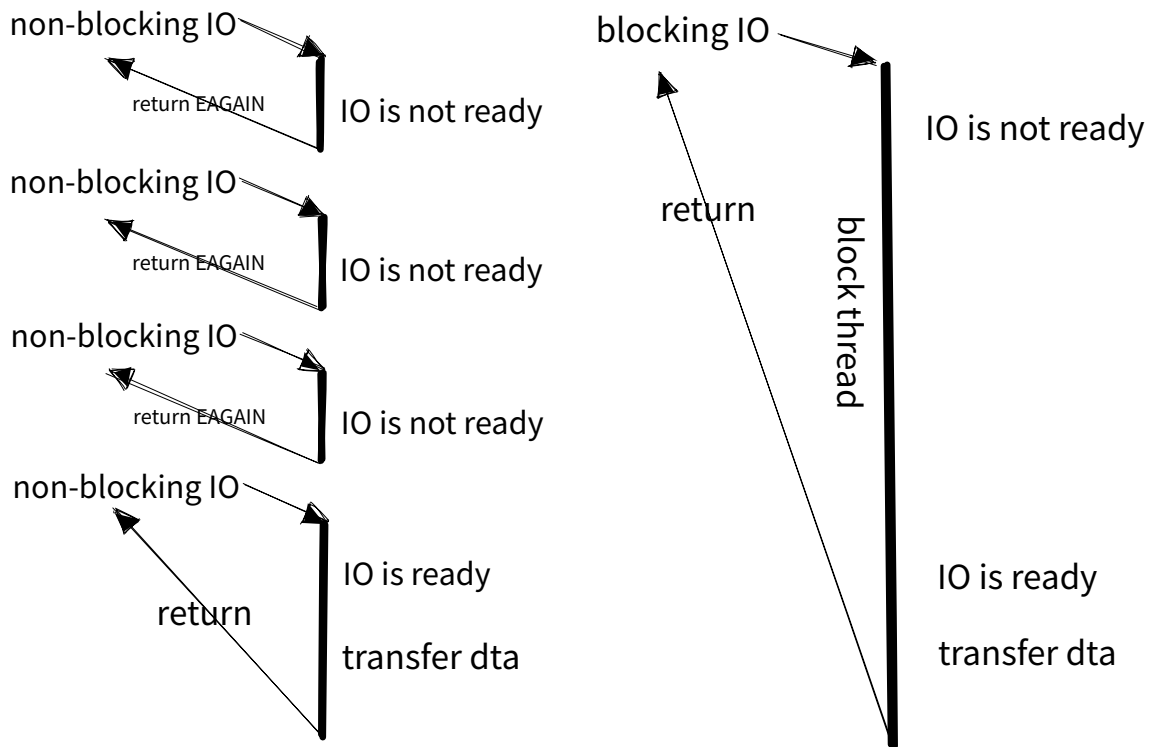


Figure2: Blocking IO and non-blocking IO

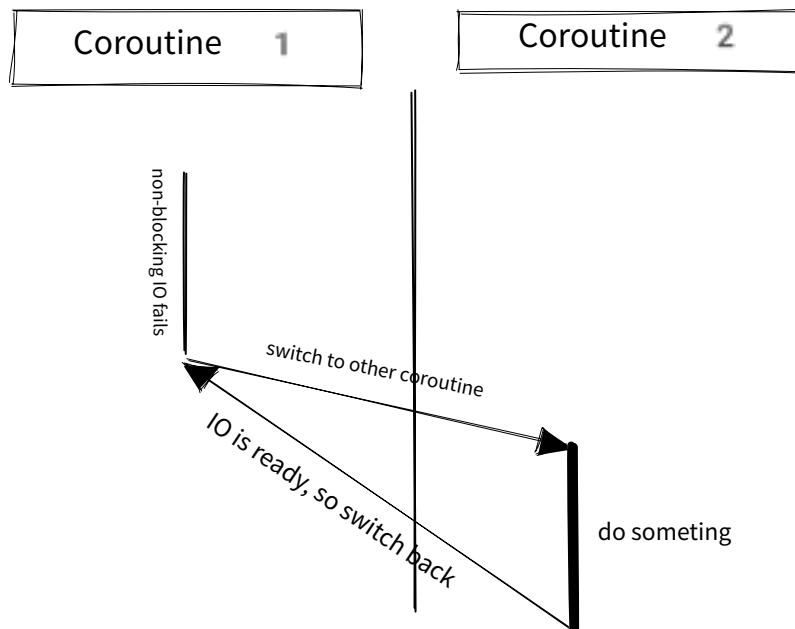


Figure 3: Coroutines with non-blocking IO

这是经典的协程模型，协程 1 调用非阻塞 IO 失败，于是主动放弃控制权，切换到协程 2，协程 2 执行其他逻辑，发现协程 1 等待的 IO 就绪后，再切换回协程 1。

Go runtime 提供的抽象是“IO 阻塞 goroutine，切换到其他 goroutine”，因此用户代码不需要“主动换出”，Go scheduler 会替用户调度 goroutine。

2.2 IO 多路复用

非阻塞 IO 高性能的关键在于避免等待 IO 就绪，因此操作系统内核提供了 IO 多路复用机制，通知用户 IO 已就绪，避免用户忙等待 IO。也就是说，轮询 IO 是否就绪的责任从用户转移到了操作系统。

Linux 平台上一般使用 `epoll` 进行 IO 多路复用。`epoll_create(2)` 创建轮询器（一个包含所有相关信息的文件描述符），`epoll_ctl(2)` 向轮询器注册/删除文件描述符上的 IO 事件，`epoll_wait(2)` 阻塞直到发生事件并返回待处理事件。

注意，`epoll` 只能处理 `socket` 和 `pipe`，不能处理磁盘文件。

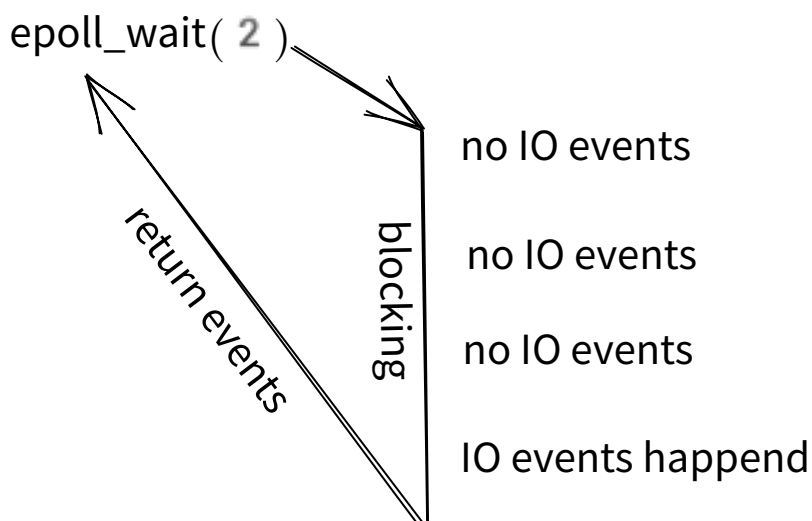


Figure 4: `epoll_wait(2)` blocks the OS thread

3 数据结构

3.1 pollDesc

Go runtime 定义了 `pollDesc` 描述用于网络轮询的文件描述符，其中包括 OS 原生文件描述符、锁、读写者 goroutine、定时器信息以及描述符状态 `atomicInfo`。

```

1 // Network poller descriptor.
2 //
3 // No heap pointers.
4 type pollDesc struct {
5     _      sys.NotInHeap
6     link *pollDesc // in pollcache, protected by pollcache.lock
7     fd    uintptr   // constant for pollDesc usage lifetime
8
9     // atomicInfo holds bits from closing, rd, and wd,
10    // which are only ever written while holding the lock,
11    // summarized for use by netpollcheckerr,
12    // which cannot acquire the lock.
13    // After writing these fields under lock in a way that
14    // might change the summary, code must call publishInfo
15    // before releasing the lock.
16    // Code that changes fields and then calls netpollunblock
17    // (while still holding the lock) must call publishInfo
18    // before calling netpollunblock, because publishInfo is what
19    // stops netpollblock from blocking anew
20    // (by changing the result of netpollcheckerr).
21    // atomicInfo also holds the eventErr bit,
22    // recording whether a poll event on the fd got an error;
23    // atomicInfo is the only source of truth for that bit.
24    atomicInfo atomic.Uint32 // atomic pollInfo
25
26    // rg, wg are accessed atomically and hold g pointers.
27    // (Using atomic.Uintptr here is similar to using guintptr
28    // elsewhere.)
29    rg atomic.Uintptr // pdReady, pdWait, G waiting for read or pdNil
30    wg atomic.Uintptr // pdReady, pdWait, G waiting for write or pdNil
31
32    lock    mutex // protects the following fields
33    closing bool
34    user    uint32    // user settable cookie
35    rseq    uintptr   // protects from stale read timers
36    rt      timer     // read deadline timer (set if rt.f != nil)
37    rd      int64     // read deadline (a nanotime in the future, -1
38    // when expired)
39    wseq    uintptr   // protects from stale write timers
40    wt      timer     // write deadline timer
41    wd      int64     // write deadline (a nanotime in the future, -1
42    // when expired)
43    self    *pollDesc // storage for indirect interface. See (*pollDesc
44    //).makeArg.
45 }

```

`link`是指向`pollDesc`的指针，所有`pollDesc`放到一个链表中管理。

`fd`是`pollDesc`使用的 OS 原生文件描述符，在`pollDesc`的整个生命周期保持不变。

`atomicInfo`维护`pollDesc`的轮询状态。

`rg/wg`是该文件描述符上阻塞的 `goroutine`。

`resq`以下的成员用于定时器处理 IO 超时。

后文会详细介绍`pollDesc`各字段。

3.2 pollCache

`pollCache`是一个单向链表，用于管理`pollDesc`的分配和释放。因为多个 `goroutine` 可能并发地注册 `poller`，所以需要互斥锁保护并发访问。

```
1 type pollCache struct {
2     lock  mutex
3     first *pollDesc
4     // PollDesc objects must be type-stable,
5     // because we can get ready notification from epoll/kqueue
6     // after the descriptor is closed/reused.
7     // Stale notifications are detected using seq variable,
8     // seq is incremented when deadlines are changed or descriptor is
9     // reused.
10 }
```

`(*pollCache).alloc()`从`pollCache`中分配一个`pollDesc`:

1. 若`pollCache`为空，分配大小为`pollBlockSize`（4K）的非 GC 内存（`pollDesc`）。
2. 摘下链表头上的`pollDesc`，返回给用户。


```

1 func (c *pollCache) alloc() *pollDesc {
2     lock(&c.lock)
3     // 分配内存，创造 pollDesc
4     if c.first == nil {
5         const pdSize = unsafe.Sizeof(pollDesc{})
6         n := pollBlockSize / pdSize
7         if n == 0 {
8             n = 1
9         }
10        // Must be in non-GC memory because can be referenced
11        // only from epoll/kqueue internals.
12        mem := persistentalloc(n*pdSize, 0, &memstats.other_sys)
13        for i := uintptr(0); i < n; i++ {
14            pd := (*pollDesc)(add(mem, i*pdSize))
15            pd.link = c.first
16            c.first = pd
17        }
18    }
19    // 摘下表头的 pollDesc 并返回
20    pd := c.first
21    c.first = pd.link
22    lockInit(&pd.lock, lockRankPollDesc) // Go runtime 内部的锁实现了
    层级，见 runtime/lockrank.go
23    unlock(&c.lock)
24    return pd
25 }

```

```

1 `(*pollCache).free()`回收`pollDesc`，将其添加到表头。
2 ``go
3 func (c *pollCache) free(pd *pollDesc) {
4     lock(&c.lock)
5     pd.link = c.first
6     c.first = pd
7     unlock(&c.lock)
8 }

```

注意，`(*pollCache).alloc()`分配的`pollDesc`的`rg/wg`均为`pdNil`，`(*pollCache).free()`没有修改字段`rg/wg`。

分配和回收有以下两个断言：

- `pollCache`中的`pollDesc`的`rg/wg`要么是`pdNil`，要么是`pdReady`。
- `(*pollCache).free()`回收的`pollDesc`的`rg/wg`要么是`pdNil`，要么是`pdReady`。

4 netpoller 基础函数

由于不同平台上的 IO 多路复用接口不同，能力也各不相同（Linux 的 `epoll` 不支持 `disk IO`，但 `FreeBSD` 的 `kqueue` 支持），Go runtime 在 `runtime/netpoll.go` 中定义了以下平台无关的接口，通

过`go:build`条件编译，调用到对应平台的实现。

```
1 // Integrated network poller (platform-independent part).
2 // A particular implementation (epoll/kqueue/port/AIX/Windows)
3 // must define the following functions:
4 //
5 // func netpollinit()
6 //     Initialize the poller. Only called once.
7 //
8 // func netpollopen(fd uintptr, pd *pollDesc) int32
9 //     Arm edge-triggered notifications for fd. The pd argument is to
10 //     pass
11 //     back to netpollready when fd is ready. Return an errno value.
12 //
13 // func netpollclose(fd uintptr) int32
14 //     Disable notifications for fd. Return an errno value.
15 //
16 // func netpoll(delta int64) gList
17 //     Poll the network. If delta < 0, block indefinitely. If delta ==
18 //     0,
19 //     poll without blocking. If delta > 0, block for up to delta
20 //     nanoseconds.
21 //     Return a list of goroutines built by calling netpollready.
22 //
23 // func netpollBreak()
24 //     Wake up the network poller, assumed to be blocked in netpoll.
25 //
26 // func netpollIsPollDescriptor(fd uintptr) bool
27 //     Reports whether fd is a file descriptor used by the poller.
```

`netpoller` 提供了种种接口声明在内部包 `internal/poll/fd_poll_runtime.go` 中，实现在 `runtime/netpoll.go` 中，以`poll_runtime_`为前缀。因此，Go 语言运行时没有给用户提供任何 `netpoller` 接口。

4.1 初始化轮询器

互斥锁`netpollInitLock`保护`netpollInited`,`netpollInited`指示是否已初始化轮询器。

```
1 var (
2     netpollInitLock mutex
3     netpollInited   atomic.Uint32
4
5     pollcache      pollCache
6     netpollWaiters atomic.Uint32
7 )
```

`netpollGenericInit()`调用`netpollinit()`初始化 `poller`, `netpollGenericInit()`相当于`Once.Do()`。

```
1 func netpollGenericInit() {
2     if netpollInited.Load() == 0 {
3         lockInit(&netpollInitLock, lockRankNetpollInit)
4         lock(&netpollInitLock)
5         if netpollInited.Load() == 0 {
6             netpollinit()
7             netpollInited.Store(1)
8         }
9         unlock(&netpollInitLock)
10    }
11 }
```

Linux 上的 `netpollinit()` 做以下几件事：

1. 创建 edge-triggered 的 `epoll` 轮询器。
2. 创建一个 `pipe` 用于打断 `netpoller`。
3. 注册读端的读操作到 `epoll` 上，其 `syscall.EpollEvent` 的 `data` 字段是 `netpollBreakRd`（`pipe` 的读端）。

该 `pipe` 用于和 `netpoller` 通信：`netpoller` 监听 `pipe` 的读端，调用 `netpollBreak()` 向该 `pipe` 写入数据，`netpoller` 监听到发生 `pipe` 上的读事件，从而打断 `netpoller`。

```
1 // runtime/netpoll.go
2 func netpollinit() {
3     var errno uintptr
4     // 创建 epoll 轮询器
5     epfd, errno = syscall.EpollCreate1(syscall.EPOLL_CLOEXEC)
6     if errno != 0 {
7         println("runtime: epollcreate failed with", errno)
8         throw("runtime: netpollinit failed")
9     }
10    // 创建用于打断 poller 的 pipe
11    r, w, errpipe := nonblockingPipe()
12    if errpipe != 0 {
13        println("runtime: pipe failed with", -errpipe)
14        throw("runtime: pipe failed")
15    }
16    // 注册 pipe 上的读事件到 poller
17    ev := syscall.EpollEvent{
18        Events: syscall.EPOLLIN,
19    }
20    *(*uintptr)(unsafe.Pointer(&ev.Data)) = &netpollBreakRd
21    errno = syscall.EpollCtl(epfd, syscall.EPOLL_CTL_ADD, r, &ev)
22    if errno != 0 {
23        println("runtime: epollctl failed with", errno)
24        throw("runtime: epollctl failed")
25    }
26    netpollBreakRd = uintptr(r)
27    netpollBreakWr = uintptr(w)
28 }
```

Go runtime 中只有一个 `epoll` 轮询器，定义为全局变量。

```
1 // runtime/netpoll.go
2 var (
3     epfd int32 = -1 // epoll descriptor
4     // ...
5 )
```

4.2 注册文件描述符

`netpollopen()`封装了`epoll_ctl(2)`，注册文件描述符上的读写事件到 `netpoller` 上，其中`epoll_event`的`data`字段设置为 `epoll fd`。

```
1 // runtime/netpoll.go
2 func netpollopen(fd uintptr, pd *pollDesc) uintptr {
3     var ev syscall.EpollEvent
4     ev.Events = syscall.EPOLLIN | syscall.EPOLLOUT | syscall.EPOLLRDHUP
5         | syscall.EPOLLET
6     *(*pollDesc)(unsafe.Pointer(&ev.Data)) = pd
7     return syscall.EpollCtl(epfd, syscall.EPOLL_CTL_ADD, int32(fd), &ev)
8 }
```

4.3 检测文件描述符是否由 `poller` 管理

`epoll_event`的`data`字段可以用于判断该描述符是否由 `poller` 管理，在 Go 1.9 中只用于测试。

```
1 // IsPollDescriptor reports whether fd is the descriptor being used by
2 // the poller.
3 // This is only used for testing.
4 func IsPollDescriptor(fd uintptr) bool {
5     return runtime_isPollServerDescriptor(fd)
6 }
```

4.4 删除文件描述符

从 `poller` 删除该描述符同理，直接调用`epoll_ctl(2)`即可。

```
1 // runtime/netpoll.go
2 func netpollclose(fd uintptr) uintptr {
3     var ev syscall.EpollEvent
4     return syscall.EpollCtl(epfd, syscall.EPOLL_CTL_DEL, int32(fd), &ev)
5 }
```

4.5 中断 netpoller

前面初始化轮询器时，把 pipe 的读端上的读事件注册到了 epoll 上，因此只要向该 pipe 写数据，就能触发 epoll，从而打断 netpoller。

```
1 // netpollBreak interrupts an epollwait.
2 func netpollBreak() {
3     // Failing to cas indicates there is an in-flight wakeup, so we're
4     // done here.
5     if !netpollWakeSig.CompareAndSwap(0, 1) {
6         return
7     }
8     for {
9         var b byte
10        n := write(netpollBreakWr, unsafe.Pointer(&b), 1)
11        if n == 1 {
12            break
13        }
14        if n == -_EINTR {
15            continue
16        }
17        if n == -_EAGAIN {
18            return
19        }
20        println("runtime: netpollBreak write failed with", -n)
21        throw("runtime: netpollBreak write failed")
22    }
23 }
```

在死循环中写 pipe 是因为 UNIX 上系统中断可被信号（signal）中断，中断的系统调用返回错误码 EINTR，for 循环处理 write(2) 被中断的情况。

4.6 文件描述符状态

pollDesc 结构体的 atomicInfo 字段维护文件描述符 fd 的轮询状态，包括 IO 超时，描述符关闭等。

```
1 // runtime/netpoll.go
2 type pollDesc struct {
3     // ...
4     fd uintptr // constant for pollDesc usage lifetime
5
6     // atomicInfo holds bits from closing, rd, and wd,
7     // which are only ever written while holding the lock,
8     // summarized for use by netpollcheckerr,
9     // which cannot acquire the lock.
10    // After writing these fields under lock in a way that
11    // might change the summary, code must call publishInfo
12    // before releasing the lock.
13    // Code that changes fields and then calls netpollunblock
14    // (while still holding the lock) must call publishInfo
15    // before calling netpollunblock, because publishInfo is what
16    // stops netpollblock from blocking anew
17    // (by changing the result of netpollcheckerr).
18    // atomicInfo also holds the eventErr bit,
19    // recording whether a poll event on the fd got an error;
20    // atomicInfo is the only source of truth for that bit.
21    atomicInfo atomic.Uint32 // atomic pollInfo
22    //...
23 }
```

`atomicInfo`维护的状态最终用于判断文件描述符是否出错。

```
1 // runtime/netpoll.go
2 // Error codes returned by runtime_pollReset and runtime_pollWait.
3 // These must match the values in internal/poll/fd_poll_runtime.go.
4 const (
5     pollNoError      = 0 // no error
6     pollErrClosing   = 1 // descriptor is closed
7     pollErrTimeout   = 2 // I/O timeout
8     pollErrNotPollable = 3 // general error polling descriptor
9 )
```

`atomicInfo`被设置为私有变量，由`pushInfo()`更新，通过`netpollcheckerr()`判断是否出错。

会修改文件描述符状态的代码，必须在操作完成后调用`pushInfo()`更新状态。

```

1 // runtime/netpoll.go
2 // publishInfo updates pd.atomicInfo (returned by pd.info)
3 // using the other values in pd.
4 // It must be called while holding pd.lock,
5 // and it must be called after changing anything
6 // that might affect the info bits.
7 // In practice this means after changing closing
8 // or changing rd or wd from < 0 to >= 0.
9 func (pd *pollDesc) publishInfo() {
10     var info uint32
11     if pd.closing {
12         info |= pollClosing
13     }
14     if pd.rd < 0 {
15         info |= pollExpiredReadDeadline
16     }
17     if pd.wd < 0 {
18         info |= pollExpiredWriteDeadline
19     }
20
21     // Set all of x except the pollEventErr bit.
22     x := pd.atomicInfo.Load()
23     for !pd.atomicInfo.CompareAndSwap(x, (x&pollEventErr)|info) {
24         x = pd.atomicInfo.Load()
25     }
26 }

```

`netpollcheckerr()` 通过位运算从 `atomicInfo` 中提取状态并返回错误。前面提到, `atomicInfo` 维护的是文件描述符的轮询状态, 而非文件描述符上系统调用的状态。也就是说, `netpollcheckerr()` 不判断其上的系统调用是否成功, 即使该描述符上的系统调用出错, `netpollcheckerr()` 也返回 `true`。

具体的系统调用的错误处理, 应当交给 Go 语言封装系统调用的函数, 而非 `netpoller`。

```

1 func netpollcheckerr(pd *pollDesc, mode int32) int {
2     info := pd.info()
3     if info.closing() {
4         return pollErrClosing
5     }
6     if (mode == 'r' && info.expiredReadDeadline()) || (mode == 'w' &&
7         info.expiredWriteDeadline()) {
8         return pollErrTimeout
9     }
10    // Report an event scanning error only on a read event.
11    // An error on a write event will be captured in a subsequent
12    // write call that is able to report a more specific error.
13    if mode == 'r' && info.eventErr() {
14        return pollErrNotPollable
15    }
16    return pollNoError
17 }

```

5 goroutine 挂起与恢复

每个 `pollDesc` 都关联了一个读者 `goroutine`、一个写者 `goroutine`，当该文件描述符不读写时挂起响应的读写者。

`pollDesc` 上读写者的挂起与恢复由一个有限状态机驱动，状态如下：

```

1 // runtime/netpoller.go
2 // pollDesc contains 2 binary semaphores, rg and wg, to park reader and
  // writer
3 // goroutines respectively. The semaphore can be in the following
  // states:
4 //
5 //     pdReady - io readiness notification is pending;
6 //             a goroutine consumes the notification by changing the
  //             state to pdNil.
7 //     pdWait - a goroutine prepares to park on the semaphore, but not
  //             yet parked;
8 //             the goroutine commits to park by changing the state to G
  //             pointer,
9 //             or, alternatively, concurrent io notification changes
  //             the state to pdReady,
10 //            or, alternatively, concurrent timeout/close changes the
  //            state to pdNil.
11 //     G pointer - the goroutine is blocked on the semaphore;
12 //                io notification or timeout/close changes the state to
  //                pdReady or pdNil respectively
13 //                and unparks the goroutine.
14 //     pdNil - none of the above.
15 const (
16     pdNil    uintptr = 0
17     pdReady  uintptr = 1
18     pdWait   uintptr = 2
19 )

```

`pdNil` 是 `rg/wg` 的初始状态，其值为 0，相当于 `Nil`。`pollDesc` 上的 `goroutine` 已经阻塞（`park`）时，设置为该 `goroutine` 的指针。

状态转换图如下。

5.1 netpollblock

`netpollblock(pd *pollDesc, mode int32, waitio bool)` 试图阻塞文件 `pd` 上的读写者，返回真值指示 IO 是否就绪。参数 `mode` 指示对文件的操作模式（读/写/读写），`waitio` 判断是否由于等待 IO 而阻塞 `goroutine`。

`netpollblock()` 接收的 `pollDesc` 上的 `goroutine` 一定处于 `pdReady`（已 `netpollunblock()`）或 `pdNil`（文件初始状态）。

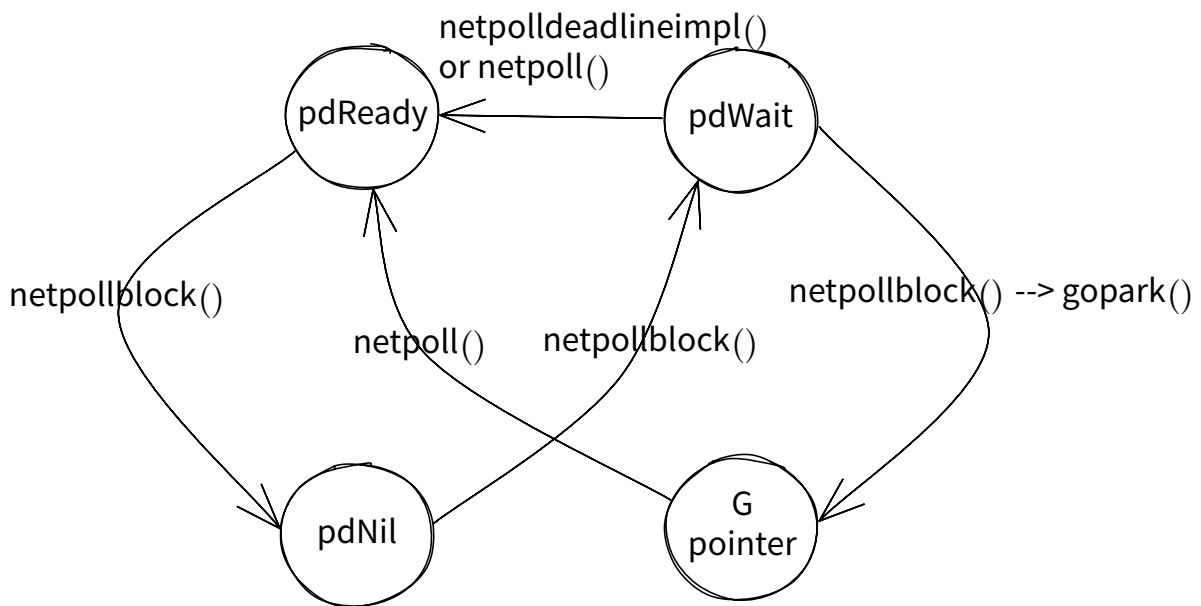


Figure 5: The FSM of goroutine

- **pdReady**状态: 该状态由`netpollunblock()`函数设置, 表示一次IO就绪。`netpollblock()`消耗一次IO就绪通告, 将状态从**pdReady**转换到**pdNil**。
- **pdNil**状态: 该状态是`pollDesc`上 goroutine 默认状态, 转换为**pdWait**。
- **pdWait**状态: 该状态表示 goroutine 将要阻塞但还未阻塞, `netpollblock()`调用`gopark()`阻塞该 goroutine。

```
1 // returns true if IO is ready, or false if timed out or closed
2 // waitio - wait only for completed IO, ignore errors
3 // Concurrent calls to netpollblock in the same mode are forbidden
4 // , as pollDesc
5 // can hold only a single waiting goroutine for each mode.
6 func netpollblock(pd *pollDesc, mode int32, waitio bool) bool {
7     gpp := &pd.rg
8     if mode == 'w' {
9         gpp = &pd.wg
10    }
11    // set the gpp semaphore to pdWait
12    for {
13        // Consume notification if already ready.
14        if gpp.CompareAndSwap(pdReady, pdNil) {
15            return true
16        }
17        if gpp.CompareAndSwap(pdNil, pdWait) {
18            break
19        }
20        // Double check that this isn't corrupt; otherwise we'd loop
21        // forever.
22        if v := gpp.Load(); v != pdReady && v != pdNil {
23            throw("runtime: double wait")
24        }
25    }
26    // need to recheck error states after setting gpp to pdWait
27    // this is necessary because runtime_pollUnblock/
28    // runtime_pollSetDeadline/deadlineimpl
29    // do the opposite: store to closing/rd/wd, publishInfo, load of
30    // rg/wg
31    if waitio || netpollcheckerr(pd, mode) == pollNoError {
32        gopark(netpollblockcommit, unsafe.Pointer(gpp),
33            waitReasonIOWait, traceEvGoBlockNet, 5)
34    }
35    // be careful to not lose concurrent pdReady notification
36    old := gpp.Swap(pdNil)
37    if old > pdWait {
38        throw("runtime: corrupted polldesc")
39    }
40    return old == pdReady
41 }
```

`gopark()` 是 Go scheduler 中负责阻塞 goroutine 的函数，原型如下：

```

1 // Puts the current goroutine into a waiting state and calls unlockf on
  the
2 // system stack.
3 //
4 // If unlockf returns false, the goroutine is resumed.
5 //
6 // unlockf must not access this G's stack, as it may be moved between
7 // the call to gopark and the call to unlockf.
8 //
9 // Note that because unlockf is called after putting the G into a
  waiting
10 // state, the G may have already been readied by the time unlockf is
  called
11 // unless there is external synchronization preventing the G from being
12 // readied. If unlockf returns false, it must guarantee that the G
  cannot be
13 // externally readied.
14 //
15 // Reason explains why the goroutine has been parked. It is displayed
  in stack
16 // traces and heap dumps. Reasons should be unique and descriptive. Do
  not
17 // re-use reasons, add new ones.
18 func gopark(unlockf func(*g, unsafe.Pointer) bool, lock unsafe.Pointer,
  reason waitReason, traceEv byte, traceskip int) {

```

简单地说，这里的 `gopark()` 阻塞 goroutine `gpp`，调用回调函数 `netpollblockcommit()`，并将该 goroutine 阻塞的原因设置为 `waitReasonIOWait`。

`netpollblockcommit()` 做两件事：1. 将 `pollDesc` 中的 `rg` 或 `wg` 设置为指向该 goroutine 的指针。2. 递增 `netpollWaiters` 原子变量。

```

1 func netpollblockcommit(gp *g, gpp unsafe.Pointer) bool {
2     r := atomic.Casuintptr((*uintptr)(gpp), pdWait, uintptr(unsafe.
  Pointer(gp)))
3     if r {
4         // Bump the count of goroutines waiting for the poller.
5         // The scheduler uses this to decide whether to block
6         // waiting for the poller if there is nothing else to do.
7         netpollWaiters.Add(1)
8     }
9     return r
10 }

```

原子变量 `netpollWaiters` 用于避免不必要的 `epoll` 阻塞。没有 goroutine 等待 netpoller 唤醒时，不轮询 IO 事件。`findRunnable()` 展示了这一优化。

```

1 // runtime/proc.go
2 func findRunnable() (gp *g, inheritTime, tryWakeP bool) {
3     // ...
4
5     // Poll network.
6     // This netpoll is only an optimization before we resort to
7     // stealing.
8     // We can safely skip it if there are no waiters or a thread is
9     // blocked
10    // in netpoll already. If there is any kind of logical race with
11    // that
12    // blocked thread (e.g. it has already returned from netpoll, but
13    // does
14    // not set lastpoll yet), this thread will do blocking netpoll
15    // below
16    // anyway.
17    if netpollnited() && netpollWaiters.Load() > 0 && sched.lastpoll.
18        Load() != 0 {
19        if list := netpoll(0); !list.empty() { // non-blocking
20            gp := list.pop()
21            injectglist(&list)
22            casgstatus(gp, _Gwaiting, _Grunnable)
23            if trace.enabled {
24                traceGoUnpark(gp, 0)
25            }
26            return gp, false, false
27        }
28    }
29    // ...
30 }

```

5.2 netpollunblock

`netpollunblock(pd *pollDesc, mode int32, ioready bool)`是`netpollblock()`的逆操作，修改对应 `goroutine` 的状态，并返回可运行的 `goroutine` 的指针。

`netpollunblock()`可以接收全部四种 `goroutine` 状态的`pollDesc`。

- **pdNil**: 该状态表示初始状态。若`ioready`为真，则切换到`pdReady`，返回`nil`；否则，不修改状态，直接返回`nil`。`ioready`为真，说明其上 IO 就绪，因此 `netpoller` 恢复其上的 `goroutine`；`ioready`为假，说明该`netpollunblock()`因 IO 超时或取消 IO 而被调用。
- **pdReady**: 该状态表示 IO 就绪。`goroutine` 未阻塞，不需要恢复该 `goroutine`，因此直接返回`nil`。
- **pdWait**: 该状态表示 `goroutine` 即将被阻塞，还未被阻塞。因此根据`ioready`切换到`pdNil`或`pdReady`状态并返回`nil` (`pdNil`)。
- **goroutine 指针**: 该 `goroutine` 阻塞后由`netpollblockcommit()`设置。根据`ioready`切换到`pdNil`或`pdReady`，返回该 `goroutine` 指针。

```
1 func netpollunblock(pd *pollDesc, mode int32, ioready bool) *g {
2     gpp := &pd.rg
3     if mode == 'w' {
4         gpp = &pd.wg
5     }
6
7     for {
8         old := gpp.Load()
9         if old == pdReady {
10             return nil
11         }
12         if old == pdNil && !ioready {
13             // Only set pdReady for ioready. runtime_pollWait
14             // will check for timeout/cancel before waiting.
15             return nil
16         }
17         var new uintptr
18         if ioready {
19             new = pdReady
20         }
21         if gpp.CompareAndSwap(old, new) {
22             if old == pdWait {
23                 old = pdNil
24             }
25             // 注意，pdNil 的值为 0，相当于 Nil。
26             return (*g)(unsafe.Pointer(old))
27         }
28     }
29 }
```

`netpollunblock()` 只修改状态并返回 goroutine 指针, 恢复 goroutine 的任务交给 `netpollgoready()`。

`netpollgoready()` 是 `netpollblockcommit()` 的逆操作, 完成以下两件事:

1. 递减 `netpollWaiters`
2. 恢复 goroutine

```
1 func netpollgoready(gp *g, traceskip int) {
2     netpollWaiters.Add(-1)
3     goready(gp, traceskip+1)
4 }
```

```

1  ## 轮询 IO 事件
2  所有文件描述符都注册到一个 epoll 轮询器上，因此`netpoll()`轮询整个系统中的 IO 事件：
3
4  1. 根据参数`delay`设置 epoll 超时时间；
5  2. 调用`epoll_wait()`轮询 IO 事件；
6  3. 处理 IO 事件：
7      1. 通过`netpollBreakRd`打断`netpoll()`；
8      2. epoll 错误，调用`(*pollDesc).setEventErr()`设置轮询错误；
9      3. IO 就绪，调用`netpollgoready()`恢复该文件描述符上阻塞的 goroutine。
10 `netpoll()`中的`syscall.EpollEvent`数组长度为 128，因此一次最多监听到 128 个事件。
11 ```go
12 // netpoll checks for ready network connections.
13 // Returns list of goroutines that become runnable.
14 // delay < 0: blocks indefinitely
15 // delay == 0: does not block, just polls
16 // delay > 0: block for up to that many nanoseconds
17 func netpoll(delay int64) gList {
18     if epfd == -1 {
19         return gList{}
20     }
21     var waitms int32
22     if delay < 0 {
23         waitms = -1
24     } else if delay == 0 {
25         waitms = 0
26     } else if delay < 1e6 {
27         waitms = 1
28     } else if delay < 1e15 {
29         waitms = int32(delay / 1e6)
30     } else {
31         // An arbitrary cap on how long to wait for a timer.
32         // 1e9 ms == ~11.5 days.
33         waitms = 1e9
34     }
35     var events [128]syscall.EpollEvent
36 retry:
37     n, errno := syscall.EpollWait(epfd, events[:], int32(len(events)), waitms)
38     if errno != 0 {
39         if errno != _EINTR {
40             println("runtime: epollwait on fd", epfd, "failed with", errno)
41             throw("runtime: netpoll failed")
42         }
43         // If a timed sleep was interrupted, just return to
44         // recalculate how long we should sleep now.
45         if waitms > 0 {
46             return gList{}
47         }
48         goto retry
49     }
50     var toRun gList
51     for i := int32(0); i < n; i++ {
52         ev := events[i]
53         if ev.Events == 0 {
54             continue
55         }
56

```

Go runtime 没有单独使用一个线程轮询 IO 事件，而是由特定操作触发 netpoller。触发网络轮询的事件包括：

- goroutine 调度
- 垃圾回收的某些阶段：Drain、MarkDone、MarkTermination、StartTheWorld。
- sysmon goroutine

调用关系如下：

```
1 schedule()
2     -> findRunnable()
3     -> netpoll()
4
5 gcStart()
6     -> startTheWorldWithSema()
7     -> netpoll()
8 gcDrain()
9     -> pollWork()
10 gcMarkDone()
11     -> startTheWorldWithSema()
12     -> netpoll()
13 gcMarkTermination()
14     -> startTheWorldWithSema()
15     -> netpoll()
16 startTheWorld()
17     -> startTheWorldWithSema()
18     -> netpoll()
19
20 sysmon()
21     -> netpoll()
```

6 “拦截”阻塞 IO

显然，goroutine 在非阻塞 IO 未就绪时被挂起，为了“拦截”阻塞，必须将 IO 系统调用都封装成非阻塞的。

6.1 封装文件描述符

Go runtime 在 internal/poll/fd_unix.go 定义了 UNIX 上的文件描述符。

```
1 // FD is a file descriptor. The net and os packages use this type as a
2 // field of a larger type representing a network connection or OS file.
3 type FD struct {
4     // Lock sysfd and serialize access to Read and Write methods.
5     fdmu fdMutex
6
7     // System file descriptor. Immutable until Close.
8     Sysfd int
9
10    // I/O poller.
11    pd pollDesc
12
13    // Writev cache.
14    iovecs *[]syscall.Iovec
15
16    // Semaphore signaled when file is closed.
17    csema uint32
18
19    // Non-zero if this file has been set to blocking mode.
20    isBlocking uint32
21
22    // Whether this is a streaming descriptor, as opposed to a
23    // packet-based descriptor like a UDP socket. Immutable.
24    IsStream bool
25
26    // Whether a zero byte read indicates EOF. This is false for a
27    // message based socket connection.
28    ZeroReadIsEOF bool
29
30    // Whether this is a file rather than a network socket.
31    isFile bool
32 }
```

结构体FD包含系统原生文件描述符Sysfd和I/O pollerpd，Sysfd用于执行系统调用，pd用于轮询该文件。

Go 语言的文件在底层都是 internal/poll/fd_unix.go 中定义的poll.FD。os.File定义于 os/types.go，其中内嵌了*file，file是特定 OS 上的文件类型，通过go:build条件编译调用到特定平台上的实现。

```
1 // File represents an open file descriptor.
2 type File struct {
3     *file // os specific
4 }
```

Linux 平台的file定义在 os/file_unix.go，其中包含poll.FD。


```
1 // file is the real representation of *File.
2 // The extra level of indirection ensures that no clients of os
3 // can overwrite this data, which could cause the finalizer
4 // to close the wrong file descriptor.
5 type file struct {
6     pfd          poll.FD
7     name         string
8     dirinfo      *dirInfo // nil unless directory being read
9     nonblock     bool      // whether we set nonblocking mode
10    stdoutOrErr  bool      // whether this is stdout or stderr
11    appendMode   bool      // whether file is opened for appending
12 }
```

再看net包的socket(), 返回的文件描述符是*netFD。

```
1 // socket returns a network file descriptor that is ready for
2 // asynchronous I/O using the network poller.
3 func socket(ctx context.Context, net string, family, sotype, proto int,
4             ipv6only bool, laddr, raddr sockaddr, ctrlCtxFn func(context.
5             Context, string, string, syscall.RawConn) error) (fd *netFD, err
6             error) {
7     // ...
8 }
```

netFD定义在 net/fd_posix.go:

```
1 // Network file descriptor.
2 type netFD struct {
3     pfd poll.FD
4
5     // immutable until Close
6     family    int
7     sotype     int
8     isConnected bool // handshake completed or use of association with
9             peer
10    net        string
11    laddr      Addr
12    raddr      Addr
13 }
```

其中同样包含poll.FD。

不同平台的 IO 多路复用能力不同，FreeBSD 的 kqueue 支持磁盘文件，Linux 的 epoll 不支持磁盘文件。

6.2 初始化文件描述符

创建文件描述符时设置为非阻塞 IO，并注册到 poller 上。

Go 语言的网络 IO 是非阻塞的，创建 socket 时设置。

```
1 // net/sock_posix.go
2 func socket(ctx context.Context, net string, family, sotype, proto int,
   ipv6only bool, laddr, raddr sockaddr, ctrlCtxFn func(context.
   Context, string, string, syscall.RawConn) error) (fd *netFD, err
   error) {
3     // sysSocket() 创建非阻塞 socket
4     s, err := sysSocket(family, sotype, proto)
5     if err != nil {
6         return nil, err
7     }
8     // ...
9     if laddr != nil && raddr == nil {
10        switch sotype {
11            case syscall.SOCK_STREAM, syscall.SOCK_SEQPACKET:
12                if err := fd.listenStream(ctx, laddr, listenerBacklog(),
13                    ctrlCtxFn); err != nil {
14                    fd.Close()
15                    return nil, err
16                }
17                return fd, nil
18            case syscall.SOCK_DGRAM:
19                if err := fd.listenDatagram(ctx, laddr, ctrlCtxFn); err !=
20                    nil {
21                    fd.Close()
22                    return nil, err
23                }
24                return fd, nil
25        }
26    }
27    if err := fd.dial(ctx, laddr, raddr, ctrlCtxFn); err != nil {
28        fd.Close()
29        return nil, err
30    }
31    return fd, nil
32 }
```

sysSocket() 创建非阻塞 socket。

```
1 // net/sock_cloexec.go
2
3 // Wrapper around the socket system call that marks the returned file
4 // descriptor as nonblocking and close-on-exec.
5 func sysSocket(family, sotype, proto int) (int, error) {
6     s, err := socketFunc(family, sotype|syscall.SOCK_NONBLOCK|syscall.
7         SOCK_CLOEXEC, proto)
8     if err != nil {
9         return -1, os.NewSyscallError("socket", err)
10    }
11    return s, nil
12 }
```

后续`fd.listenStream(...)`、`fd.listenDatagram(...)`和`fd.Dial(...)`将该`netFD`注册到`netpoller`。

`os.File`也同样试图创建非阻塞的文件描述符，但因为不同平台的IO多路复用能力不同，FreeBSD的`kqueue`支持磁盘IO，但Linux`epoll`只支持`socket`和`pipe`，所以创建`os.File`时需要判断文件是否可轮询。在可轮询的平台上（FreeBSD），`os.File`也是非阻塞的，被注册到`netpoller`上；在不可轮询的平台上（Linux），`os.File`是阻塞的，会阻塞线程。

`(*os.File).Create()`创建描述符并初始化`poll.FD`的调用链如下：

```
1 Create()
2     -> OpenFile()
3     -> openFileNonlog()
4     -> newFile() 根据平台设置非阻塞/阻塞 IO，初始化 poll.FD（阻塞则为空值），（非阻塞）并注册到 netpoller。
```

6.3 封装系统调用

Go 文件描述符的本质是`poll.FD`，自然其上的各种读写操作最终也归结于`poll.FD`上的读写操作。

以网络IO为例，查看`socket`上`Write()`的实现。`socket`实际上是一个`*netFD`，`(*netFD).Write()`的实现如下：

```
1 // net/fd_unix.go
2 func (fd *netFD) Write(p []byte) (nn int, err error) {
3     nn, err = fd.pfd.Write(p)
4     runtime.KeepAlive(fd)
5     return nn, wrapSyscallError(writeSyscallName, err)
6 }
```

查看`(*os.File).Write()`的实现，`(*os.File).Write()`包装了`(*os.File).write()`。

```
1 // os/file.go
2 func (f *File) Write(b []byte) (n int, err error) {
3     if err := f.checkValid("write"); err != nil {
4         return 0, err
5     }
6     n, e := f.write(b)
7     if n < 0 {
8         n = 0
9     }
10    if n != len(b) {
11        err = io.ErrShortWrite
12    }
13
14    epipecheck(f, e)
15
16    if e != nil {
17        err = f.wrapErr("write", e)
18    }
19
20    return n, err
21 }
```

(`*os.File`).`Write()`是平台无关的接口，(`*os.File`).`write()`是特定 OS 上的实现，Linux 的实现在 `os/file_posix.go` 中。

```
1 // os/file_posix.go
2 func (f *File) write(b []byte) (n int, err error) {
3     n, err = f.pfd.Write(b)
4     runtime.KeepAlive(f)
5     return n, err
6 }
```

可见，不论是磁盘文件还是网络 IO，最终都归结于 `poll.FD` 上的读写。(`*poll.FD`).`Write()` 定义如下：

```

1 // internal/poll/fd_unix.go
2 func (fd *FD) Write(p []byte) (int, error) {
3     if err := fd.writeLock(); err != nil {
4         return 0, err
5     }
6     defer fd.writeUnlock()
7     if err := fd.pd.prepareWrite(fd.isFile); err != nil {
8         return 0, err
9     }
10    var nn int
11    for {
12        max := len(p)
13        if fd.IsStream && max-nn > maxRW {
14            max = nn + maxRW
15        }
16        n, err := ignoringEINTRIO(syscall.Write, fd.Sysfd, p[nn:max])
17        if n > 0 {
18            nn += n
19        }
20        if nn == len(p) {
21            return nn, err
22        }
23        if err == syscall.EAGAIN && fd.pd.pollable() {
24            if err = fd.pd.waitWrite(fd.isFile); err == nil {
25                continue
26            }
27        }
28        if err != nil {
29            return nn, err
30        }
31        if n == 0 {
32            return nn, io.ErrUnexpectedEOF
33        }
34    }
35 }

```

(`*poll.FD`).`Write()`做了两件事:

1. 调用(`*poll.FD`).`prepareWrite()`注册此描述符到 netpoller;
2. 在循环中调用`syscall.Write`, 可轮询(非阻塞 IO) 文件描述符未就绪则调用(`*poll.FD`).`waitWrite()`等待(阻塞 goroutine)。

对于可轮询的文件描述符, (`*poll.FD`).`prepareWrite()`调用, `runtime_pollReset()`重置`pollDesc`的读写者为`pdNil`。

```

1 // internal/poll/fd_poll_runtime.go
2 func (pd *pollDesc) prepareWrite(isFile bool) error {
3     return pd.prepare('w', isFile)
4 }
5
6 // ...
7
8 func (pd *pollDesc) prepare(mode int, isFile bool) error {
9     if pd.runtimeCtx == 0 {
10         return nil
11     }
12     res := runtime_pollReset(pd.runtimeCtx, mode)
13     return convertErr(res, isFile)
14 }

```

internal/poll/fd_poll_runtime.go 中形如 `runtime_XXXXXX()` 的函数只有函数原型，没有函数体。

```

1 // internal/poll/fd_poll_runtime.go
2 func runtime_pollServerInit()
3 func runtime_pollOpen(fd uintptr) (uintptr, int)
4 func runtime_pollClose(ctx uintptr)
5 func runtime_pollWait(ctx uintptr, mode int) int
6 func runtime_pollWaitCanceled(ctx uintptr, mode int) int
7 func runtime_pollReset(ctx uintptr, mode int) int
8 func runtime_pollSetDeadline(ctx uintptr, d int64, mode int)
9 func runtime_pollUnblock(ctx uintptr)
10 func runtime_isPollServerDescriptor(fd uintptr) bool

```

这是因为这些函数都实现在 `runtime/netpoll.go` 中，符号名为 `poll_runtime_XXXXXX()`，通过 `go:linkname` 将其导出到 `internal/poll/fd_poll_runtime.go` 中，包含以下函数：

```

1 func poll_runtime_isPollServerDescriptor(fd uintptr) bool
2 func poll_runtime_pollOpen(fd uintptr) (*pollDesc, int)
3 func poll_runtime_pollClose(pd *pollDesc)
4 func poll_runtime_pollReset(pd *pollDesc, mode int) int
5 func poll_runtime_pollWait(pd *pollDesc, mode int) int
6 func poll_runtime_pollSetDeadline(pd *pollDesc, d int64, mode int)
7 func poll_runtime_pollUnblock(pd *pollDesc)
8 func poll_runtime_pollServerInit()

```

`go:linkname <source> <target>` 的作用是指示编译器在链接时使用符号名 `<source>` 替代 `<target>`。

`(*pollFD).Write()` 中的 `fd.pd.waitWrite(fd.isFile)` 实际上实在调用 `poll_runtime_pollWait()`。

```

1 // internal/poll/fd_poll_runtime.go
2 func (pd *pollDesc) waitWrite(isFile bool) error {
3     return pd.wait('w', isFile)
4 }
5 // ...
6 func (pd *pollDesc) wait(mode int, isFile bool) error {
7     if pd.runtimeCtx == 0 {
8         return errors.New("waiting for unsupported file type")
9     }
10    res := runtime_pollWait(pd.runtimeCtx, mode)
11    return convertErr(res, isFile)
12 }

```

`runtime_pollWait()` 实际上是定义在 `runtime/netpoll.go` 中的 `poll_runtime_pollWait()`。

`//go:linkname poll_runtime_pollWait internal/poll.runtime_pollWait` 指示编译器使用 `poll_runtime_pollWait` 作为符号 `internal/poll.runtime_pollWait`。

```

1 // poll_runtime_pollWait, which is internal/poll.runtime_pollWait,
2 // waits for a descriptor to be ready for reading or writing,
3 // according to mode, which is 'r' or 'w'.
4 // This returns an error code; the codes are defined above.
5 //
6 //go:linkname poll_runtime_pollWait internal/poll.runtime_pollWait
7 func poll_runtime_pollWait(pd *pollDesc, mode int) int {
8     errcode := netpollcheckerr(pd, int32(mode))
9     if errcode != pollNoError {
10         return errcode
11     }
12     // As for now only Solaris, illumos, and AIX use level-triggered IO
13     .
14     if GOOS == "solaris" || GOOS == "illumos" || GOOS == "aix" {
15         netpollarm(pd, mode)
16     }
17     for !netpollblock(pd, int32(mode), false) {
18         errcode = netpollcheckerr(pd, int32(mode))
19         if errcode != pollNoError {
20             return errcode
21         }
22         // Can happen if timeout has fired and unblocked us,
23         // but before we had a chance to run, timeout has been reset.
24         // Pretend it has not happened and retry.
25     }
26     return pollNoError
27 }

```

for 循环是唤醒 goroutine 的关键。假设 `netpollblock()` 阻塞 goroutine，`netpoll()` 唤醒该 goroutine，该 goroutine 退出 `netpollblock()` 返回 **true**，跳出 **for** 循环。

除了 `netpoll()` 以 IO notification 唤醒 goroutine，还可以通过关闭文件描述符（对应错

误`pollErrClosing`)、IO 超时（对应错误`pollErrTimeout`）跳出循环并返回错误。

`pollErrClosing`错误对应文件描述符被关闭，撤销还未就绪的 IO，调用链如下：

```
1 (*poll.FD).Close()
2     -> (*poll.FD).evit()
3         -> (*poll.FD).runtime_pollUnblock()
4             -> (*pollDesc).pushInfo() 设置 pollClosing 错误
5             -> (*pollDesc).netpollunblock() 修改 rg/wg 状态
6             -> (*pollDesc).netpollgoready() 恢复该 goroutine
```

对应 `goroutine` 从`netpollblock()`退出后，`netpollcheckerr()`返回错误，从`for`循环跳出，向上传播错误。

`pollErrTimeout`同理，但有一个定时器，定时器超时调用回调函数，回调函数中执行类似`pollErrClosing`的操作。

7 实例

剖析了 `netpoller` 源代码，再回过头看 TCP echo server：


```
1 import (
2     "fmt"
3     "io"
4     "log"
5     "net"
6 )
7
8 func worker(conn net.Conn) {
9     defer conn.Close()
10    b := make([]byte, 512)
11    for {
12        size, err := conn.Read(b)
13        if err == io.EOF {
14            break
15        }
16        if err != nil {
17            log.Fatal(err)
18        }
19        size, err = conn.Write(b[0:size])
20        if err != nil {
21            log.Fatal(err)
22        }
23    }
24 }
25
26 func main() {
27     listener, err := net.Listen("tcp", "127.0.0.1:8080")
28     if err != nil {
29         log.Fatal(err)
30     }
31     for {
32         conn, err := listener.Accept()
33         if err != nil {
34             log.Fatal(err)
35         }
36         go worker(conn)
37     }
38 }
```

`conn.Read()` 读取 socket，在 Go runtime 层面发生了什么？

1. `conn.Read(b)` 尝试读取 socket 上的数据，`(net.Conn).Read()` 最底层是 `(*internal/poll.FD).Read()`；
2. IO 就绪则系统调用 `read(2)` 阻塞线程，直到完成数据传输后返回；
3. IO 未就绪则系统调用 `read(2)` 返回 `EAGAIN`，调用 `(*internal/poll.FD).waitRead()` 阻塞（`gopark()`）该 goroutine；调用链为：

```
1 (*internal/poll.FD).Read()
2 -> (*internal/poll.pollDesc).waitRead()
3     -> (*internal/poll.pollDesc).wait()
4         -> internal/poll.runtime_pollWait()
5             -> runtime.netpollblock()
6                 -> runtime.gopark()
```

4. Go runtime 在适当的时机（GC、调度 goroutine 以及 sysmon goroutine）轮询（`netpoll()`）IO 事件，发现该 socket 上发生读事件 `EOLLIN`，设置该文件描述符的 `rg` 为 `pdReady`，恢复（`netpollgoready()`）该 goroutine。
5. `work goroutine` 恢复执行，退出 `runtime.netpollblock()`，退出 `internal/poll.runtime_pollWait()`，返回底层系统调用 `read(2)` 的返回值，`conn.Read(b)` 退出。