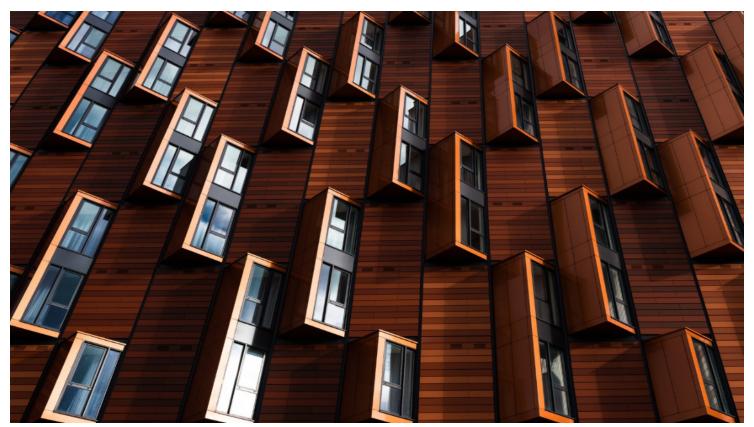
第3讲 | ifconfig: 最熟悉又陌生的命令行

笔记本: P.趣谈网络协议 **创建时间:** 2018/5/23 9:14

作者: hongfenghuoju

第3讲 | ifconfig: 最熟悉又陌生的命令行

2018-05-23 刘超



更新时间:

2018/5/23 9:15

上一节结尾给你留的一个思考题是,你知道怎么查看 IP 地址吗?

当面试听到这个问题的时候,面试者常常会觉得走错了房间。我面试的是技术岗位啊,怎么问这么简单的问题?

的确,即便没有专业学过计算机的人,只要倒腾过电脑,重装过系统,大多也会知道这个问题的答案:在 Windows 上是 ipconfig,在 Linux 上是 ifconfig。

那你知道在 Linux 上还有什么其他命令可以查看 IP 地址吗?答案是 ip addr。如果回答不上来这个问题,那你可能没怎么用过 Linux。

那你知道 ifconfig 和 ip addr 的区别吗?这是一个有关 net-tools 和 iproute2 的 "历史"故事,你刚来到第三节,暂时不用了解这么细,但这也是一个常考的知识点。

想象一下,你登录进入一个被裁剪过的非常小的 Linux 系统中,发现既没有 ifconfig 命令,也没有 ip addr 命令,你是不是感觉这个系统压根儿没法用?这个时候,你可以自行安装 net-tools 和 iproute2 这两个工具。当然,大多数时候这两个命令是系统自带的。

安装好后,我们来运行一下 ip addr。不出意外,应该会输出下面的内容。

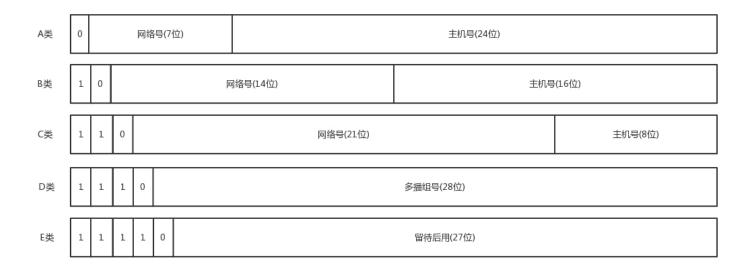
```
root@test:~# ip addr
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN group default
    link/loopback 00:00:00:00:00 brd 00:00:00:00:00
    inet 127.0.0.1/8 scope host lo
        valid_lft forever preferred_lft forever
    inet6 ::1/128 scope host
        valid_lft forever preferred_lft forever
2: eth0: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc pfifo_fast state UP group default qlen 1000
    link/ether fa:16:3e:C7:79:75 brd ff:ff:ff:ff:ff
    inet 10.100.122.2/24 brd 10.100.122.255 scope global eth0
        valid_lft forever preferred_lft forever
    inet6 fe80::f816:3eff:fec7:7975/64 scope link
        valid_lft forever preferred_lft forever
```

这个命令显示了这台机器上所有的网卡。大部分的网卡都会有一个 IP 地址, 当然, 这不是必须的。在后面的分享中, 我们会遇到没有 IP 地址的情况。

IP 地址是一个网卡在网络世界的通讯地址,相当于我们现实世界的门牌号码。既然是门牌号码,不能大家都一样,不然就会起冲突。比方说,假如大家都叫六单元 1001 号,那快递就找不到地方了。所以,有时候咱们的电脑弹出网络地址冲突,出现上不去网的情况,多半是 IP 地址冲突了。

如上输出的结果,<mark>10.100.122.2 就是一个 IP 地址</mark>。<mark>这个地址被点分隔为四个部分,每个部分 8 个 bit, 所以 IP 地址总共是 32 位</mark>。这样产生的 IP 地址的数量很快就不够用了。因为当时设计 IP 地址的时候,哪知道今天会有这么多的计算机啊!因为不够用,于是就有了 IPv6,也就是上面输出结果里面 inet6 fe80::f816:3eff:fec7:7975/64。这个有 128 位,现在看来是够了,但是未来的事情谁知道呢?

本来 32 位的 IP 地址就不够,还被分成了 5 类。现在想想,当时分配地址的时候,真是太奢侈了。



在网络地址中,至少在当时设计的时候,<mark>对于 A、B、 C 类主要分两部分</mark>,<mark>前面一部分是网络号,后面一部分是主机号。</mark>这很好理解,大家都是六单元 1001 号,我是小区 A 的六单元 1001 号,而你是小区 B 的六单元 1001 号。

下面这个表格,详细地展示了 A、B、C 三类地址所能包含的主机的数量。在后文中,我也会多次借助这个表格来讲解。

类别	IP地址范围	最大主机数	私有IP地址范围
Α	0.0.0.0-127.255.255.255	16777214	10.0.0.0-10.255.255.255
В	128.0.0.0-191.255.255.255	65534	172.16.0.0-172.31.255.255
С	192.0.0.0-223.255.255.255	254	192.168.0.0-192.168.255.255

这里面有个尴尬的事情,就是 C 类地址能包含的最大主机数量实在太少了,只有 254 个。当时设计的时候恐怕没想到,现在估计一个网吧都不够用吧。而 B 类地址能包含的最大主机数量又太多了。6 万多台机器放在一个网络下面,一般的企业基本达不到这个规模,闲着的地址就是浪费。

无类型域间选路 (CIDR)

于是有了一个折中的方式叫作无类型域间选路,简称CIDR。这种方式打破了原来设计的几类地址的做法,将32位的IP地址一分为二,前面是网络号,后面是主机号。从哪里分呢?你如果注意观察的话可以看到,10.100.122.2/24,这个IP地址中有一个斜杠,斜杠后面有个数字24。这种地址表示形式,就是CIDR。后面24的意思是,32位中,前24位是网络号,后8位是主机号。

伴随着 CIDR 存在的,一个是<mark>广播地址</mark>,<mark>10.100.122.255。如果发送这个地址,所有 10.100.122 网络里面的机器都可以收到。另一个是子网掩码,255.255.255.0。</mark>

将子网掩码和 IP 地址进行 AND 计算。前面三个 255, 转成二进制都是 1。1 和任何数值取 AND, 都是原来数值,因而前三个数不变,为 10.100.122。后面一个 0,转换成二进制是 0,0 和任何数值取 AND,都是 0,因而最后一个数变为 0,合起来就是 10.100.122.0。这就是网络号。将子网掩码和 IP 地址按位计算 AND,就可得到网络号。

公有 IP 地址和私有 IP 地址

在日常的工作中,几乎不用划分 A 类、B 类或者 C 类,所以时间长了,很多人就忘记了这个分类,而只记得 CIDR。但是有一点还是要注意的,就是公有 IP 地址和私有 IP 地址。

类别	IP地址范围	最大主机数	私有IP地址范围
Α	0.0.0.0-127.255.255.255	16777214	10.0.0.0-10.255.255.255
В	128.0.0.0-191.255.255.255	65534	172.16.0.0-172.31.255.255
С	192.0.0.0-223.255.255.255	254	192.168.0.0-192.168.255.255

我们继续看上面的表格。表格最右列是私有 IP 地址段。平时我们看到的数据中心里,办公室、家里或学校的 IP 地址,一般都是私有 IP 地址段。因为这些地址允许组织内部的 IT 人员自己管理、自己分配,而且可以重复。因此,你学校的某个私有 IP 地址段和我学校的可以是一样的。

这就像每个小区有自己的楼编号和门牌号,你们小区可以叫 6 栋,我们小区也叫 6 栋,没有任何问题。但是一旦出了小区,就需要使用公有 IP 地址。就像人民路 888 号,是国家统一分配的,不能两个小区都叫人民路 888 号。

公有 IP 地址有个组织统一分配,你需要去买。如果你搭建一个网站,给你学校的人使用,让你们学校的 IT 人员给你一个 IP 地址就行。但是假如你要做一个类似网易 163 这样的网站,就需要有公有 IP 地址,这样全世界的人才能访问。

表格中的 192.168.0.x 是最常用的私有 IP 地址。你家里有 Wi-Fi,对应就会有一个 IP 地址。一般你家里地上网设备不会超过 256 个,所以 /24 基本就够了。有时候我们也能见到 /16 的 CIDR,这两种是最常见的,也是最容易理解的。

不需要将十进制转换为二进制 32 位,就能明显看出 192.168.0 是网络号,后面是主机号。而整个网络里面的第一个地址 192.168.0.1,往往就是你这个私有网络的出口地址。例如,你家里的电脑连接 Wi-Fi, Wi-Fi 路由器的地址就是 192.168.0.1,而 192.168.0.255 就是广播地址。一旦发送这个地址,整个192.168.0 网络里面的所有机器都能收到。

但是也不总都是这样的情况。因此,其他情况往往就会很难理解,还容易出错。

举例:一个容易"犯错"的 CIDR

我们来看 16.158.165.91/22 这个 CIDR。 求一下这个网络的第一个地址、子网掩码和广播地址。

你要是上来就写 16.158.165.1, 那就大错特错了。

/22 不是 8 的整数倍,不好办,只能先变成二进制来看。16.158 的部分不会动,它占了前 16 位。中间的 165,变为二进制为10100101。除了前面的 16 位,还剩 6 位。所以,这 8 位中前 6 位是网络号,16.158.<101001>,而<01>.91 是机器号。

第一个地址是 16.158.<101001><00>.1,即 16.158.164.1。子网掩码是 255.255.<111111><00>.0,即 255.255.252.0。广播地址为 16.158.<101001><11>.255,即 16.158.167.255。

这五类地址中,还有一类 D 类是组播地址。使用这一类地址,属于某个组的机器都能收到。这有点类似在公司里面大家都加入了一个邮件组。发送邮件,加入这个组的都能收到。组播地址在后面讲述 VXLAN协议的时候会提到。

讲了这么多,才讲了上面的输出结果中很小的一部分,是不是觉得原来并没有真的理解 ip addr 呢?我们接着来分析。

在 IP 地址的后面有个 <mark>scope,</mark>对于 eth0 这张网卡来讲,<mark>是 global,说明这张网卡是可以对外的,可以</mark> 接收来自各个地方的包。对于 lo 来讲,是 host,说明这张网卡仅仅可以供本机相互通信。

lo 全称是loopback,又称环回接口,往往会被分配到 127.0.0.1 这个地址。这个地址用于本机通信,经过内核处理后直接返回,不会在任何网络中出现。

MAC 地址

在 IP 地址的上一行是 link/ether fa:16:3e:c7:79:75 brd ff:ff:ff:ff:ff; 这个被称为MAC 地址,是一个网卡的物理地址,用十六进制,6 个 byte 表示。

MAC 地址是一个很容易让人"误解"的地址。因为 MAC 地址号称全局唯一,不会有两个网卡有相同的 MAC 地址,而且网卡自生产出来,就带着这个地址。很多人看到这里就会想,既然这样,整个互联网的通信,全部用 MAC 地址好了,只要知道了对方的 MAC 地址,就可以把信息传过去。

这样当然是不行的。一个网络包要从一个地方传到另一个地方,除了要有确定的地址,还需要有定位功能。而有门牌号码属性的 IP 地址,才是有远程定位功能的。

例如,你去杭州市网商路 599 号 B 楼 6 层找刘超,你在路上问路,可能被问的人不知道 B 楼是哪个,但是可以给你指网商路怎么去。但是如果你问一个人,你知道这个身份证号的人在哪里吗?可想而知,没有人知道。

MAC 地址更像是身份证,是一个唯一的标识。它的唯一性设计是为了组网的时候,不同的网卡放在一个网络里面的时候,可以不用担心冲突。从硬件角度,保证不同的网卡有不同的标识。

MAC 地址是有一定定位功能的,只不过范围非常有限。你可以根据 IP 地址,找到杭州市网商路 599 号 B 楼 6 层,但是依然找不到我,你就可以靠吼了,大声喊身份证 XXXX 的是哪位?我听到了,我就会站起来说,是我啊。但是如果你在上海,到处喊身份证 XXXX 的是哪位,我不在现场,当然不会回答,因为我在杭州不在上海。

所以,MAC 地址的通信范围比较小,局限在一个子网里面。例如,从 192.168.0.2/24 访问 192.168.0.3/24 是可以用 MAC 地址的。一旦跨子网,即从 192.168.0.2/24 到 192.168.1.2/24,MAC 地址就不行了,需要 IP 地址起作用了。

网络设备的状态标识

解析完了 MAC 地址,我们再来看 < BROADCAST,MULTICAST,UP,LOWER_UP > 是干什么的?这个叫作net_device flags,网络设备的状态标识。

UP 表示网卡处于启动的状态;BROADCAST 表示这个网卡有广播地址,可以发送广播包;MULTICAST表示网卡可以发送多播包;LOWER_UP 表示 L1 是启动的,也即网线插着呢。MTU1500 是指什么意思呢?是哪一层的概念呢?最大传输单元 MTU 为 1500,这是以太网的默认值。

上一节,我们讲过网络包是层层封装的。MTU 是二层 MAC 层的概念。MAC 层有 MAC 的头,以太网规定连 MAC 头带正文合起来,不允许超过 1500 个字节。正文里面有 IP 的头、TCP 的头、HTTP 的头。如果放不下,就需要分片来传输。

qdisc pfifo_fast 是什么意思呢? qdisc 全称是queueing discipline, 中文叫排队规则。内核如果需要通过某个网络接口发送数据包,它都需要按照为这个接口配置的 qdisc (排队规则) 把数据包加入队列。

最简单的 qdisc 是 pfifo,它不对进入的数据包做任何的处理,数据包采用先入先出的方式通过队列。 pfifo fast 稍微复杂一些,它的队列包括三个波段(band)。在每个波段里面,使用先进先出规则。

三个波段 (band) 的优先级也不相同。band 0 的优先级最高, band 2 的最低。如果 band 0 里面有数据包,系统就不会处理 band 1 里面的数据包,band 1 和 band 2 之间也是一样。

数据包是按照服务类型(Type of Service,TOS)被分配多三个波段(band)里面的。TOS 是 IP 头里面的一个字段,代表了当前的包是高优先级的,还是低优先级的。

队列是个好东西,后面我们讲云计算中的网络的时候,会有很多用户共享一个网络出口的情况,这个时候如何排队,每个队列有多粗,队列处理速度应该怎么提升,我都会详细为你讲解。

小结

怎么样,看起来很简单的一个命令,里面学问很大吧?通过这一节,希望你能记住以下的知识点,后面都能用得上:

• IP 是地址, 有定位功能; MAC 是身份证, 无定位功能;

- CIDR 可以用来判断是不是本地人;
- IP 分公有的 IP 和私有的 IP。后面的章节中我会谈到"出国门",就与这个有关。

最后,给你留两个思考题。

- 1. 你知道 net-tools 和 iproute2 的 "历史" 故事吗?
- 2. 这一节讲的是如何查看 IP 地址, 那你知道 IP 地址是怎么来的吗?