

状态机制

Table of Contents

- 概述
- conntrack记录
- 数据包在用户空间的状态
 - TCP 连接
 - UDP连接
 - ICMP 连接
- 缺省的连接操作
- 复杂协议和连接跟踪

本章将详细介绍状态机制

概述

状态机制是 iptables 中特殊的一部分，其实它不应该叫状态机制，因为它只是一种 **连接跟踪** 机制。但是，很多人都认可状态机制这个名字。也或多或少地用这个名字来表示和连接跟踪相同的意思。这不应该引起什么混乱的。连接跟踪可以让 Netfilter 知道某个特定连接的状态。运行连接跟踪的防火墙称作 **带有状态机制** 的防火墙，以下简称 **状态防火墙**。状态防火墙比非状态防火墙要安全，因为它允许编写更严格的规则。

在iptables里，包分为被跟踪连接的四种不同状态有关的。它们是 **NEW**，**ESTABLISHED**，**RELATED** 和 **INVALID**。后面会深入讨论每一个状态。使用 **state** 匹配操作，能很容易地控制 **谁或什么能发起新的会话**

所有在内核中由 Netfilter 的特定框架做的连接跟踪称作 **conntrack**。conntrack 可以作为模块安装，也可以作为内核的一部分。大部分情况下，需要更详细的连接跟踪。这是相比于缺省的conntrack而言。也因此，conntrack中有许多用来处理 **ICMP**，**UDP** 或 **ICMP** 协议的部分。这些模块从数据包中提取详细的、唯一的 **信息**，因此能保持对每一个数据流的跟踪。这些信息也告知 conntrack 流当前的状态。例如，UDP 流一般由他们的 **目的地地址**、**源地址**、**目的端口** 和 **源端口** 唯一确定

在以前的内核里，可以打开或关闭重组功能。然而，自从iptables和Netfilter，尤其是连接跟踪被引入内核，这个选项就被取消了。因为如果没有的重组，连接跟踪就不能正常工作。现在重组已经整合入 conntrack，并且在conntrack启动时自动启动

不要关闭重组功能，除非你要关闭连接跟踪

除了本地产生的包由 **OUTPUT** 链处处理外，所有连接跟踪都是在 **PREROUTING** 链里进行处理的。iptables会在PREROUTING链里从新计算所有的状态：

- 如果发送一个流的初始化包，状态就会在OUTPUT链里被设置为 **NEW**
 - 当收到回应的包时，状态就会在PREROUTING链里被设置为 **ESTABLISHED**
 - 如果第一个包不是本地产生的，那就会在PREROUTING链里被设置为 **RELATED** 状态
- 综上，所有状态的变化和计算都是在 **nat** 表中的 **PREROUTING** 链和 **OUTPUT** 链里完成的

conntrack记录

先来看看怎样阅读 /proc/net/ip_conntrack 里的conntrack记录。这些记录表示的是当前被跟踪的连接。如果安装了 ip_conntrack 模块：

```
$ cat /proc/net/ip_conntrack

tcp        6 117 SYN_SENT src=192.168.1.6 dst=192.168.1.9 sport=32775 \
dport=22 [UNREPLIED] src=192.168.1.9 dst=192.168.1.6 sport=22 \
dport=32775 use=2
```

首先显示的是协议，这里是 **tcp**，接着是十进制的6（**tcp**的协议类型代码是6）。之后的117是 **这条conntrack记录的生存时间**，它会有规律地被消耗，直到收到这个连接的更多的包。那时，这个值就会被设为当时那个状态的缺省值。接下来的是这个连接在当前时间点的状态。上面的例子说明这个包处在状态 **SYN_SENT**，这个值是iptables显示的，以便好理解，而内部用的值稍有所不同。**SYN_SENT**说明正在观察的这个连接只在一个方面发送了**一个TCP SYN包**。再下面是 **源地址**、**目的地地址** 和 **目的端口**。其中有个特殊的词 **UNREPLIED**，说明 **这个连接还没有收到任何回应**。最后，是希望接收的应答包的信息，他们的地址和端口和前面是相反的

连接跟踪记录的信息依据IP所包含的协议不同而不同，所有相应的值都是在头文件linux/include/netfilter-ipv4/ip_conntrack.h中定义的。IP、TCP、UDP、ICMP协议的缺省值是在linux/include/netfilter-ipv4/ip_conntrack.h里定义的。具体的值可以查看相应的协议，但这里用不到它们，因为它们大都只在conntrack内部使用。随着状态的变化，生存时间也会改变

数据包在用户空间的状态

包的状态依据IP所包含的协议不同而不同，但在内核外部，也就是用户空间里，只有4种状态：**NEW**，**ESTABLISHED**，**RELATED** 和 **INVALID**。它们主要是和状态匹配一起使用。下面简要地介绍以下几种状态：

Table 1: 数据包在用户空间的状态

State	Comment
NEW	NEW说明这个包是我们看到第一个包。意思就是，这是 conntrack模块看到的某个连接第一个包 ，它即将被匹配了。比如，我们看到一个SYN包，是我们所感兴趣的连接第一个包，就要匹配它。第一个包也可能不是SYN包，但它仍会被认为是NEW状态。这样有时会引发一些问题，但对某些情况是有非常大的帮助的。例如，在我们想恢复某条从其他防火墙丢失的连接时，或者某个连接已经超时，但实际上并未关闭时
ESTABLISHED	ESTABLISHED已经注意到两个方向上的数据传输，而且会继续匹配这个连接的包。处于ESTABLISHED状态的连接是非常容易理解的。ftp是个很好的例子，FTP-data连接就是和FTP-control有RELATED的。还有一个连接要从NEW变为ESTABLISHED， 只需要接到应答包即可，不管这个包是发往防火墙的，还是由防火墙转发的 。ICMP的错误和重定向等信息包也被看作是ESTABLISHED，只要它们是我们所发出的信息的应答
RELATED	RELATED是个比较麻烦的状态。当一个连接和某个已处于ESTABLISHED状态的连接有关系时，就被认为是RELATED的了。换句话说，一个连接要是RELATED的， 首先要有一个ESTABLISHED的连接。这个ESTABLISHED的连接再产生一个主连接之外的连接 ，这个新的连接就是RELATED的了，当然前提是conntrack模块要能理解RELATED。ftp是个很好的例子，FTP-data连接就是和FTP-control有RELATED的。还有其他的例子，比如，通过IRC的DCC连接。有了这个状态，ICMP应答、FTP传输、DCC等才能穿过防火墙正常工作。注意，大部分还有一些UDP协议都依赖这个机制。这些协议是很复杂的，它们把连接信息放在数据包里，并且要求这些信息能被正确理解
INVALID	INVALID说明 数据包不能被识别属于哪个连接或没有任何状态 。有几个原因可以产生这种情况，比如， 内存溢出 ， 收到不属于这个连接的ICMP错误信息 。一般地，我们DROP这个状态的任何东西

这些状态可以一起使用，以便匹配数据包。这可以使防火墙非常强壮和有效。以前，我们经常打开1024以上的所有端口来放行应答的数据。现在，有了状态机制，就不需要这样了。可以只开放那些有应答数据的端口，其他的都可以关闭。这样就安全多了

TCP 连接

一个TCP连接是经过三次握手协商连接信息才建立起来的。整个会话由一个 **SYN** 包开始，然后是一个 **SYN/ACK** 包，最后是一个 **ACK** 包，此时，会话才建立成功，能够发送数据。最大的问题在于 **连接跟踪怎样控制这个过程**

默认情况下，连接跟踪基本上对所有的连接要做同样的操作。看看下面的图，就能明白在连接的不同阶段，流是处于什么状态的。就如你看到的，连接跟踪的代码是从用户的观点来看待TCP连接建立的流程的。连接跟踪一看到 **SYN** 包，就认为这个连接是 **NEW** 状态，一看到返回的 **SYN/ACK** 包，就认为连接是 **ESTABLISHED** 状态



如果仔细想想第二步，应该能理解为什么。有了这个特殊处理：

- NEW** 和 **ESTABLISHED** 包 就可以发出本地网络
- 只有 **ESTABLISHED** 的连接才能有回应信息

如果把整个建立连接的过程中传输的数据包都看作NEW，那么三次握手用的包都是NEW状态的，这样我们就不能阻止从外部到本防火墙即使连接是从外到内的，但它使用的包也是NEW状态的，而且为了其他连接能正常工作，不得不允许NEW状态的包返回并进入。更复杂的是，针对TCP内核使用了很多内部状态，它们的定义在 RFC 793 - Transmission Control Protocol的2.1-2.5节中

以用户的观点来看，这是很简单的。但是，从内核的角度看这一块还有点困难的。来看一个例子。认真考虑一下在 /proc/net/ip_conntrack 里，连接的状态是如何改变的：

```
tcp        6 117 SYN_SENT src=192.168.1.5 dst=192.168.1.35 sport=1031 \
dport=23 [UNREPLIED] src=192.168.1.35 dst=192.168.1.5 sport=23 \
dport=1031 use=1
```

从上面的记录可以看出，**SYN_SENT** 状态被设置了，这说明连接已经发出一个 **SYN** 包，但应答还没发过来，这可从 **[UNREPLIED]** 标志看出

```
tcp        6 57 SYN_RECV src=192.168.1.5 dst=192.168.1.35 sport=1031 \
dport=23 src=192.168.1.35 dst=192.168.1.5 sport=23 dport=1031 \
use=1
```

现在已经收到了相应的 **SYN/ACK** 包，状态也变为 **SYN_RECV**，这说明最初发出的 **SYN** 包已正确传输，并且 **SYN/ACK** 包也到达了防火墙。这就意味着在连接的双方都有数据传输，因此可以认为两个方向都有相应的回应。当然这是假设的

```
tcp        6 431999 ESTABLISHED src=192.168.1.5 dst=192.168.1.35 \
sport=1031 dport=23 src=192.168.1.35 dst=192.168.1.5 \
sport=23 dport=1031 use=1
```

现在发出了三步握手的最后一个包，即 **ACK** 包，连接也就进入 **ESTABLISHED** 状态了。再传输几个数据包，连接就是 **[ASSURED]** 的了

下面介绍TCP连接在关闭过程中的状态：

