

记录Tun 透明代理的多种实现方式, 以及如何避免 routing loop #57

©pentaikulawo opened this issue Aug 1, 2021 · 7 comments

taikulawo commented Aug 1, 2021 • edited

Atun出来后再写入tun, 下次还会将自己所写入的packet读出来, 如果设置默认路由是tun网卡, 会导致死循环, 下文会介绍解决routing loop的多种方法

<https://www.kernel.org/doc/Documentation/networking/tun/tap.txt>

How does Virtual network device actually work ?
Virtual network device can be viewed as a simple Point-to-Point or Ethernet device, which instead of receiving packets from a physical media, receives them from user space program and instead of sending packets via physical media sends them to the user space program.

Let's say that you configured IPv6 on the tap0, then whenever the kernel sends an IPv6 packet to tap0, it is passed to the application (Vtun for example). The application encrypts, compresses and sends it to the other side over TCP or UDP. The application on the other side decompresses and decrypts the data received and writes the packet to the TAP device, the kernel handles the packet like it came from real physical device.

一共两个方法, read 和 write

read 从tun读取数据包

write 将数据包写入tun, tun直接得userspace的packet注入内核, 就好像内核刚从物理网卡读取出来一样

<https://github.com/0freezy/seeker/blob/5a1b83a24c48b696fb26cc3d540d0d3cd7159f1/README.adoc#%E6%8C%87%E5%AE%A3%E3%9A%9C%E6%8B%96%E6%9F%90%E7%BD%91%E6%AE%85%E8%87%B0%E4%B8%A3%E7%90%86>

<https://github.com/0freezy/seeker/blob/5a1b83a24c48b696fb26cc3d540d0d3cd7159f1/README.adoc#%E5%AE%A6%E6%87%E8%80%E4%85%E8%8F%90%E6>

除此之外有各种在虚拟机, 搞坏了还能还原。

NAT, 转发到socks server

1. tun读取ip packet, 本地启动TcpServer/UdpServer

2. 解析出ip src, dest 信息

3. 通过fake by ip, 利用NAT, 获取一个内部网段的唯一ip

4. 将数据包包的ip dest改为自身ip, 并将dest port改为TcpServer监听的端口

5. TcpServer/UdpServer收到请求, 根据预设的规则进行判断, 最后决定将数据发给 socks server还是直连

这里的关键是ip packet的dest port 指向了 TcpServer/UdpServer 监听的端口, 通过这种方式, 让 IP packet 重新走了一遍OS自身的 tcp/ip stack(好处是借用OS自己的协议栈 [shadowsocks/shadowsocks-swt#119 \(comment\)](#)) , 同时获取到 src ip 和 dest ip

实现佳处是通过 session 实现类似NAT的功能, packet走完OS自己的协议栈后, 能够还原出original destination.

看下面 Go 和 Rust 的代码, 代码是不是非常棒?

https://github.com/ydrew/kone/blob/3af6d647d7435134d6039172bd4dc8c5e9fd45f3/1tcp_relay.go#L144

https://github.com/wildeeper/seeker/blob/4973f6355d549207e3b454bbaa36d6d887fcd/tun_nat/src/lib.rs#L28

类似还有最近 (2021/08/10) shadowsocks-rust 搞的tun模式

随便在tun网段bind某个端口 tcp, daddr

<https://github.com/shadowsocks/shadowsocks-rust/blob/4287c2aab8f1826446b6888b58462b4342a0a53/crates/shadowsocks-service/src/local/tun/tcp.rs#L67>

拦截流量后直接connection, 找到就得 dest 改为 tcp_daddr, 这样我们不再需要自己的userspace tcp/ip stack, 而是复用OS自身的网络栈。

<https://github.com/shadowsocks/shadowsocks-rust/blob/4287c2aab8f1826446b6888b58462b4342a0a53/crates/shadowsocks-service/src/local/tun/tcp.rs#L155>

直接从ip packet构建TcpListener

leaf使用自己修改的wip, Wip packet直接封装出 TcpListener, 这里的重点是, TcpListener往dest port是自身, 而无法获取真正的 dest port, 另外TcpServer只能收到bind的端口的请求。

leaf修改wip代码后, 使得任意端口的流量都会流经TcpListener

https://github.com/wildeeper/leaf/blob/9497af9c348be2a2661a7f736a7628790835ddb/leaf/src/proxy/tun/netstack/tcp_listener_impl.rs#L57

现在clash和其他很多项目都用到leaf作者封装的wip, 开发出来的go-tun2socks

leaf也有nat, 但只给UDP用, 主要是fake by ip

https://github.com/wildeeper/leaf/blob/9497af9c348be2a2661a7f736a7628790835ddb/leaf/src/proxy/tun/netstack/stack_impl.rs#L268

为什么leaf只有UDP用到nat_managers?

应该只是业务必要罢了, 修改的 wip 已经能对 TcpListener 提供很好的支持, 所以只给UDP做了 fake by ip, 用到nat, 这有区别上面说的 NAT, 转发到 socks server, leaf 不依赖 OS 自己的 tcp/ip 重組形式, 也就不需要维护 session 来进行 NAT

如何避免重新读出刚写入tun设备的数据包

同时, 如果自己设置默认路由全部流量都发到tun设备, 自己read出来后交给userspace网络栈处理, 再将处理后的数据包写入tun, 之后的read就不会重新read出自己刚写入的包呢?

我理解会, 内核根据默认路由会重新路由, 但那样死循环了呀

我看这些利用tun的transparent proxy从未遇到过这问题, 就好奇是什么原理

seeker有这句话

== 指定 IP 或某网段走代理

修改路由表, 将希望走代理的 IP 或者网段路由到虚拟网卡, 如果使用了本机 socks5 代理, 则必须确保 socks5 不会添加入路由表的网段, 否则会死循环。

设置默认路由后绝大多数流量都会发到tun, 为了防止自己写到的网段的流量路由又回来, 需要再设置更短的路由表, 发包的ip dest 用这个ip网段, 这样就不会再回来。

Active Routes:				
Network Destination Netmask Gateway Interface Metric				
0.0.0.0	0.0.0.0	192.168.3.1	192.168.3.21	25
0.0.0.0	0.0.0.0	On-link	192.16.0.1	0
127.0.0.0	255.0.0.0	On-link	127.0.0.1	331
192.16.0.0	255.255.0.0	On-link	192.16.0.1	

clash为默认网关, clash使用wiredguard, 在windows上使用 SO_UNICAST, IP绑定 outgoing 网卡流量, 不会再走routing decision

透明代理的tcp listener都是寄生在ip之上, Atun没各收到数据写入Userspace stack, 自己的socks代理直接从wip处理后的数据里就能解析出socket, leaf就这么做的。

net stack写入tun

<https://github.com/wildeeper/leaf/blob/04ff9c5bb652c53e197066af863af6a083ecff/leaf/src/proxy/tun/inbound.rs#L113>

tun写入 net stack

<https://github.com/mellow-io/go-tun2socks/blob/6789c4be9d7c2915a73ba0d303d71a3a135e1574/cmd/tun2socks/main.go#L109>

https://github.com/mellow-io/mellow/blob/7716e54768ded3c4c480eb706646c6bbaac08/src/helper/linux/config_route4L1

```
# ip
CMD=ifl
# tun gateway
TUN_CMD=ifl
# 原本的默认网关
ORIG_CMD=ifl
# 原本的路由表及接口
# https://github.com/mellow-io/mellow/blob/7716e54768ded3c4c480eb706646c6bbaac08/src/main.go#L66
ORIG_CMD=ifl
# 新的路由表
# https://man7.org/linux/man-pages/man8/iproute.8.html
# https://github.com/mellow-io/mellow/blob/7716e54768ded3c4c480eb706646c6bbaac08/src/main.go#L79
ORIG_CMD=ifl
```

```
"CMD" route del default table main
# 注意: 这删除了 main table, 默认使用的路由表
"CMD" route add default via $TUN_CMD table main
# 注意: 这里加到default table, default table是默认路由表, main路由表会匹配到会使用default
"CMD" route add default via $ORIG_CMD dev $ORIG_CMD $ORIG_CMD table default
# 策略路由, 从原本的default interface 接口来的ip数据包转发到default表, 这样就能送出去流量, 不会死循环了
"CMD" rule add from $ORIG_CMD table default
# 看来wiredguard等入tun后并不会发出routing loop, SO_NGIFCONF+确保不会重新读出刚写入的数据包
# https://www.kernel.org/doc/Documentation/networking/routing.html
```

wiredguard 之前通知 go_wonnet 来避免routing loop, <https://lists.ozonwall.net/devel/2016/02/08/222>

但看类似防护为 ip rule 就够用, 不应该在内核做太多magic的功能(直接不走路由), 最后不了了之

wiredguard实现

<https://github.com/WireGuard/wiredguard-go>

<https://www.mellow.io/mellow#10>

https://www.policyrouting.org/iproute2_doc.html#x46

防止routing loop的方式

为需要直连的ip设置单独的路由(删除掉默认路由)

同学, 我又来了[无辜笑], 我看leaf的全局模式可以转发整个系统的流量, 最近自己在搞一个类似的透明代理工具, 有个问题想请教下

如果自己设置默认路由全部流量都发到tun设备, 自己read出来后交给userspace网络栈处理, 再将处理后的数据包写入tun, 之后的read就不会重新read出自己刚写入的包呢

我理解会, 内核根据默认路由会重新路由, 但那样死循环了呀

但我看leaf的全局模式transparent proxy从未遇到过这问题, 这是怎么做到的呢

回复: read出来不会直接写入tun, 是重新封装发包给VPN server

vpn server是单独加了一条路由, 而且发送包不是直接写tun包, 是直接通过tcp/udp发出, socket可以指定网卡发送

```
ip route del default
ip route add default dev wg0
ip route add 163.172.161.0/32 via 192.168.1.1 dev eth0
```

[The Classic Solutions](#)

为直连ip设置单独路由(不删除默认路由, 只覆盖)

默认路由的作用是匹配不到默认路由, 通过设置 0.0.0.0/0, 让这些路由总是先于 default 命中。

再对要直连的 ip(这里是 163.172.161.0) 设置单独的路由, 不需要删除原来的默认路由

```
ip route add 0.0.0.0/1 dev wg0
ip route add 128.0.0.1 dev wg0
ip route add 163.172.161.0/32 via 192.168.1.1 dev eth0
```

这样也叫做 0/1, 128/1 trick。

但这 trick 有局限[问题Q1](#)

推荐阅读

[Overriding The Default Route](#)

iptables

iptables REDIRECT

<https://github.com/jammyw/loopproxy>

iptables TPROXY

iptables with fwmark

iptables 配合 fwmark

https://flylib.com/books/en/2_783.1.50/

策略路由(ip rule with fwmark)

使用 ip rule 删除掉某个网卡流量

Rule-based Routing <https://www.wiredguard.com/netns/#routing-all-your-traffic>

ip rule <https://man7.org/linux/man-pages/man8/ip-rule.8.html>

https://www.reddit.com/r/WireGuard/comments/m8wntf/i_dont_understand_how_wgquick_adds_routes/prkmpw?utm_source=share&utm_medium=web2a&context=3

或者 rule 还有更强大的防止routing decision

这方法没用过, 但看若能行

[tailscale/tailscale#144](#)

namespace solution

<https://superuser.com/questions/166405/tun-device-how-to-avoid-routing-dead-loop-when-write-a-transparent-proxy>

bind before connect

bind之后connect, routing不会起作用. 这样就能解决设置默认网关后导致的 routing loop

通过测试 [leaf](#) 我已经能半十分确认上面这句话的正确性

if I bind an interface before to connect, Does that mean the connect for outgoing traffic will use that interface I bind without follow the routing decision?

[@victor](#) yes, if you bind() to an interface before connect()ing, the connection will go out through that interface. That is the whole point.

listen之前需要bind, 决定listen到哪个网卡.

如果客户端去connect, 在调用connect时bind会隐式发生, 你也可以主动bind before connect, 绕过路由选择, 强逼出流量使用某个network interface

<https://stackoverflow.com/a/4297381/7529562>

windows 平台

IP_UNICAST_IF
wireguard也依赖tun device, 但这tun和Unix上的tun不太像.

windows并没有tun的概念, 为弥补这块, win tunt 就是个 windows 内核驱动, 也是个userspace tunnel, 前者从内核拉取, 向内核注入packet, 后者将前者的数据包传回 userspace 处理

<https://git.zx2c4.com/wireguard-windows/about/docs/attacksurface.md#win tunt>

windows 并没有策略路由, 所以通过 bind interface IP_UNICAST_IF 来避免routing loop (从而印证了 bind before connect 可以解决routing loop)

以下三封邮件解释了为什么使用 bind.

<https://lists.zx2c4.com/pipermail/wireguard/2019-September/004493.html>

<https://lists.zx2c4.com/pipermail/wireguard/2019-September/004547.html>

<https://lists.zx2c4.com/pipermail/wireguard/2019-September/004542.html>

其中 <https://lists.zx2c4.com/pipermail/wireguard/2019-September/004541.html> 谈到的 IP_UNICAST_IF [IPv6_UNICAST_IF 可理解为指定 outgoing 流量的网络接口.

而里面说的 WFP 是 [Windows Filtering Platform](#), 看起来是 windows 平台内核级别的流量过滤API(类似 iptables)

windows 没有Linux 的类似 ip rule 这种策略路由. wireguard-windows 用来自 IP_UNICAST_IF 来开发

下面是 wireguard-go 在 windows 上用 IP_UNICAST_IF 的实现

wireguard-go IP_UNICAST_IF的实现

https://github.com/WireGuard/wireguard-go/blob/5845b6523837e04d0c35b15349ced7f666397520/core/bind_windows.go#L567

31是windows header里定义的

https://www.pinvoke.net/search.aspx?search=IP_UNICAST_IF&namespace=All

此外, wireguard 会帮你自动设置路由, bind 来避免routing loop

https://www.reddit.com/r/WireGuard/comments/m8uirt/i_dont_understand_how_wgquick_adds_routes/

wireguard 对自己已使用 IP_UNICAST_IF 的解释

<https://git.zx2c4.com/wireguard-windows/about/docs/network.md>

VPN service 直连支持被做进VPN的商由

<https://www.cisco.com/c/en/us/topo/papers/wireguard-networking-vpn-vpnrouteassignment-ipsec-discussions-view.html#wp1-22621>

Linux 平台

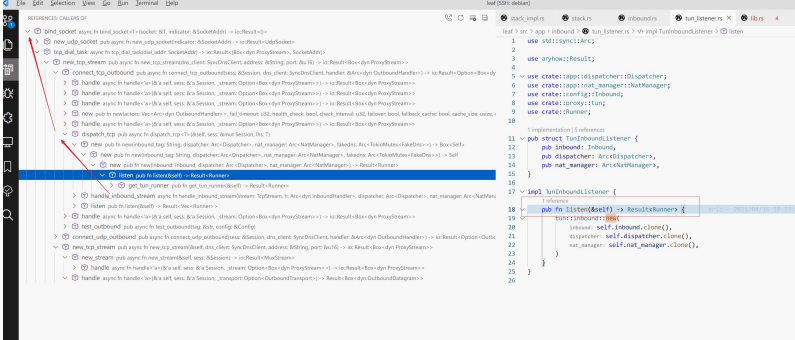
通过 setsockopt syscall 时传递 SO_BINDTODEVICE

<https://stackoverflow.com/questions/4584908/how-do-i-send-udp-packet-from-a-specific-interface-on-linux>

<https://lwn.net/Articles/1326657/132665717.4716.4.camel@bedumazet-laptop/?>

阅读代码时确认使用过上述选项

<https://github.com/willdeeper/leaf/blob/0cf00c9db6552c53e19706fa863af6a083ecf/leaf/src/proxy/mod.rs#L235>



shadowsocks-rust 新添加的 tun mode 也使用了

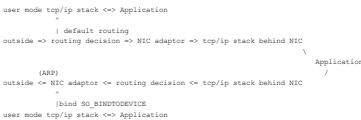
<https://github.com/shadowsocks/shadowsocks-rust/blob/4287c2aa8df1826446b6888f584620b432a0a53/bin/main.rs#L116>

<https://github.com/shadowsocks/shadowsocks-rust/pull/586#issuecomment-999196940>

经测试(leaf bind, 不添加 ip rule), bind SO_BINDTODEVICE 之后再将数据写入socket, 数据会直接到达网卡的发送队列, 不会再次走routing decision.

路由选择的核心在于找到一个网卡, 并 bind SO_BINDTODEVICE 直接绕过了这一步.

对于从外网接收数据.



当 leaf 使用bind到默认网卡时, 不需要ip rule策略路由.

即使socks代理在本地, 也不会导致routing loop.

这是因为路由表中 local 表优先级最高, 我们修改的是main和default表.

配置 socks outbound 为其他机器时正常工作, 是因为流量到了原始网卡直接发其他机器了.

但如果配置的socks outbound 为本地环回地址又会问题

例子

本地开clash, listen 7890, leaf配置 socks outbound 127.0.0.1:7890

socket SO_BINDTODEVICE 到 eth0, 用这 socket dial 到 127.0.0.1:7890

在vmm8, 也就是虚拟机的默认网卡上抓到 ip dst 为 127.0.0.1 的流量, 127.0.0.1属于环回地址, 不应该在该以太网卡上抓到, 所以连接失败.

No.	Time	Source	Source Port	Destination	Dest Port	Protocol	Length	Info
52379	1880.460375	192.168.58.3		59874	127.0.0.1	7890 TCP	74	59874 → 7890 [SYN] Seq=0 Win=64240 Len=0 MSS=1460 SACK_PERM=1 TSval=1347645507 TSecr=0
52380	1887.494284	127.0.0.1		7890	192.168.58.3	59874 TCP	58	7890 → 59874 [SYN, ACK] Seq=0 Ack=1 Win=64240 Len=0 MSS=1460
52404	1887.194127	127.0.0.1		7890	192.168.58.3	59874 TCP	58	(TCP Retransmission) 7890 → 59874 [SYN, ACK] Seq=0 Ack=1 Win=64240 Len=0 MSS=1460
52412	1887.204720	127.0.0.1		7890	192.168.58.3	59874 TCP	58	(TCP Retransmission) 7890 → 59874 [SYN, ACK] Seq=0 Ack=1 Win=64240 Len=0 MSS=1460
52413	1887.395647	127.0.0.1		7890	192.168.58.3	59874 TCP	58	(TCP Retransmission) 7890 → 59874 [SYN, ACK] Seq=0 Ack=1 Win=64240 Len=0 MSS=1460
52422	1888.102096	192.168.58.3		59874	127.0.0.1	7890 TCP	74	(TCP Retransmission) 59874 → 7890 [SYN] Seq=0 Win=64240 Len=0 MSS=1460 SACK_PERM=1
52423	1888.395765	127.0.0.1		7890	192.168.58.3	59874 TCP	58	(TCP Retransmission) 7890 → 59874 [SYN, ACK] Seq=0 Ack=1 Win=64240 Len=0 MSS=1460
52440	1878.188887	192.168.58.3		59874	127.0.0.1	7890 TCP	74	(TCP Retransmission) 59874 → 7890 [SYN] Seq=0 Win=64240 Len=0 MSS=1460 SACK_PERM=1
52441	1878.395974	127.0.0.1		7890	192.168.58.3	59874 TCP	58	(TCP Retransmission) 7890 → 59874 [SYN, ACK] Seq=0 Ack=1 Win=64240 Len=0 MSS=1460
52446	1878.396136	127.0.0.1		7890	192.168.58.3	59874 TCP	58	(TCP Retransmission) 7890 → 59874 [SYN, ACK] Seq=0 Ack=1 Win=64240 Len=0 MSS=1460
52372	1882.396127	127.0.0.1		7890	192.168.58.3	59874 TCP	58	(TCP Retransmission) 7890 → 59874 [SYN, ACK] Seq=0 Ack=1 Win=64240 Len=0 MSS=1460

clash issue <https://github.com/Dreamacro/clash/issues/135> 谈到了bind, listen 的问题

Mac

Network Extension

excludedRoutes

<https://developer.apple.com/documentation/networkextension/nepolicy/settings/1404294-excludedroutes>

wireguard 就使用到了 includedRoutes

<https://github.com/WireGuard/wireguard-apple/blob/23618994f718a48f2465479b4a1e8a0830b334/Source/WireGuardKit/PacketTunnelSettingsGenerator.swift#L117>

<https://www.v2ex.com/t/590555#r-9004049>

设置路由表

Android [https://developer.android.com/reference/android/net/VpnService#protect\(java.net.Socket\)](https://developer.android.com/reference/android/net/VpnService#protect(java.net.Socket))

Protect a socket from VPN connections. After protecting, data sent through this socket will go directly to the underlying network, so its traffic will not be forwarded through the VPN.

每个OS都有自己的特色. Linux 下提供 SO_BINDTODEVICE, 而Mac, Android则提供高级的API, 用于解决routing loop这个问题. 核心都是绕过了 routing decision 这项

总结

每个平台都有特殊的实现方式.

Linux

- iptables或者配合fwmark: 很多软件
- ip rule: 很多软件

Mac

- network extension: wireguard-apple
- 为libproxy单独设置路由: 经真VPN实现, VPN server挂有公网IP

Android

- VpnServiceProtect: shadowsocks-android
- IP_UNICAST_IF: wireguard-go

而又有最后相似的, 古典的(classic) 处理方式:

libproxy创建半独的route, 避免 routing loop. <https://www.wireguard.com/netns/#the-classic-solution>

参考的开源项目

C实现的tun2socks

<https://github.com/nussid/tunsocks>

Go实现的tun2socks, 支持windows

<https://github.com/mika-linux-network/Tun-2-Socks>

Rust实现的透明代理工具

<https://github.com/willdeeper/leaf>

