

## 1.tcp\_ip\_model

### 2.1 TCP/IP 网络模型有哪几层？

问大家，为什么要有 TCP/IP 网络模型？

对于同一台设备上的进程间通信，有很多种方式，比如有管道、消息队列、共享内存、信号等方式，而对于不同设备上的进程间通信，就需要网络通信，而设备是多样性的，所以要兼容多种多样的设备，就协商出了一套**通用的网络协议**。

这个网络协议是分层的，每一层都有各自的作用和职责，接下来就根据「TCP/IP 网络模型」分别对每一层进行介绍。

#### 应用层

最上层的，也是我们能直接接触到的就是**应用层**（*Application Layer*），我们电脑或手机使用的应用软件都是在应用层实现。那么，当两个不同设备的应用需要通信的时候，应用就把应用数据传给下一层，也就是传输层。

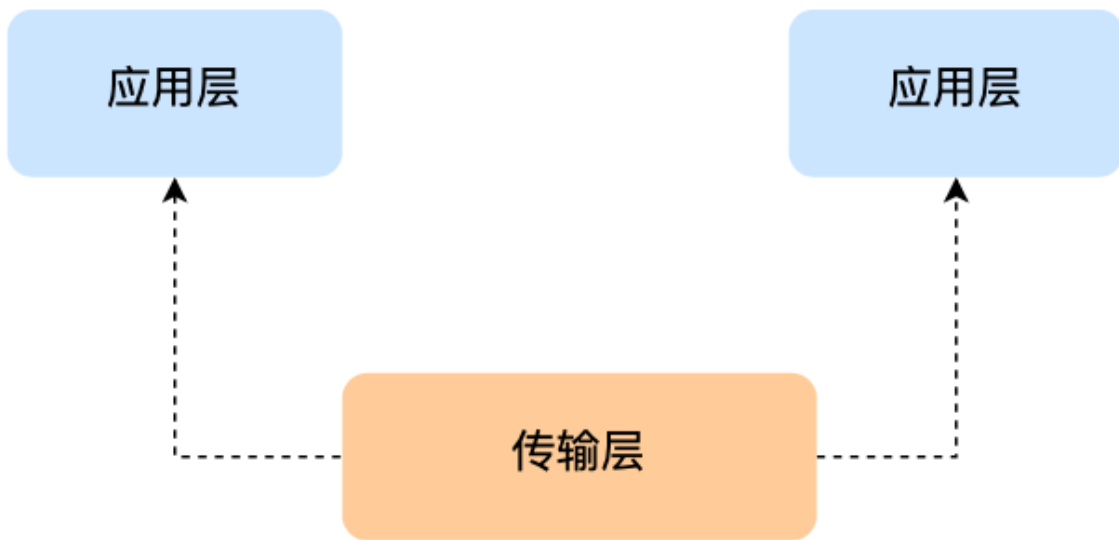
所以，应用层只需要专注于为用户提供应用功能，比如 HTTP、FTP、Telnet、DNS、SMTP 等。

应用层是不用去关心数据是如何传输的，就类似于，我们寄快递的时候，只需要把包裹交给快递员，由他负责运输快递，我们不需要关心快递是如何被运输的。

而且应用层是工作在操作系统中的用户态，传输层及以下则工作在内核态。

#### 传输层

应用层的数据包会传给传输层，**传输层**（*Transport Layer*）是为应用层提供网络支持的。

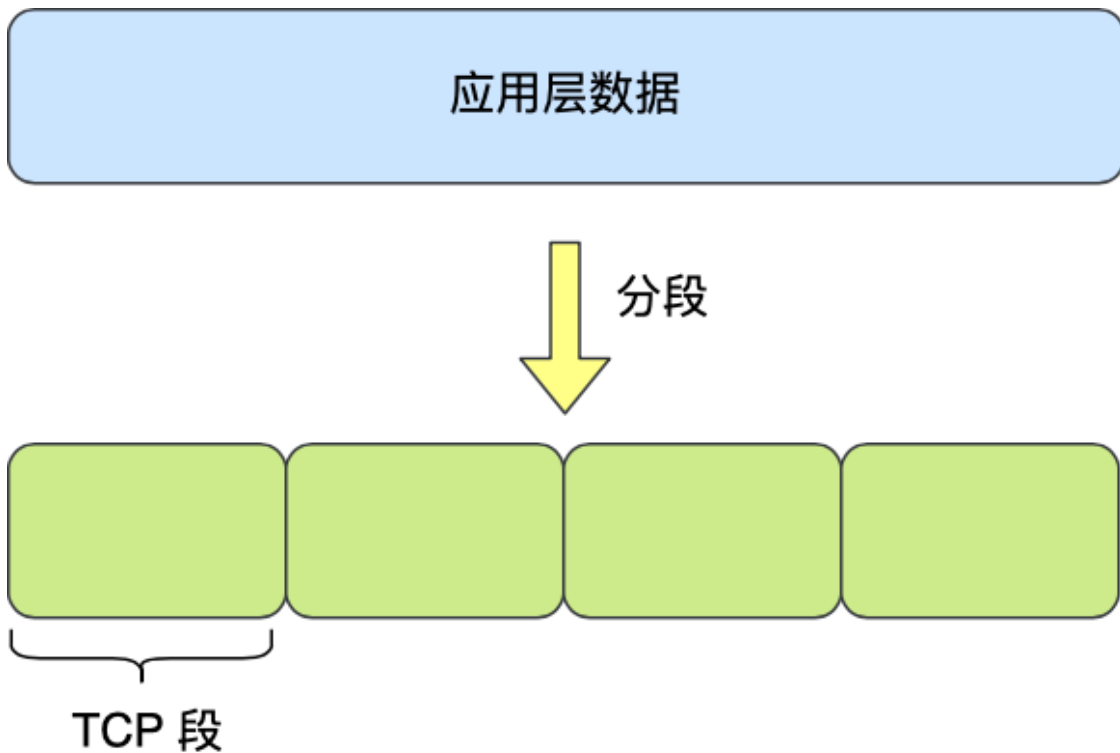


在传输层会有两个传输协议，分别是 TCP 和 UDP。

TCP 的全称叫传输控制协议（*Transmission Control Protocol*），大部分应用使用的正是 TCP 传输层协议，比如 HTTP 应用层协议。TCP 相比 UDP 多了很多特性，比如流量控制、超时重传、拥塞控制等，这些都是为了保证数据包能可靠地传输给对方。

UDP 相对来说就很简单，简单到只负责发送数据包，不保证数据包是否能抵达对方，但它实时性相对更好，传输效率也高。当然，UDP 也可以实现可靠传输，把 TCP 的特性在应用层上实现就可以，不过要实现一个商用的可靠 UDP 传输协议，也不是一件简单的事情。

应用需要传输的数据可能会非常大，如果直接传输就不好控制，因此当传输层的数据包大小超过 MSS（TCP 最大报文段长度），就要将数据包分块，这样即使中途有一个分块丢失或损坏了，只需要重新发送这一个分块，而不用重新发送整个数据包。在 TCP 协议中，我们把每个分块称为一个 **TCP 段**（*TCP Segment*）。



当设备作为接收方时，传输层则要负责把数据包传给应用，但是一台设备上可能会有很多应用在接收或者传输数据，因此需要用一个编号将应用区分开来，这个编号就是**端口**。

比如 80 端口通常是 Web 服务器用的，22 端口通常是远程登录服务器用的。而对于浏览器（客户端）中的每个标签栏都是一个独立的进程，操作系统会为这些进程分配临时的端口号。

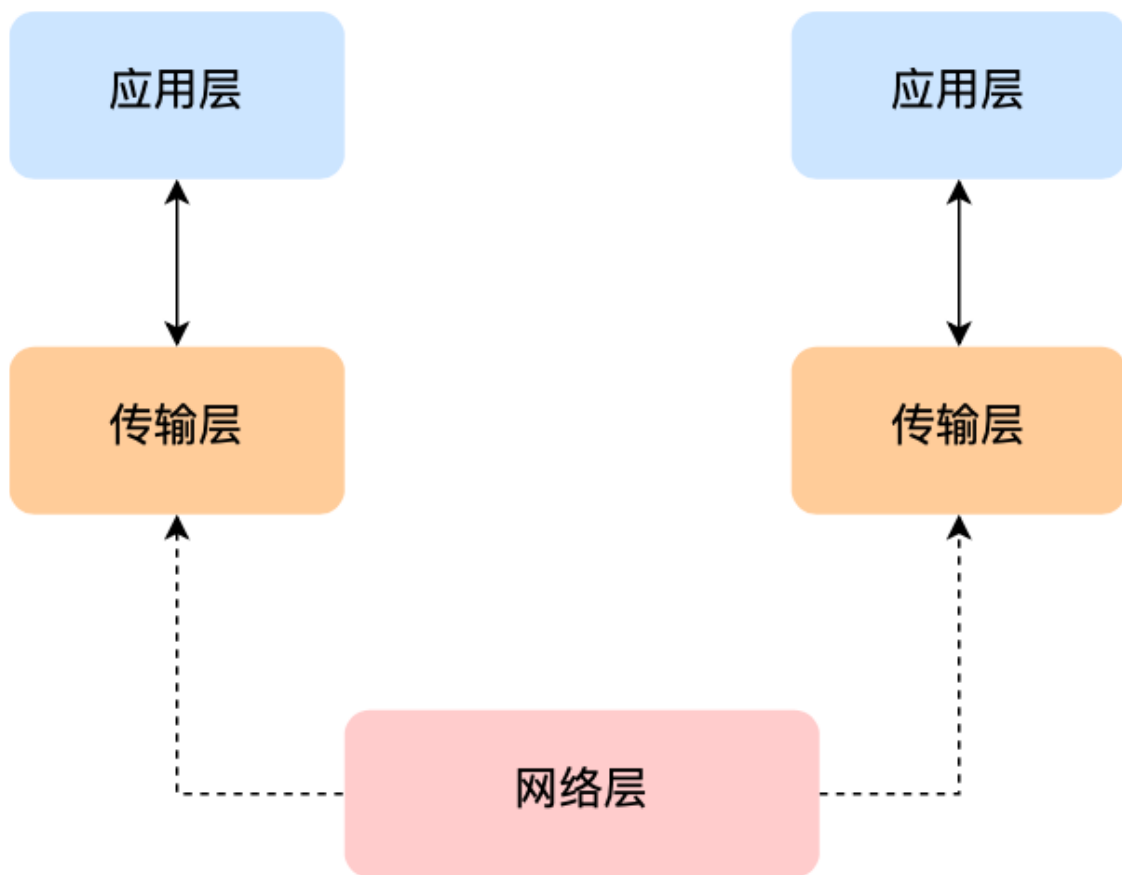
由于传输层的报文中会携带端口号，因此接收方可以识别出该报文是发送给哪个应用。

## 网络层

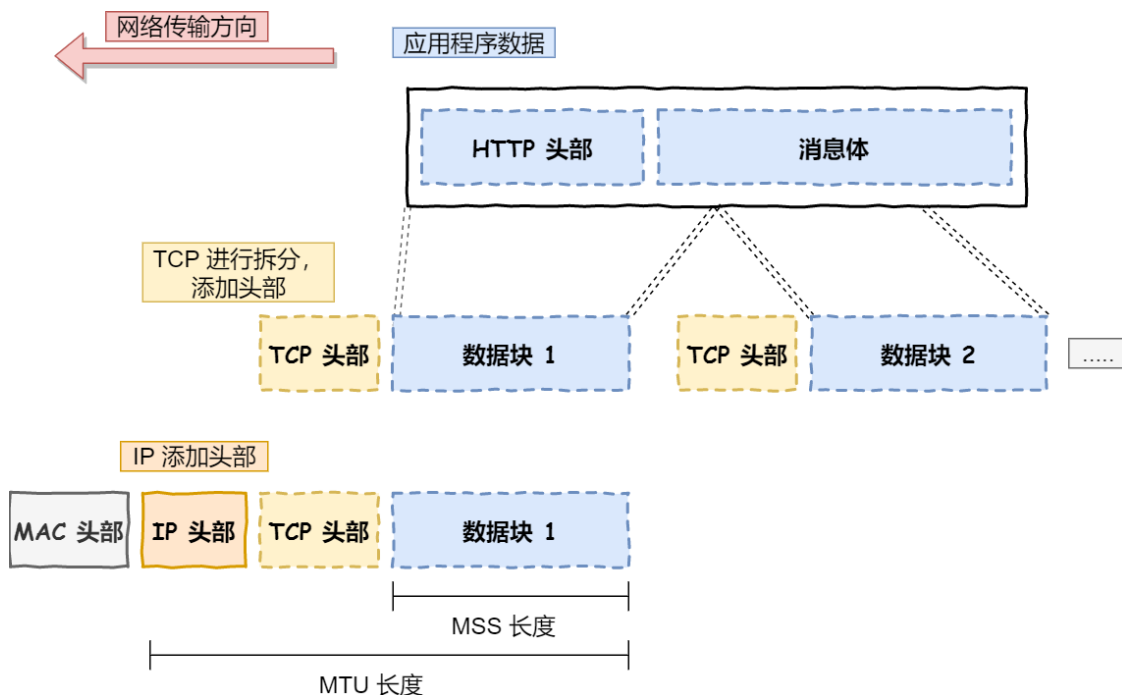
传输层可能大家刚接触的时候，会认为它负责将数据从一个设备传输到另一个设备，事实上它并不负责。

实际场景中的网络环节是错综复杂的，中间有各种各样的线路和分叉路口，如果一个设备的数据要传输给另一个设备，就需要在各种各样的路径和节点进行选择，而传输层的设计理念是简单、高效、专注，如果传输层还负责这一块功能就有点违背设计原则了。

也就是说，我们不希望传输层协议处理太多的事情，只需要服务好应用即可，让其作为应用间数据传输的媒介，帮助实现应用到应用的通信，而实际的传输功能就交给下一层，也就是**网络层**（*Internet Layer*）。



网络层最常使用的是 IP 协议（*Internet Protocol*），IP 协议会将传输层的报文作为数据部分，再加上 IP 报头组装成 IP 报文，如果 IP 报文大小超过 MTU（以太网中一般为 1500 字节）就会**再次进行分片**，得到一个即将发送到网络的 IP 报文。



网络层负责将数据从一个设备传输到另一个设备，世界上那么多设备，又该如何找到对方呢？因此，网络层需要有区分设备的编号。

我们一般用 IP 地址给设备进行编号，对于 IPv4 协议，IP 地址共 32 位，分成了四段（比如，192.168.100.1），每段是 8 位。只有一个单纯的 IP 地址虽然做到了区分设备，但是寻址起来就特别麻烦，全世界那么多台设备，难道一个一个去匹配？这显然不科学。

因此，需要将 IP 地址分成两种意义：

- 一个是**网络号**，负责标识该 IP 地址是属于哪个「子网」的；
- 一个是**主机号**，负责标识同一「子网」下的不同主机；

怎么分的呢？这需要配合**子网掩码**才能算出 IP 地址 的网络号和主机号。

举个例子，比如 10.100.122.0/24，后面的/24 表示就是 255.255.255.0 子网掩码，255.255.255.0 二进制是「11111111-11111111-11111111-00000000」，大家数数一共多少个 1？不用数了，是 24 个 1，为了简化子网掩码的表示，用/24 代替 255.255.255.0。

知道了子网掩码，该怎么计算出网络地址和主机地址呢？

将 10.100.122.2 和 255.255.255.0 进行**按位与运算**，就可以得到网络号，如下图：

IP 地址 : 10.100.122.2

00001010 1100100 1111010 00000010

子网掩码 : 255.255.255.0

11111111 11111111 11111111 00000000



IP 地址 和 子网掩码  
做 AND 运算

网络号 : 10.100.122.0

00001010 1100100 1111010

00000000

网络号

将 255.255.255.0 取反后与 IP 地址进行按位与运算，就可以得到主机号。

大家可以去搜索下子网掩码计算器，自己改变下「掩码位」的数值，就能体会到子网掩码的作用了。

网络IP地址计算器

输入IP

192

168

100

1

掩码位

24

计算

清空

可用IP

254

掩码

255

255

255

0

网络

192

168

100

0

第一可用

192

168

100

1

最后可用

192

168

100

254

广播

192

168

100

255

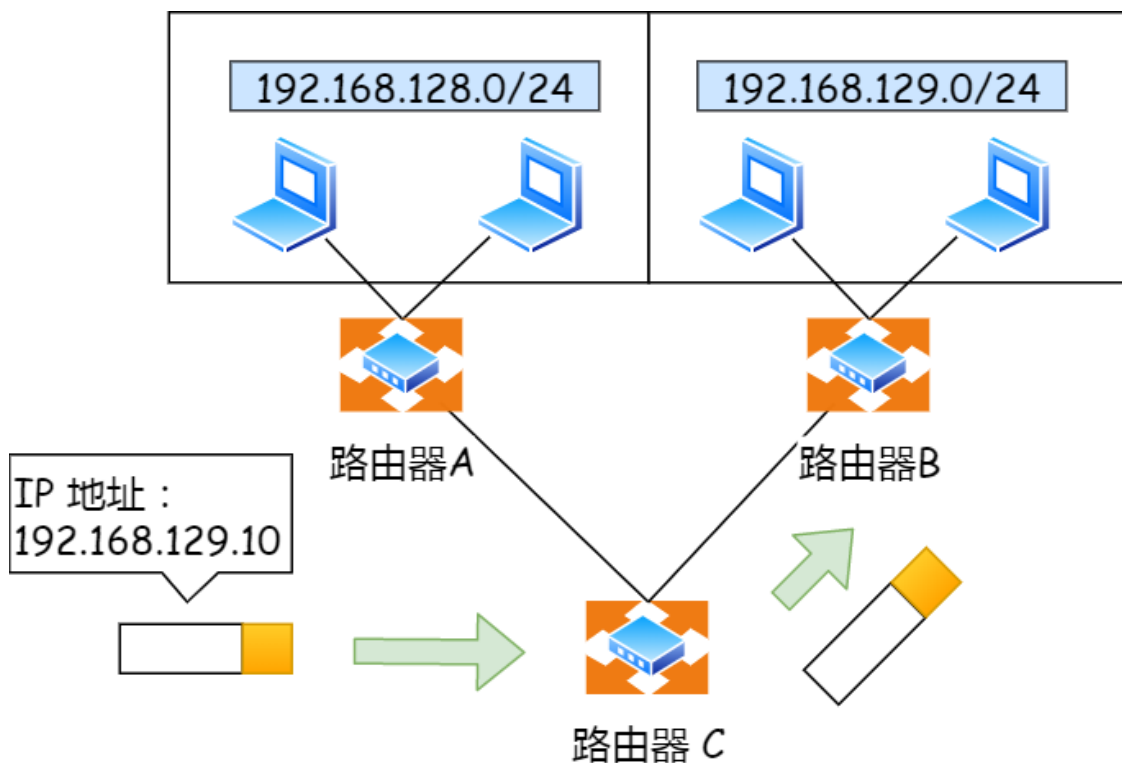
显示网络，广播，第一次和最后一个给定的网络地址

在网络掩码“位格式”也被称为CIDR格式（CIDR=无类别域间路由选择）。

那么在寻址的过程中，先匹配到相同的网络号（表示要找到同一个子网），才会去找对应的主机。

除了寻址能力，IP 协议还有另一个重要的能力就是**路由**。实际场景中，两台设备并不是用一条网线连接起来的，而是通过很多网关、路由器、交换机等众多网络设备连接起来的，那么就会形成很多条网络的路径，因此当数据包到达一个网络节点，就需要通过路由算法决定下一步走哪条路径。

路由器寻址工作中，就是要找到目标地址的子网，找到后进而把数据包转发给对应的网络内。



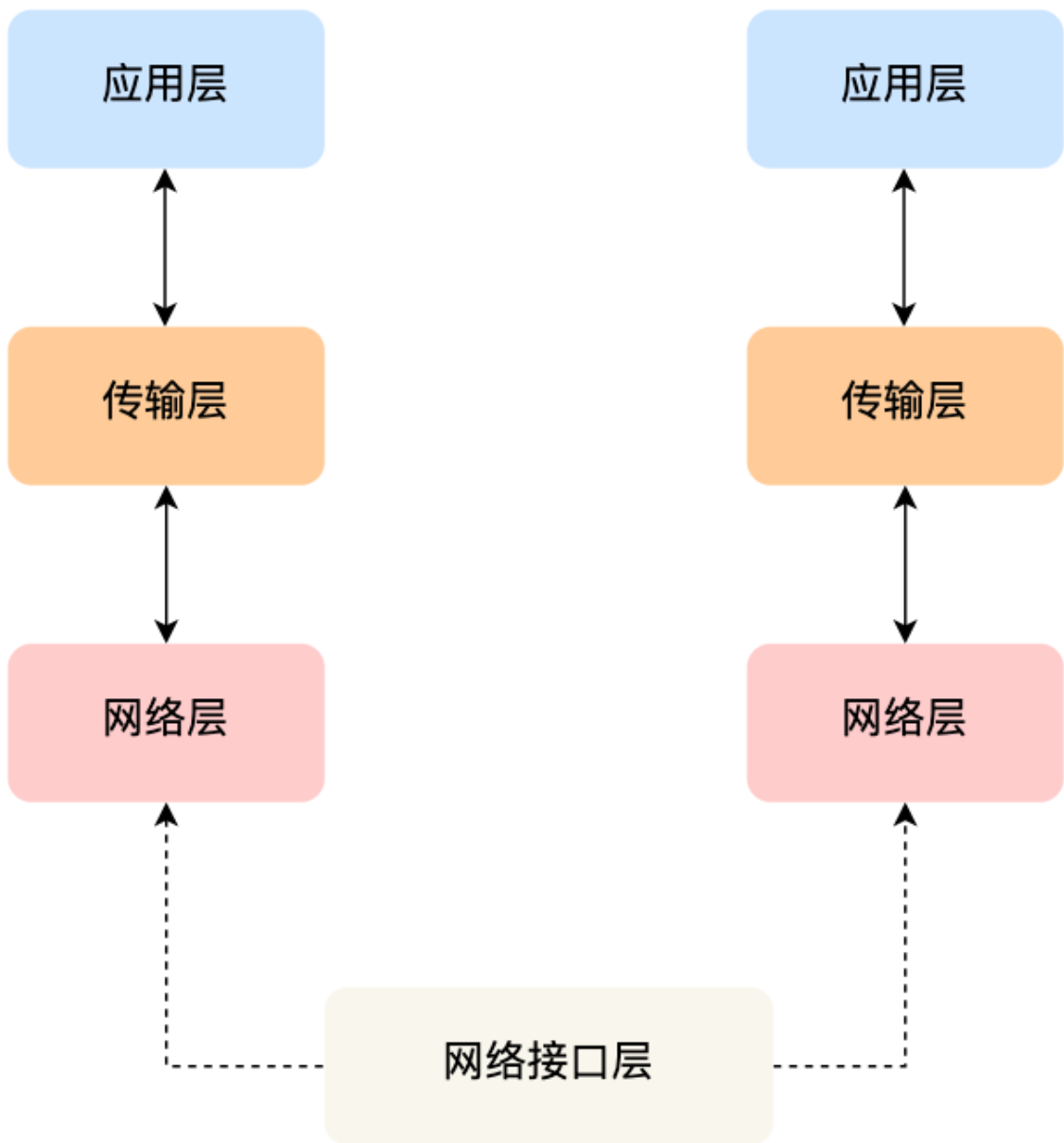
路由器只要一看到 IP 地址的  
网络号就可以进行转发

所以，IP 协议的寻址作用是告诉我们去往下一个目的地该朝哪个方向走，路由则是根据「下一个目的地」选择路径。寻址更像在导航，路由更像在操作方向盘。

## 网络接口层

生成了 IP 头部之后，接下来要交给网络接口层（*Link Layer*）在 IP 头部的前面加上 MAC 头部，并封装成数据帧（Data frame）发送到网络上。





IP 头部中的接收方 IP 地址表示网络包的目的地，通过这个地址我们就可以判断要将包发到哪里，但在以太网的世界中，这个思路是行不通的。

什么是以太网呢？电脑上的以太网接口，Wi-Fi 接口，以太网交换机、路由器上的千兆，万兆以太网口，还有网线，它们都是以太网的组成部分。以太网就是一种在「局域网」内，把附近的设备连接起来，使它们之间可以进行通讯的技术。

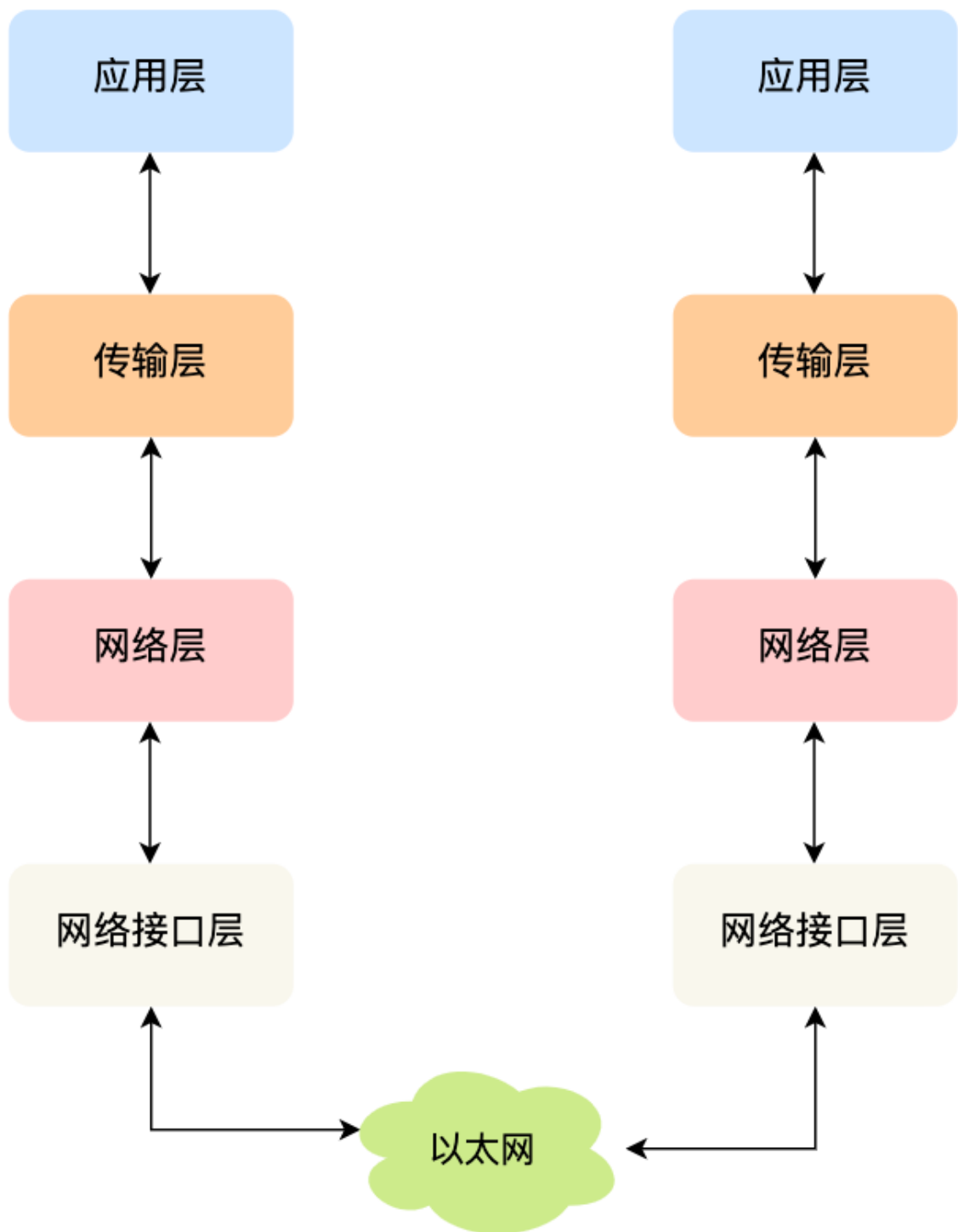
以太网在判断网络包目的地时和 IP 的方式不同，因此必须采用相匹配的方式才能在以太网中将包发往目的地，而 MAC 头部就是干这个用的，所以，在以太网进行通讯要用到 MAC 地址。

MAC 头部是以太网使用的头部，它包含了接收方和发送方的 MAC 地址等信息，我们可以通过 ARP 协议获取对方的 MAC 地址。

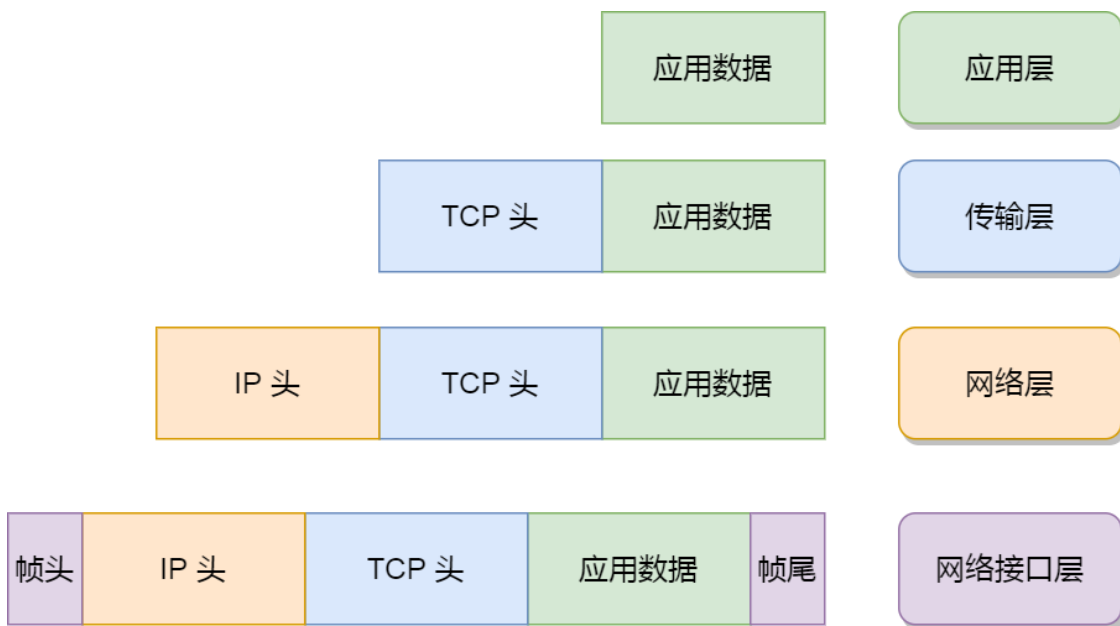
所以说，网络接口层主要为网络层提供「链路级别」传输的服务，负责在以太网、WiFi 这样的底层网络上发送原始数据包，工作在网卡这个层次，使用 MAC 地址来标识网络上的设备。

## 总结

综上所述，TCP/IP 网络通常是由上到下分成 4 层，分别是应用层，传输层，网络层和网络接口层。



再给大家贴一下每一层的封装格式：



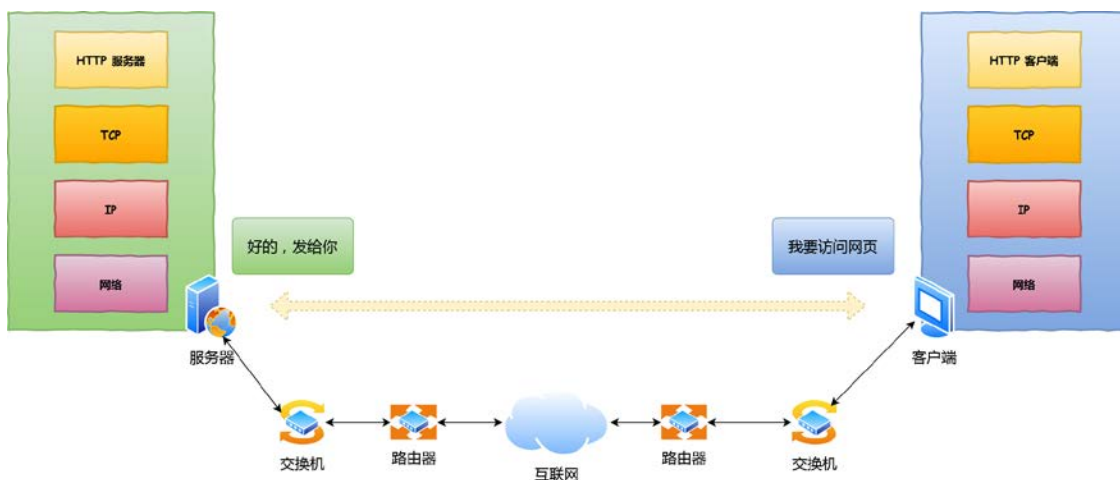
网络接口层的传输单位是帧（frame），IP 层的传输单位是包（packet），TCP 层的传输单位是段（segment），HTTP 的传输单位则是消息或报文（message）。但这些名词并没有什么本质的区分，可以统称为数据包。

## 2.2 键入网址到网页显示，期间发生了什么？

想必不少小伙伴面试过程中，会遇到「当键入网址后，到网页显示，其间发生了什么」的面试题。

还别说，这问题真挺常问的，前几天坐在我旁边的主管电话面试应聘者的时候，也问了这个问题。

接下来以下图较简单的网络拓扑模型作为例子，探究探究其间发生了什么？

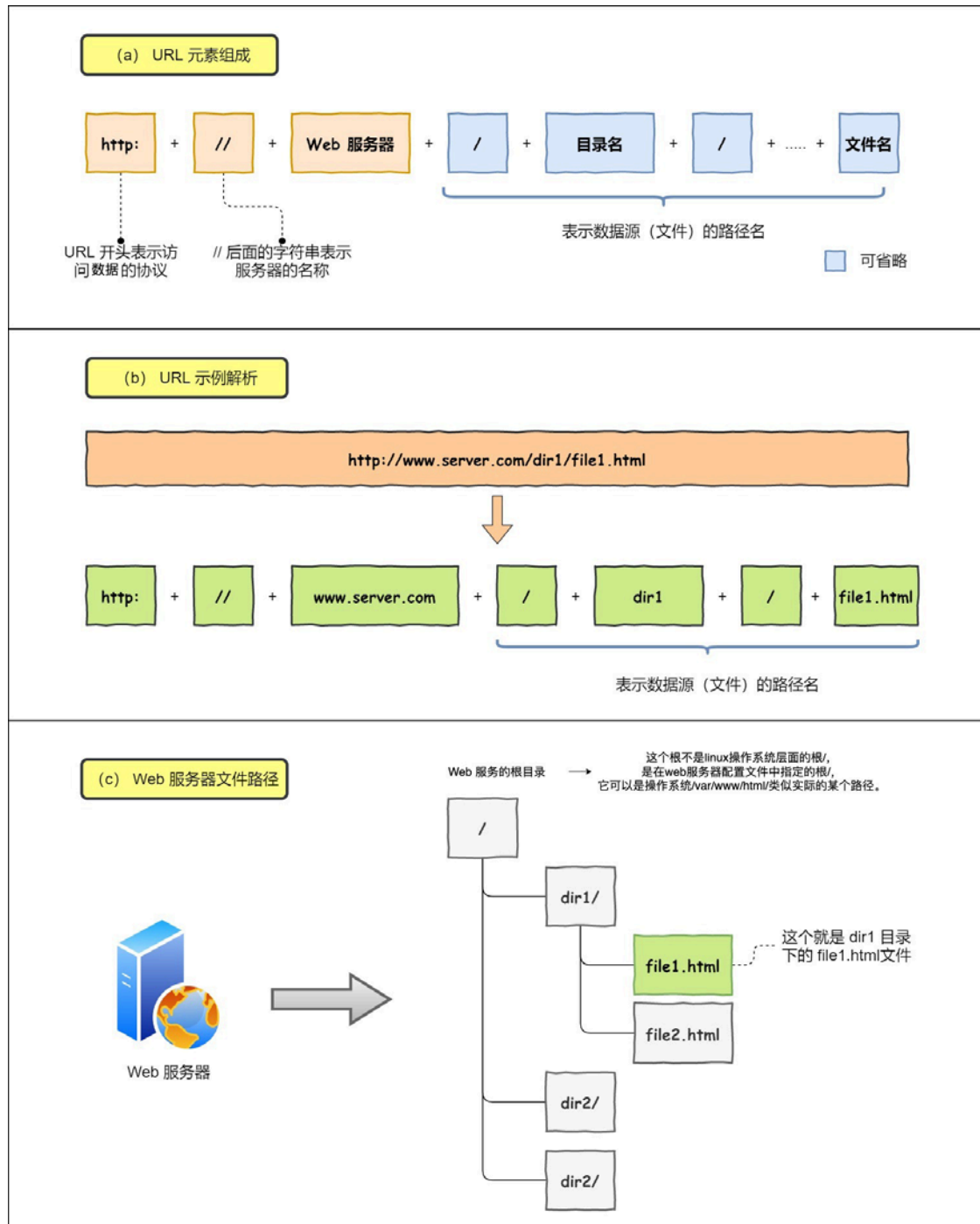


## 孤单小弟 —— HTTP

浏览器做的第一步工作是解析 URL

首先浏览器做的第一步工作就是要对 URL 进行解析，从而生成发送给 Web 服务器的请求信息。

让我们看看一条长长的 URL 里的各个元素的代表什么，见下图：



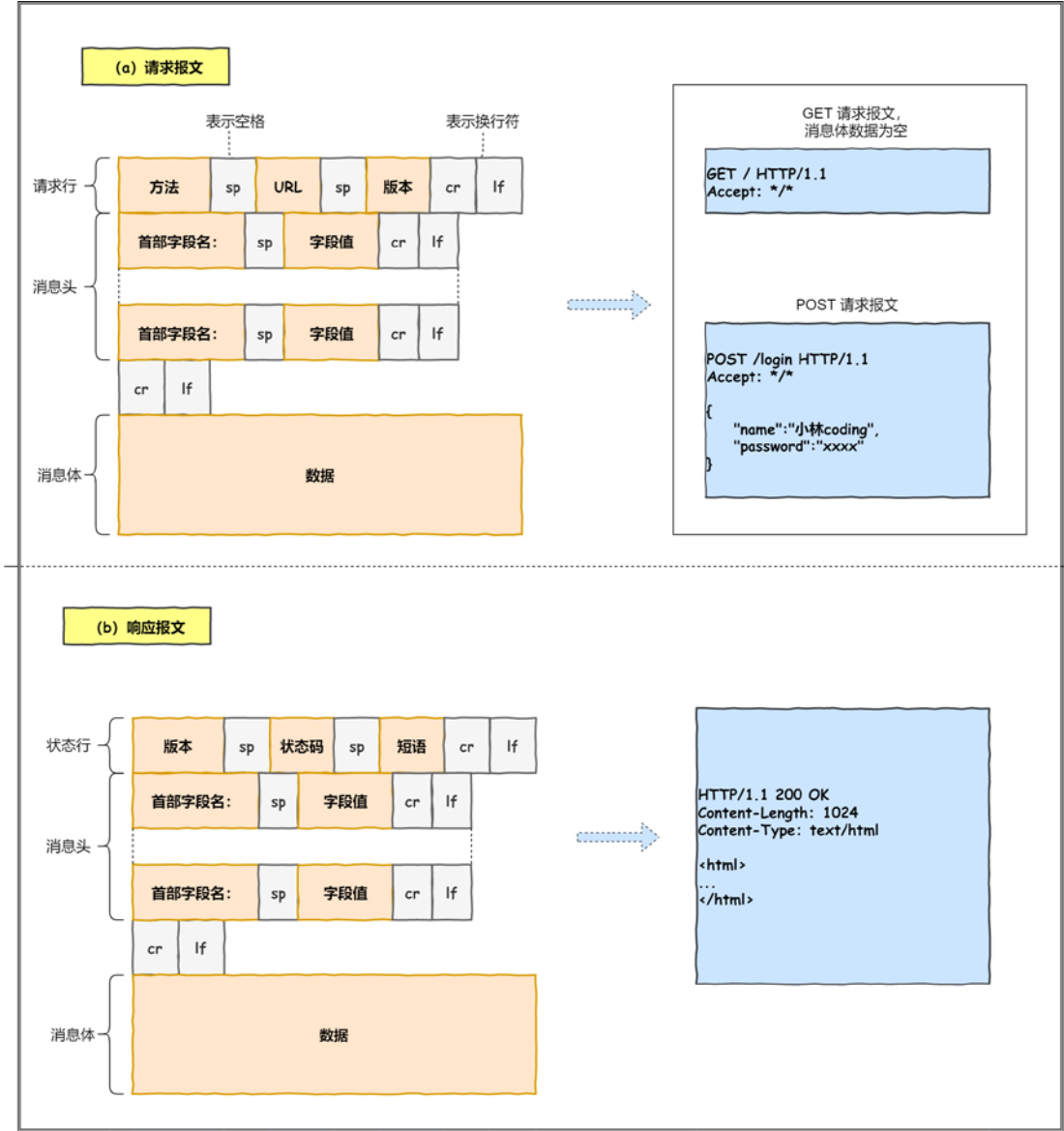
所以图中的长长的 URL 实际上是请求服务器里的文件资源。

要是上图中的蓝色部分 URL 元素都省略了，那应该是请求哪个文件呢？

当没有路径名时，就代表访问根目录下事先设置的**默认文件**，也就是 `/index.html` 或者 `/default.html` 这些文件，这样就不会发生混乱了。

生产 HTTP 请求信息

对 URL 进行解析之后，浏览器确定了 Web 服务器和文件名，接下来就是根据这些信息来生成 HTTP 请求消息了。



一个孤单 HTTP 数据包表示：“我这么一个小小的数据包，没亲没友，直接发到浩瀚的网络，谁会知道我呢？谁能载我一程呢？谁能保护我呢？我的目的地在哪呢？”充满各种疑问的它，没有停滞不前，依然踏上了征途！

---

## 真实地址查询 —— DNS

通过浏览器解析 URL 并生成 HTTP 消息后，需要委托操作系统将消息发送给 Web 服务器。

但在发送之前，还有一项工作需要完成，那就是**查询服务器域名对应的 IP 地址**，因为委托操作系统发送消息时，必须提供通信对象的 IP 地址。

比如我们打电话的时候，必须要知道对方的电话号码，但由于电话号码难以记忆，所以通常我们会将对方电话号 + 姓名保存在通讯录里。

所以，有一种服务器就专门保存了 Web 服务器域名与 IP 的对应关系，它就是 DNS 服务器。

### 域名的层级关系

DNS 中的域名都是用**句点**来分隔的，比如 `www.server.com`，这里的句点代表了不同层次之间的**界限**。

在域名中，**越靠右**的位置表示其层级**越高**。

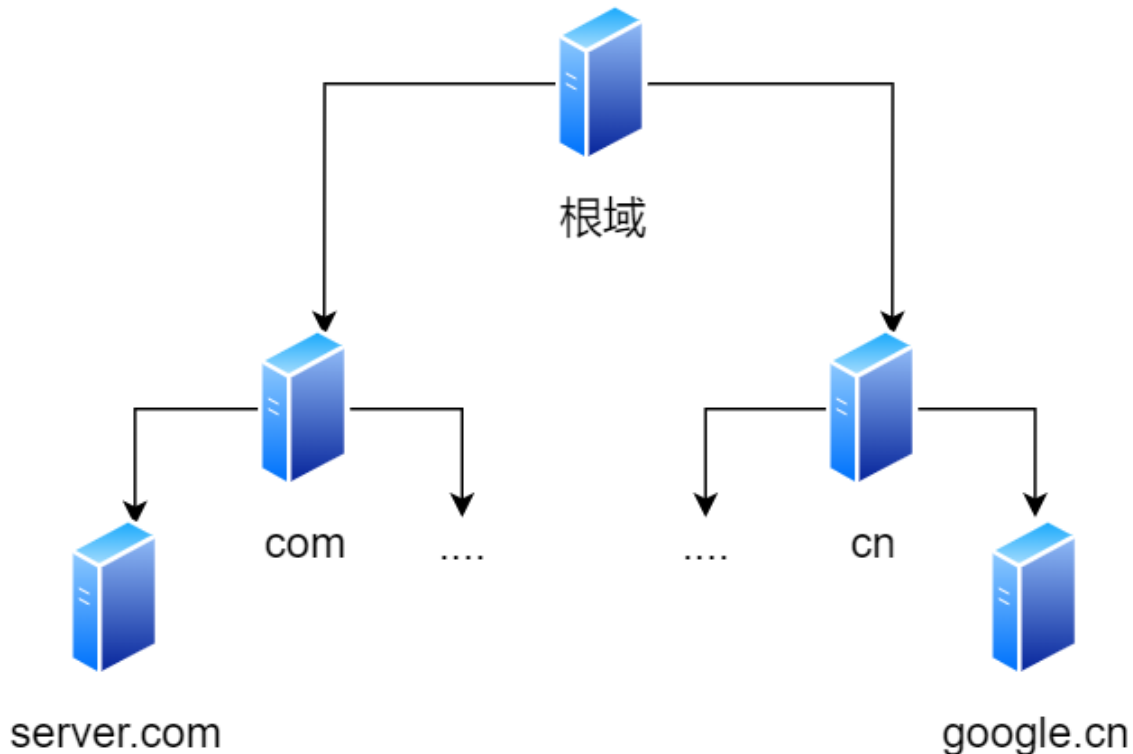
毕竟域名是外国人发明，所以思维和中国人相反，比如说一个城市地点的时候，外国喜欢从小到大的方式顺序说起（如 XX 街道 XX 区 XX 市 XX 省），而中国则喜欢从大到小的顺序（如 XX 省 XX 市 XX 区 XX 街道）。

实际上域名最后还有一个点，比如 `www.server.com.`，这个最后的一个点代表根域名。

也就是，`.` 根域是在最顶层，它的下一层就是 `.com` 顶级域，再下面是 `server.com`。

所以域名的层级关系类似一个树状结构：

- 根 DNS 服务器（`.`）
- 顶级域 DNS 服务器（`.com`）
- 权威 DNS 服务器（`server.com`）



根域的 DNS 服务器信息保存在互联网中所有的 DNS 服务器中。

这样一来，任何 DNS 服务器就都可以找到并访问根域 DNS 服务器了。

因此，客户端只要能够找到任意一台 DNS 服务器，就可以通过它找到根域 DNS 服务器，然后再一路顺藤摸瓜找到位于下层的某台目标 DNS 服务器。

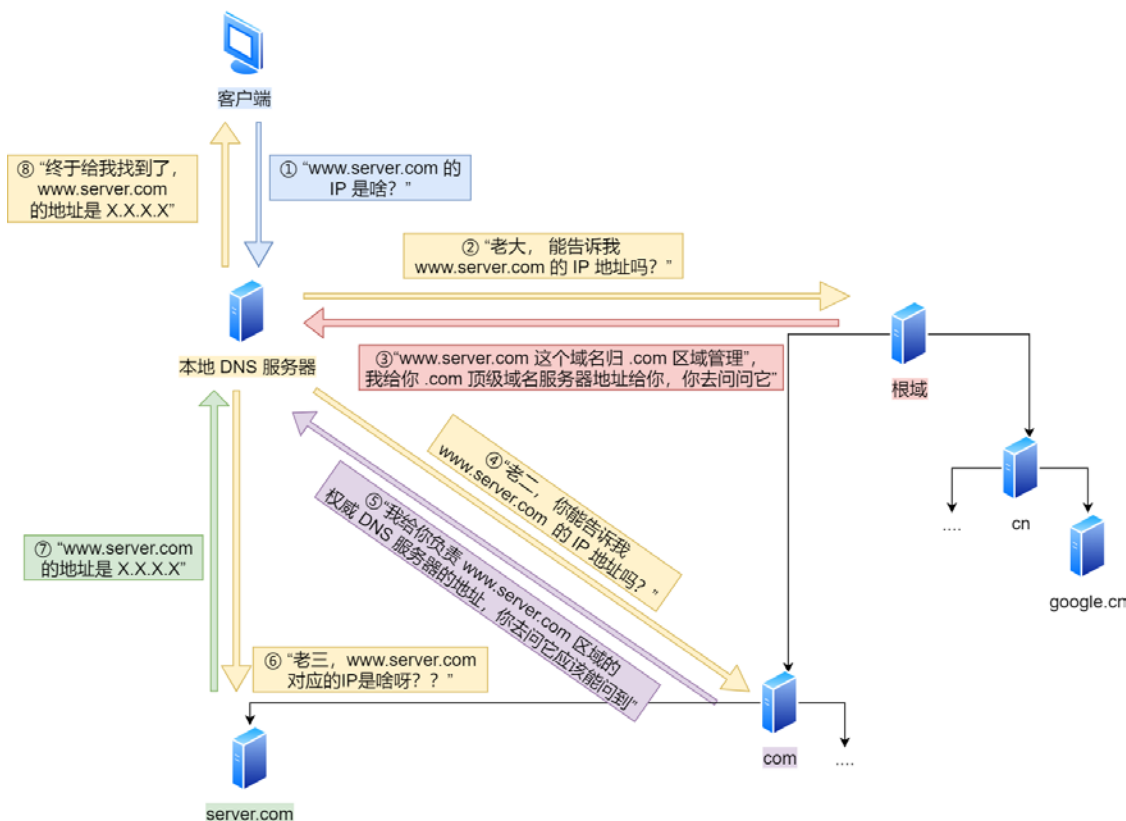
#### 域名解析的工作流程

1. 客户端首先会发出一个 DNS 请求，问 [www.server.com](http://www.server.com) 的 IP 是啥，并发给本地 DNS 服务器（也就是客户端的 TCP/IP 设置中填写的 DNS 服务器地址）。
2. 本地域名服务器收到客户端的请求后，如果缓存里的表格能找到 [www.server.com](http://www.server.com)，则它直接返回 IP 地址。如果没有，本地 DNS 会去问它的根域名服务器：“老大，能告诉我 [www.server.com](http://www.server.com) 的 IP 地址吗？”根域名服务器是最高层次的，它不直接用于域名解析，但能指明一条道路。
3. 根 DNS 收到来自本地 DNS 的请求后，发现后置是 .com，说：“[www.server.com](http://www.server.com) 这个域名归 .com 区域管理”，我把 .com 顶级域名服务器的地址给你，你去问问它吧。”
4. 本地 DNS 收到顶级域名服务器的地址后，发起请求问“老二，你能告诉我 [www.server.com](http://www.server.com) 的 IP 地址吗？”
5. 顶级域名服务器说：“我给你负责 [www.server.com](http://www.server.com) 区域的权威 DNS 服务器的地址，你去问它应该能问到”。



- 本地 DNS 于是转向问权威 DNS 服务器：“老三，[www.server.com](http://www.server.com) 对应的 IP 是啥呀？”[server.com](http://server.com) 的权威 DNS 服务器，它是域名解析结果的原出处。为啥叫权威呢？就是我的域名我做主。
- 权威 DNS 服务器查询后将对应的 IP 地址 X.X.X.X 告诉本地 DNS。
- 本地 DNS 再将 IP 地址返回客户端，客户端和目标建立连接。

至此，我们完成了 DNS 的解析过程。现在总结一下，整个过程我画成了一个图。



DNS 域名解析的过程蛮有意思的，整个过程就和我们日常生活中找人问路的过程类似，只指路不带路。

那是不是每次解析域名都要经过那么多的步骤呢？

当然不是了，还有缓存这个东西的嘛。

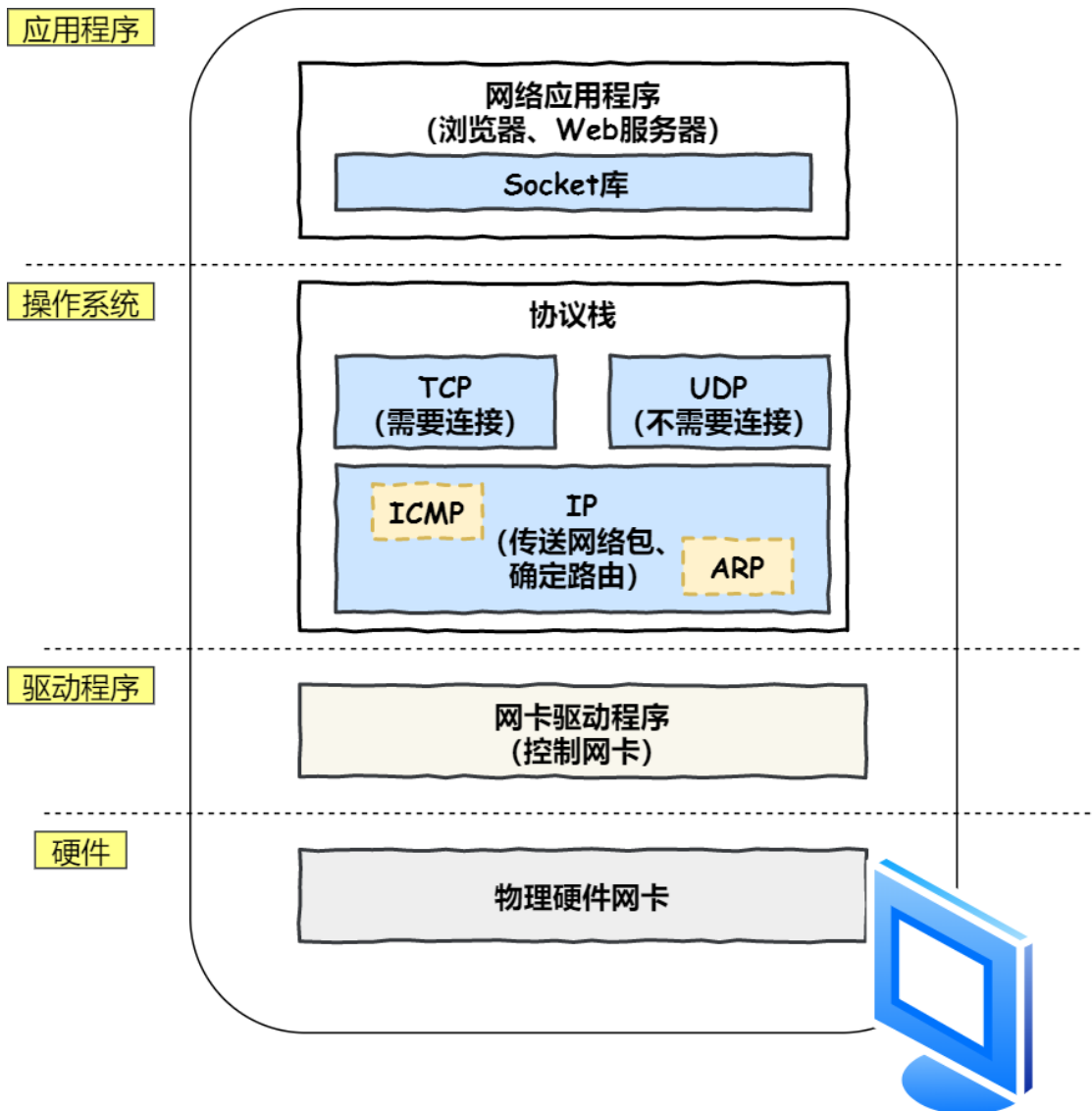
浏览器会先看自身有没有对这个域名的缓存，如果有，就直接返回，如果没有，就去问操作系统，操作系统也会去看自己的缓存，如果有，就直接返回，如果没有，再去 `hosts` 文件看，也没有，才会去问「本地 DNS 服务器」。

数据包表示：“DNS 老大哥厉害呀，找到了目的地了！我还是很迷茫呀，我要发出去，接下来我需要谁的帮助呢？”

## 指南好帮手 —— 协议栈

通过 DNS 获取到 IP 后，就可以把 HTTP 的传输工作交给操作系统中的**协议栈**。

协议栈的内部分为几个部分，分别承担不同的工作。上下关系是有一定的规则的，上面的部分会向下面的部分委托工作，下面的部分收到委托的工作并执行。



应用程序（浏览器）通过调用 Socket 库，来委托协议栈工作。协议栈的上半部分有两块，分别是负责收发数据的 TCP 和 UDP 协议，这两个传输协议会接受应用层的委托执行收发数据的操作。

协议栈的下面一半是用 IP 协议控制网络包收发操作，在互联网上传数据时，数据会被切分成一块块的网络包，而将网络包发送给对方的操作就是由 IP 负责的。

此外 IP 中还包括 ICMP 协议和 ARP 协议。

- ICMP 用于告知网络包传送过程中产生的错误以及各种控制信息。
- ARP 用于根据 IP 地址查询相应的以太网 MAC 地址。

IP 下面的网卡驱动程序负责控制网卡硬件，而最下面的网卡则负责完成实际的收发操作，也就是对网线中的信号执行发送和接收操作。

数据包看了这份指南表示：“原来我需要那么多大佬的协助啊，那我先去找找 TCP 大佬！”

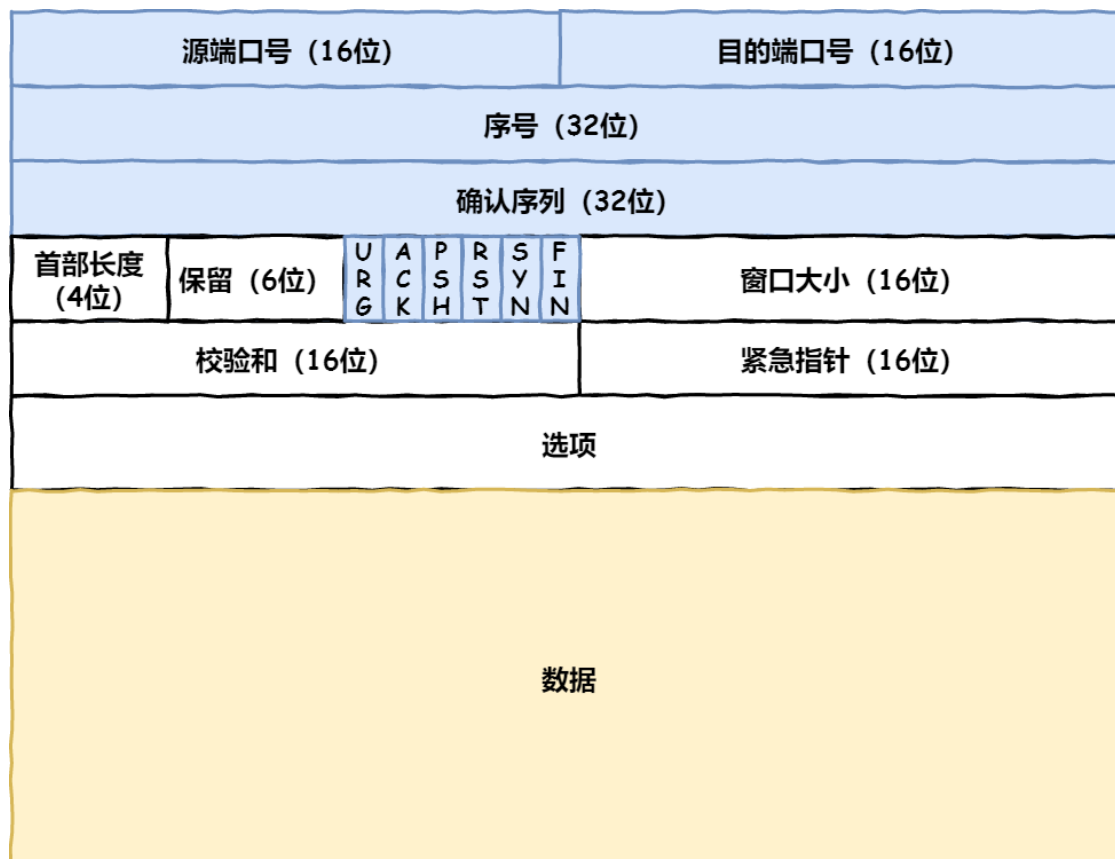
---

## 可靠传输 —— TCP

HTTP 是基于 TCP 协议传输的，所以在这我们先了解下 TCP 协议。

TCP 包头格式

我们先看看 TCP 报文头部的格式：



首先，**源端口号**和**目标端口号**是不可少的，如果没有这两个端口号，数据就不知道应该发给哪个应用。

接下来有包的**序号**，这个是为了解决包乱序的问题。

还有应该有的是**确认号**，目的是确认发出去对方是否有收到。如果没有收到就应该重新发送，直到送达，这个是为了解决不丢包的问题。

接下来还有一些**状态位**。例如 **SYN** 是发起一个连接，**ACK** 是回复，**RST** 是重新连接，**FIN** 是结束连接等。**TCP** 是面向连接的，因而双方要维护连接的状态，这些带状态位的包的发送，会引起双方的状态变更。

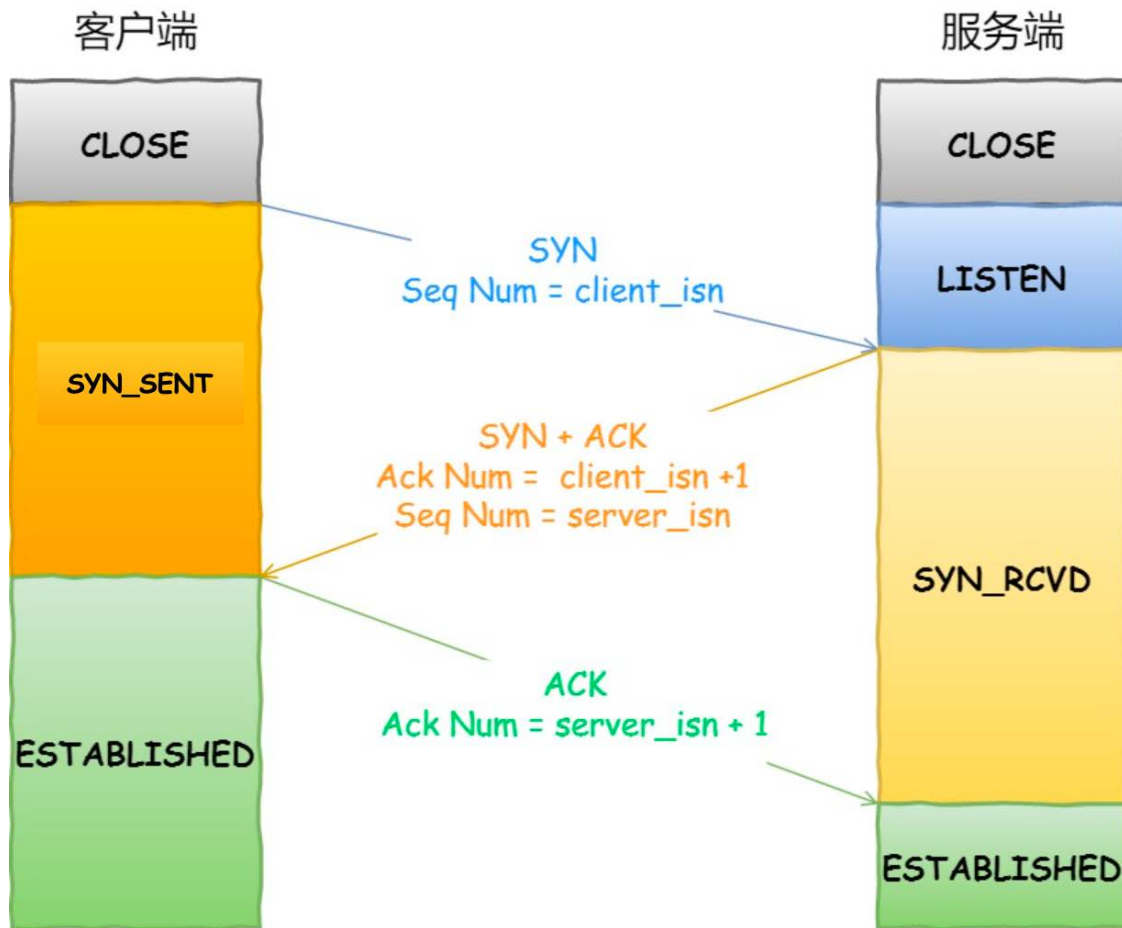
还有一个重要的就是**窗口大小**。**TCP** 要做**流量控制**，通信双方各声明一个窗口（缓存大小），标识自己当前能够的处理能力，别发送的太快，撑死我，也别发的太慢，饿死我。

除了做流量控制以外，**TCP** 还会做**拥塞控制**，对于真正的通路堵车不堵车，它无能为力，唯一能做的就是控制自己，也即控制发送的速度。不能改变世界，就改变自己嘛。

**TCP** 传输数据之前，要先三次握手建立连接

在 **HTTP** 传输数据之前，首先需要 **TCP** 建立连接，**TCP** 连接的建立，通常称为**三次握手**。

这个所谓的「连接」，只是双方计算机里维护一个状态机，在连接建立的过程中，双方的状态变化时序图就像这样。



- 一开始，客户端和服务端都处于 **CLOSED** 状态。先是服务端主动监听某个端口，处于 **LISTEN** 状态。
- 然后客户端主动发起连接 **SYN**，之后处于 **SYN-SENT** 状态。
- 服务端收到发起的连接，返回 **SYN**，并且 **ACK** 客户端的 **SYN**，之后处于 **SYN-RCVD** 状态。
- 客户端收到服务端发送的 **SYN** 和 **ACK** 之后，发送对 **SYN** 确认的 **ACK**，之后处于 **ESTABLISHED** 状态，因为它一发一收成功了。
- 服务端收到 **ACK** 的 **ACK** 之后，处于 **ESTABLISHED** 状态，因为它也一发一收了。

所以三次握手目的是保证双方都有发送和接收的能力。

如何查看 TCP 的连接状态？

TCP 的连接状态查看，在 Linux 可以通过 `netstat -napt` 命令查看。

```
[root@lincoding ~]# netstat -npt
```

Active Internet connections (servers and established)

Proto	Recv-Q	Send-Q	Local Address	Foreign Address	State	PID/Program name
tcp	0	0	::ffff:192.168.3.100:80	::ffff:192.168.3.20:55288	ESTABLISHED	3391/httpd

TCP 协议

源地址 + 端口

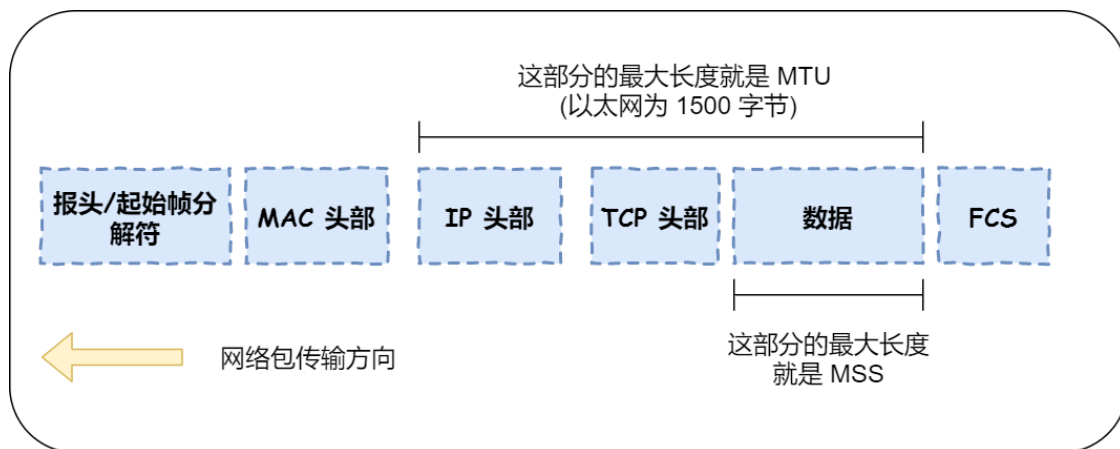
目标地址 + 端口

连接状态

Web 服务的  
进程 PID 和 进程名称

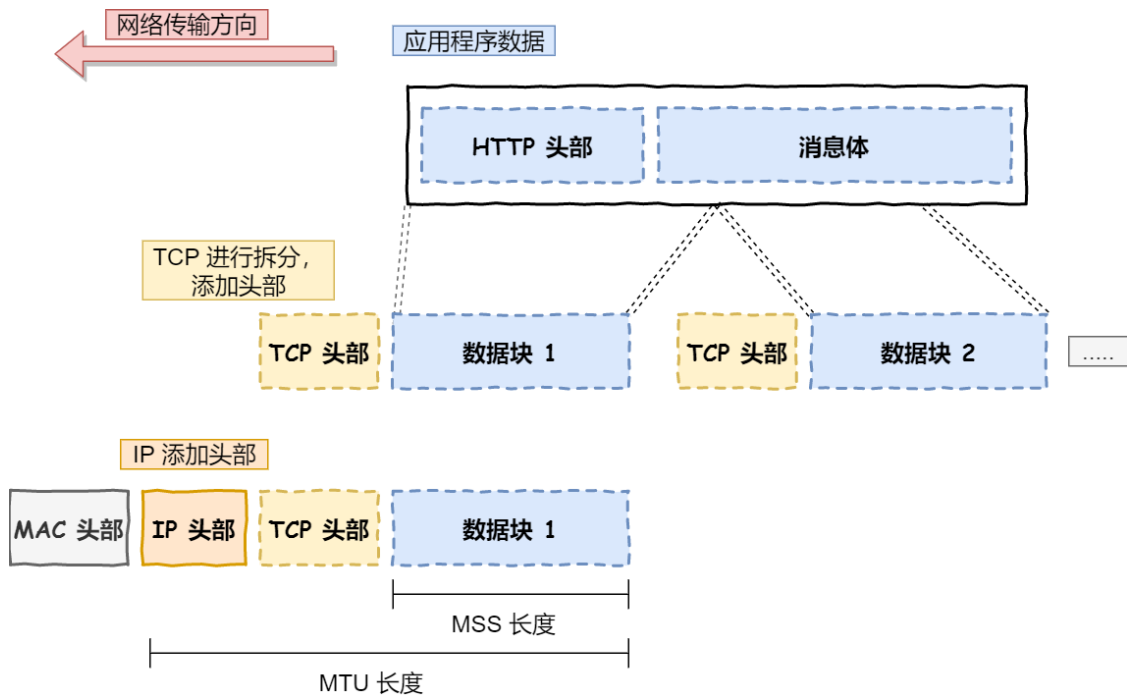
## TCP 分割数据

如果 HTTP 请求消息比较长，超过了 MSS 的长度，这时 TCP 就需要把 HTTP 的数据拆解成一块块的数据发送，而不是一次性发送所有数据。



- **MTU**: 一个网络包的最大长度，以太网中一般为 **1500** 字节。
- **MSS**: 除去 IP 和 TCP 头部之后，一个网络包所能容纳的 TCP 数据的最大长度。

数据会被以 MSS 的长度为单位进行拆分，拆分出来的每一块数据都会被放进单独的网络包中。也就是在每个被拆分的数据加上 TCP 头信息，然后交给 IP 模块来发送数据。

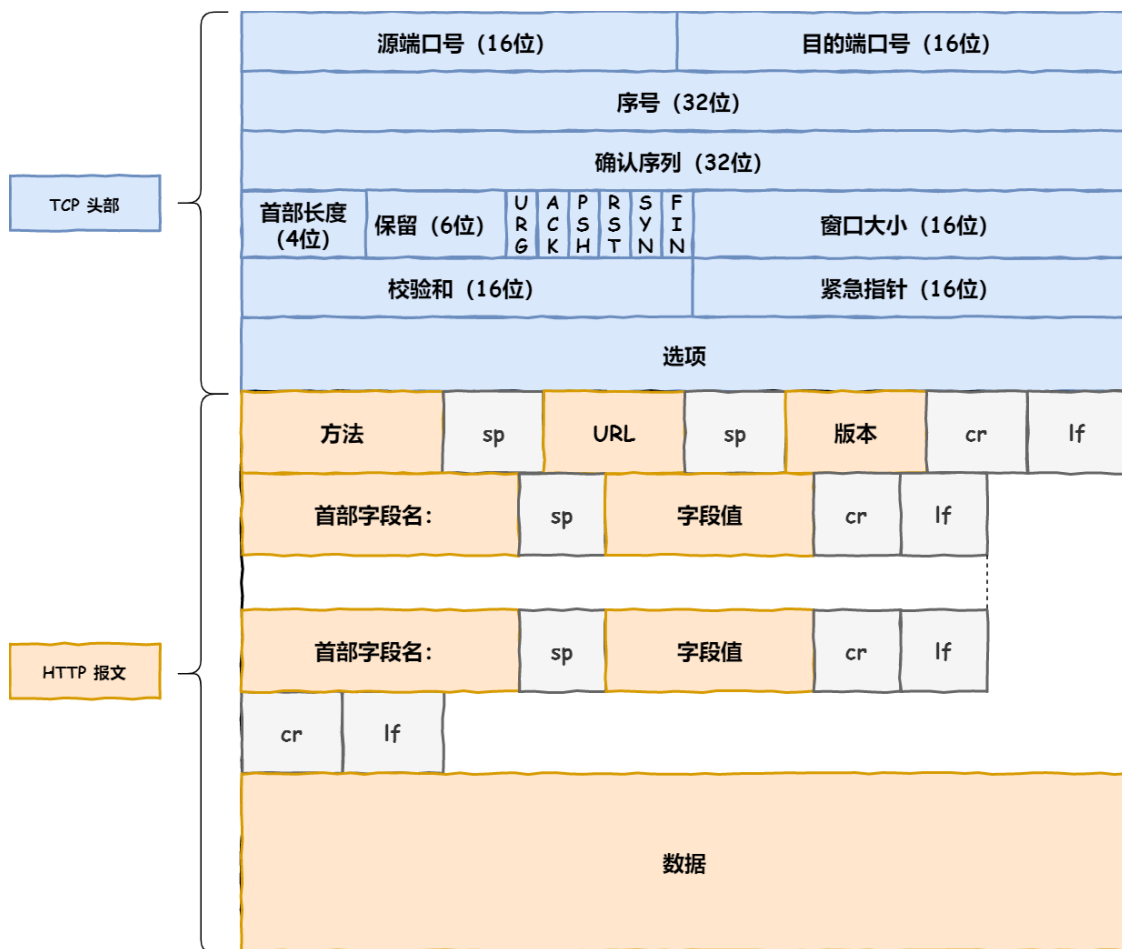


### TCP 报文生成

TCP 协议里面会有两个端口，一个是浏览器监听的端口（通常是随机生成的），一个是 Web 服务器监听的端口（HTTP 默认端口号是 80，HTTPS 默认端口号是 443）。

在双方建立了连接后，TCP 报文中的数据部分就是存放 HTTP 头部 + 数据，组装好 TCP 报文之后，就需交给下面的网络层处理。

至此，网络包的报文如下图。



此时，遇上了 TCP 的数据包激动表示：“太好了，碰到了可靠传输的 TCP 传输，它给我加上 TCP 头部，我不再孤单了，安全感十足啊！有大佬可以保护我的可靠送达！但我应该往哪走呢？”

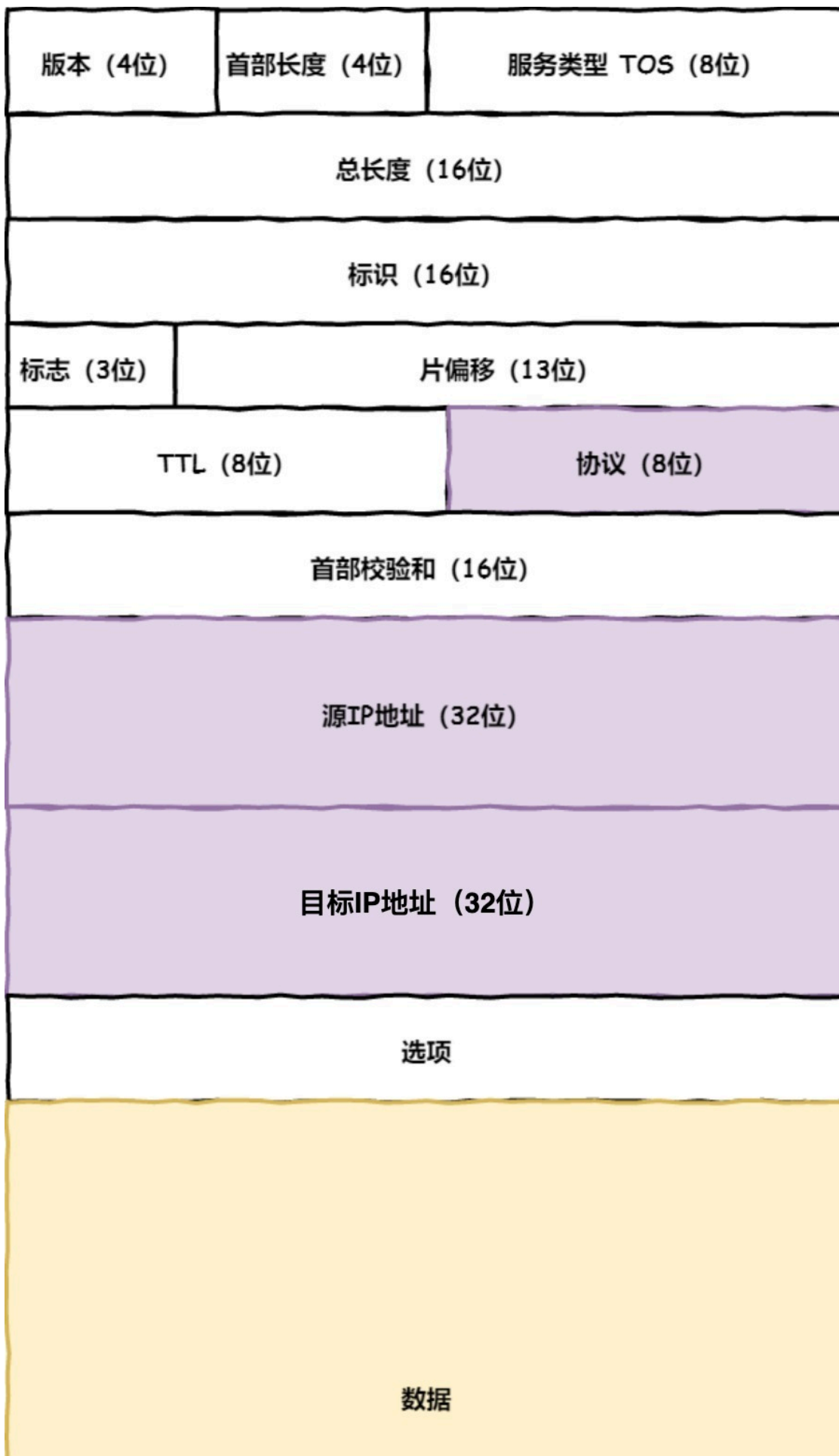
## 远程定位 —— IP

TCP 模块在执行连接、收发、断开等各阶段操作时，都需要委托 IP 模块将数据封装成**网络包**发送给通信对象。

### IP 包头格式

我们先看看 IP 报文头部的格式：





在 IP 协议里面需要有**源地址 IP** 和 **目标地址 IP**：

- 源地址 IP，即是客户端输出的 IP 地址；
- 目标地址，即通过 DNS 域名解析得到的 Web 服务器 IP。

因为 HTTP 是经过 TCP 传输的，所以在 IP 包头的**协议号**，要填写为 **06**（十六进制），表示协议为 TCP。

假设客户端有多个网卡，就会有多个 IP 地址，那 IP 头部的源地址应该选择哪个 IP 呢？

当存在多个网卡时，在填写源地址 IP 时，就需要判断到底应该填写哪个地址。这个判断相当于在多块网卡中判断应该使用哪一块网卡来发送包。

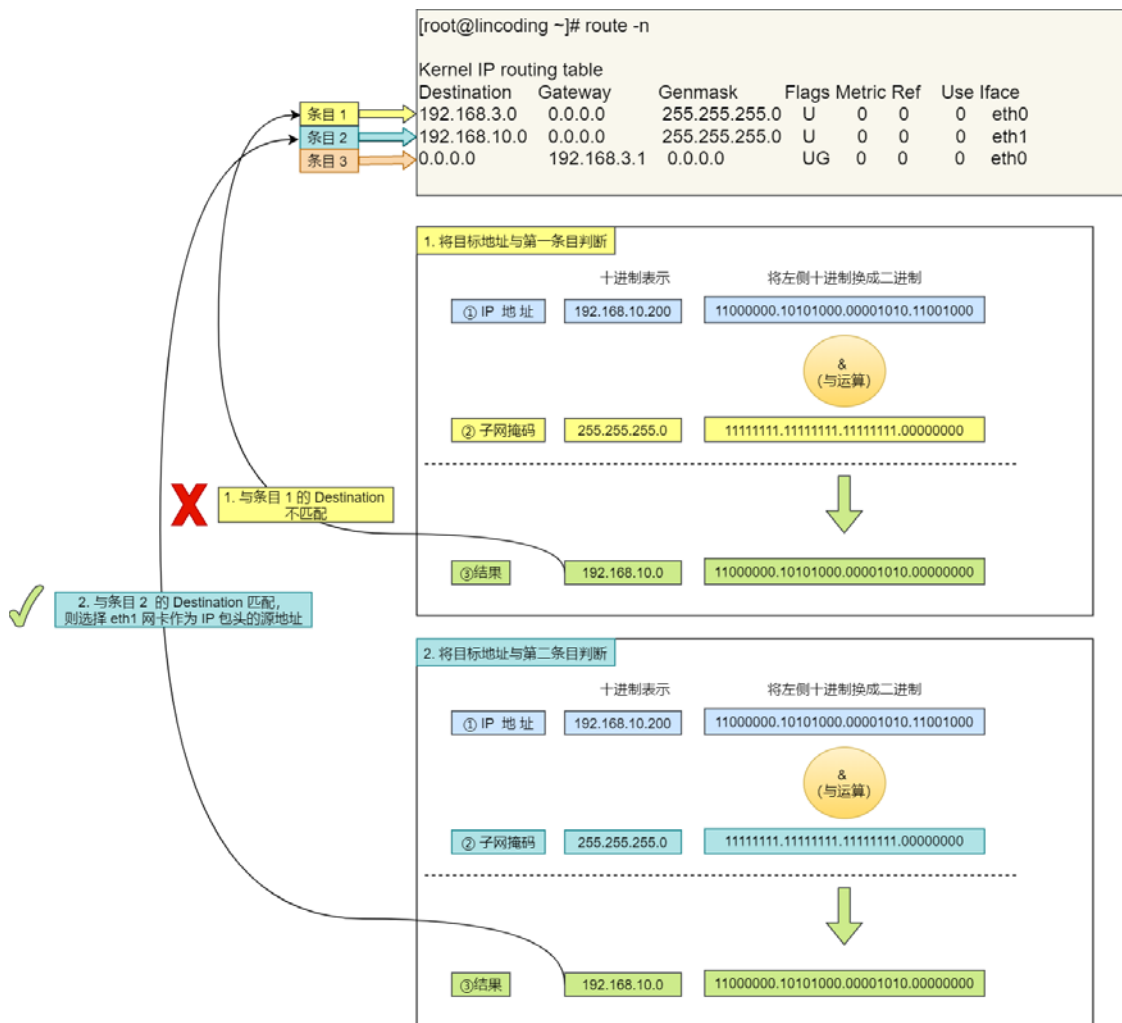
这个时候就需要根据**路由表**规则，来判断哪一个网卡作为源地址 IP。

在 Linux 操作系统，我们可以使用 `route -n` 命令查看当前系统的路由表。

```
[root@lincoding ~]# route -n
```

Kernel IP routing table							
Destination	Gateway	Genmask	Flags	Metric	Ref	Use	Iface
192.168.3.0	0.0.0.0	255.255.255.0	U	0	0	0	eth0
192.168.10.0	0.0.0.0	255.255.255.0	U	0	0	0	eth1
0.0.0.0	192.168.3.1	0.0.0.0	UG	0	0	0	eth0

举个例子，根据上面的路由表，我们假设 Web 服务器的目标地址是 **192.168.10.200**。



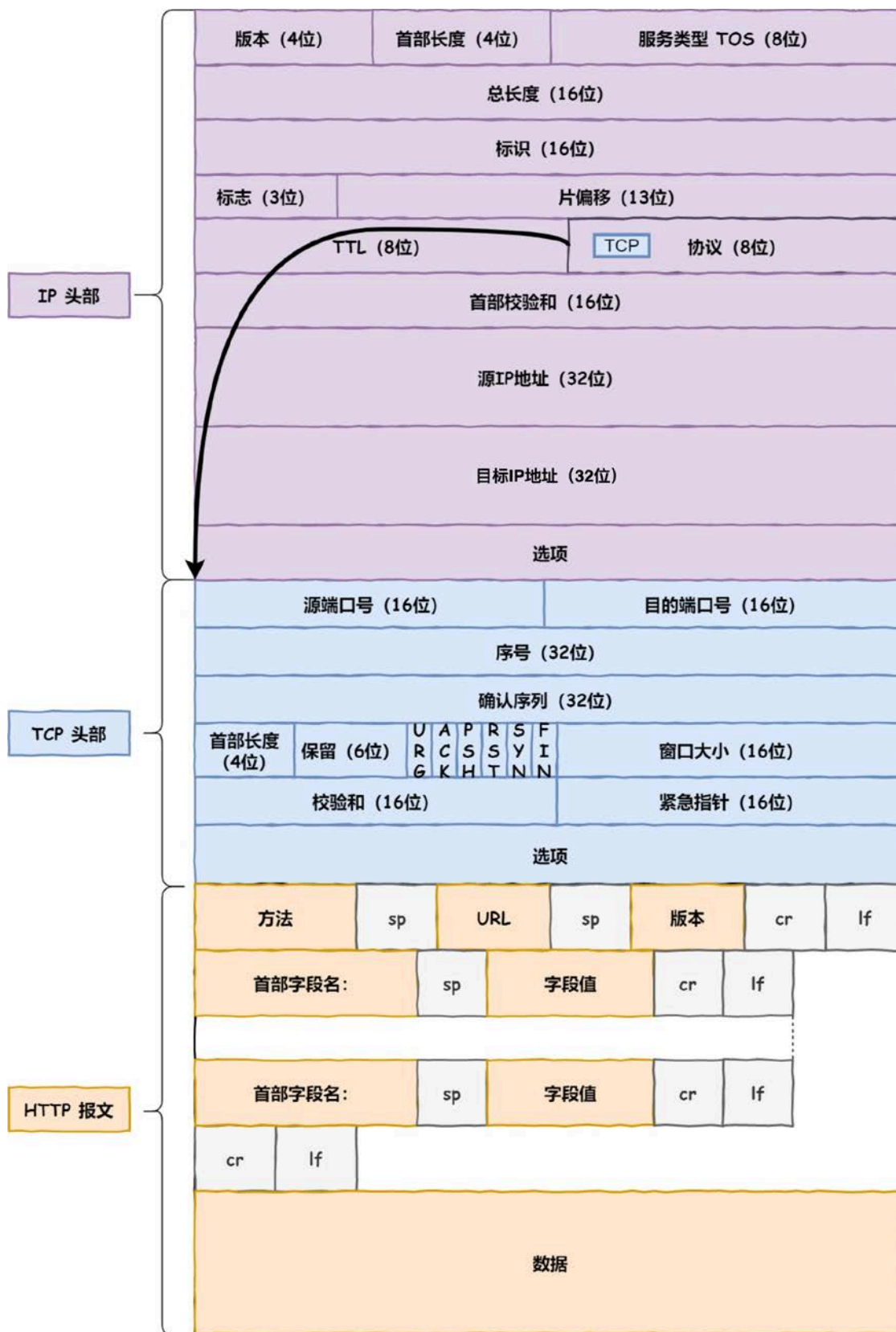
1. 首先和第一条目的子网掩码（Genmask）进行 **与运算**，得到结果为 **192.168.10.0**，但是第一个条目的 Destination 是 **192.168.3.0**，两者不一致所以匹配失败。
2. 再与第二条目的子网掩码进行 **与运算**，得到的结果为 **192.168.10.0**，与第二条目的 Destination **192.168.10.0** 匹配成功，所以将使用 **eth1** 网卡的 IP 地址作为 IP 包头的源地址。

那么假设 Web 服务器的目标地址是 **10.100.20.100**，那么依然依照上面的路由表规则判断，判断后的结果是和第三条目匹配。

第三条目比较特殊，它目标地址和子网掩码都是 **0.0.0.0**，这表示**默认网关**，如果其他所有条目都无法匹配，就会自动匹配这一行。并且后续就把包发给路由器，Gateway 即是路由器的 IP 地址。

### IP 报文生成

至此，网络包的报文如下图。



此时，加上了 IP 头部的数据包表示：“有 IP 大佬给我指路了，感谢 IP 层给我加上了 IP 包头，让我有了远程定位的能力，不会害怕在浩瀚的互联网迷茫了！可是目的地好远啊，我下一站应该去哪呢？”

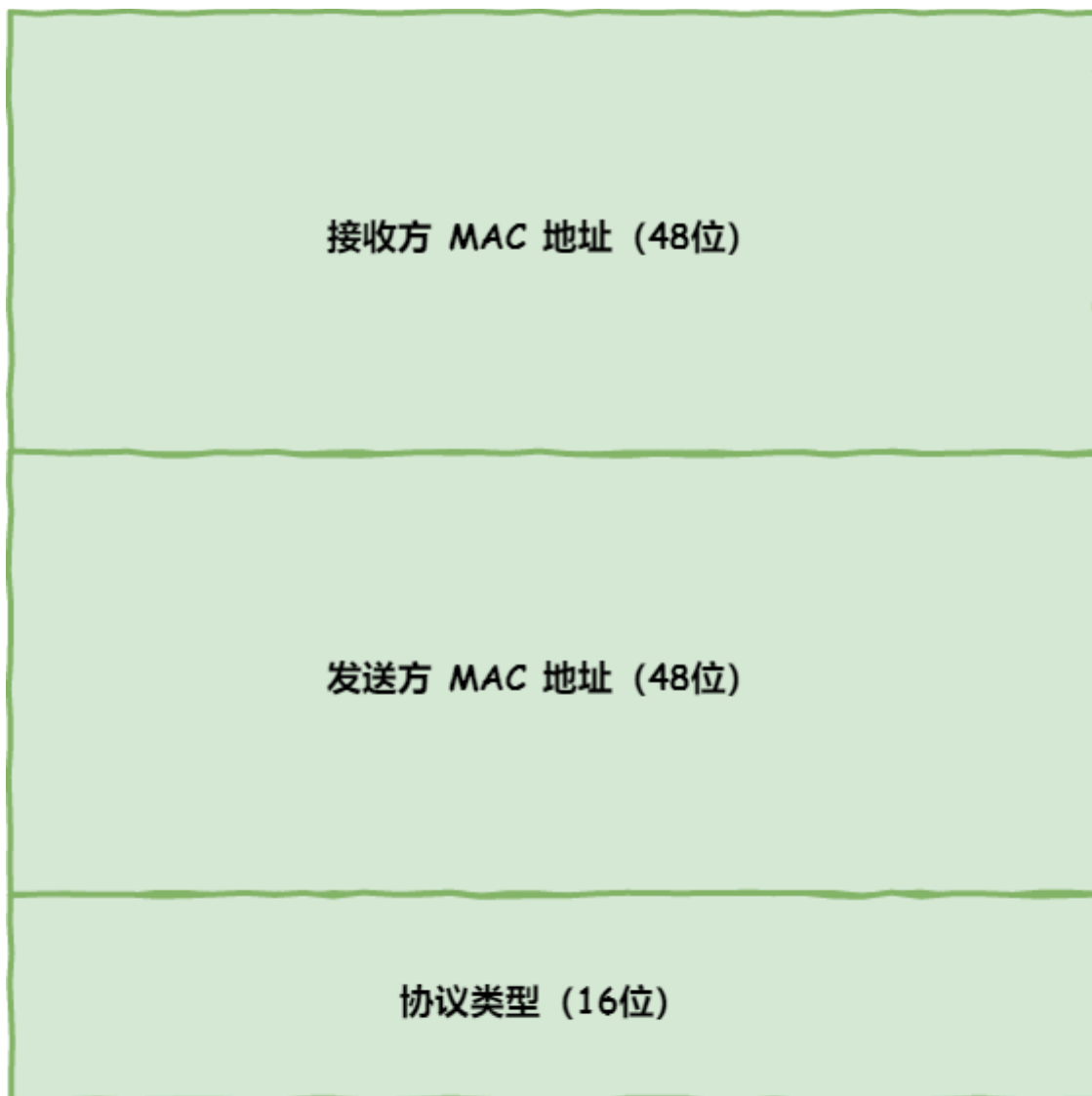
---

## 两点传输 —— MAC

生成了 IP 头部之后，接下来网络包还需要在 IP 头部的前面加上 **MAC 头部**。

### MAC 包头格式

MAC 头部是以太网使用的头部，它包含了接收方和发送方的 MAC 地址等信息。



在 MAC 包头里需要发送方 **MAC 地址**和接收方目标 **MAC 地址**，用于两点之间的传输。

一般在 TCP/IP 通信里，MAC 包头的协议类型只使用：

- 0800：IP 协议
- 0806：ARP 协议

MAC 发送方和接收方如何确认？

发送方的 MAC 地址获取就比较简单了，MAC 地址是在网卡生产时写入到 ROM 里的，只要将这个值读取出来写入到 MAC 头部就可以了。

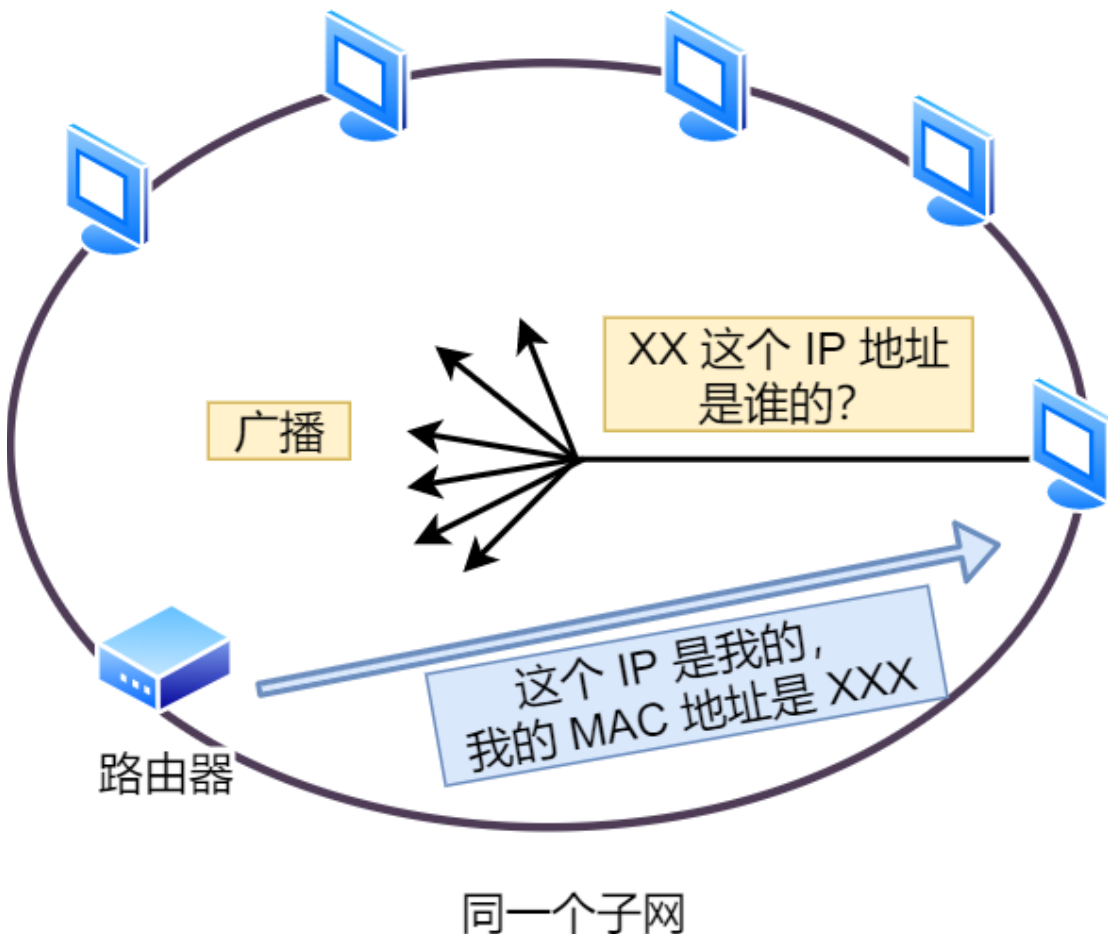
接收方的 MAC 地址就有点复杂了，只要告诉以太网对方的 MAC 的地址，以太网就会帮我们把包发送过去，那么很显然这里应该填写对方的 MAC 地址。

所以先得搞清楚应该把包发给谁，这个只要查一下路由表就知道了。在路由表中找到相匹配的条目，然后把包发给 Gateway 列中的 IP 地址就可以了。

既然知道要发给谁，那如何获取对方的 MAC 地址呢？

不知道对方 MAC 地址？不知道就喊呗。

此时就需要 ARP 协议帮我们找到路由器的 MAC 地址。



ARP 协议会在以太网中以**广播**的形式，对以太网所有的设备喊出：“这个 IP 地址是谁的？请把你的 MAC 地址告诉我”。

然后就会有人回答：“这个 IP 地址是我的，我的 MAC 地址是 XXXX”。

如果对方和自己处于同一个子网中，那么通过上面的操作就可以得到对方的 MAC 地址。然后，我们将这个 MAC 地址写入 MAC 头部，MAC 头部就完成了。

好像每次都要广播获取，这不是很麻烦吗？

放心，在后续操作系统会把本次查询结果放到一块叫做 **ARP 缓存** 的内存空间留着以后用，不过缓存的时间就几分钟。

也就是说，在发包时：

- 先查询 ARP 缓存，如果其中已经保存了对方的 MAC 地址，就不需要发送 ARP 查询，直接使用 ARP 缓存中的地址。
- 而当 ARP 缓存中不存在对方 MAC 地址时，则发送 ARP 广播查询。

查看 ARP 缓存内容

在 Linux 系统中，我们可以使用 `arp -a` 命令来查看 ARP 缓存的内容。

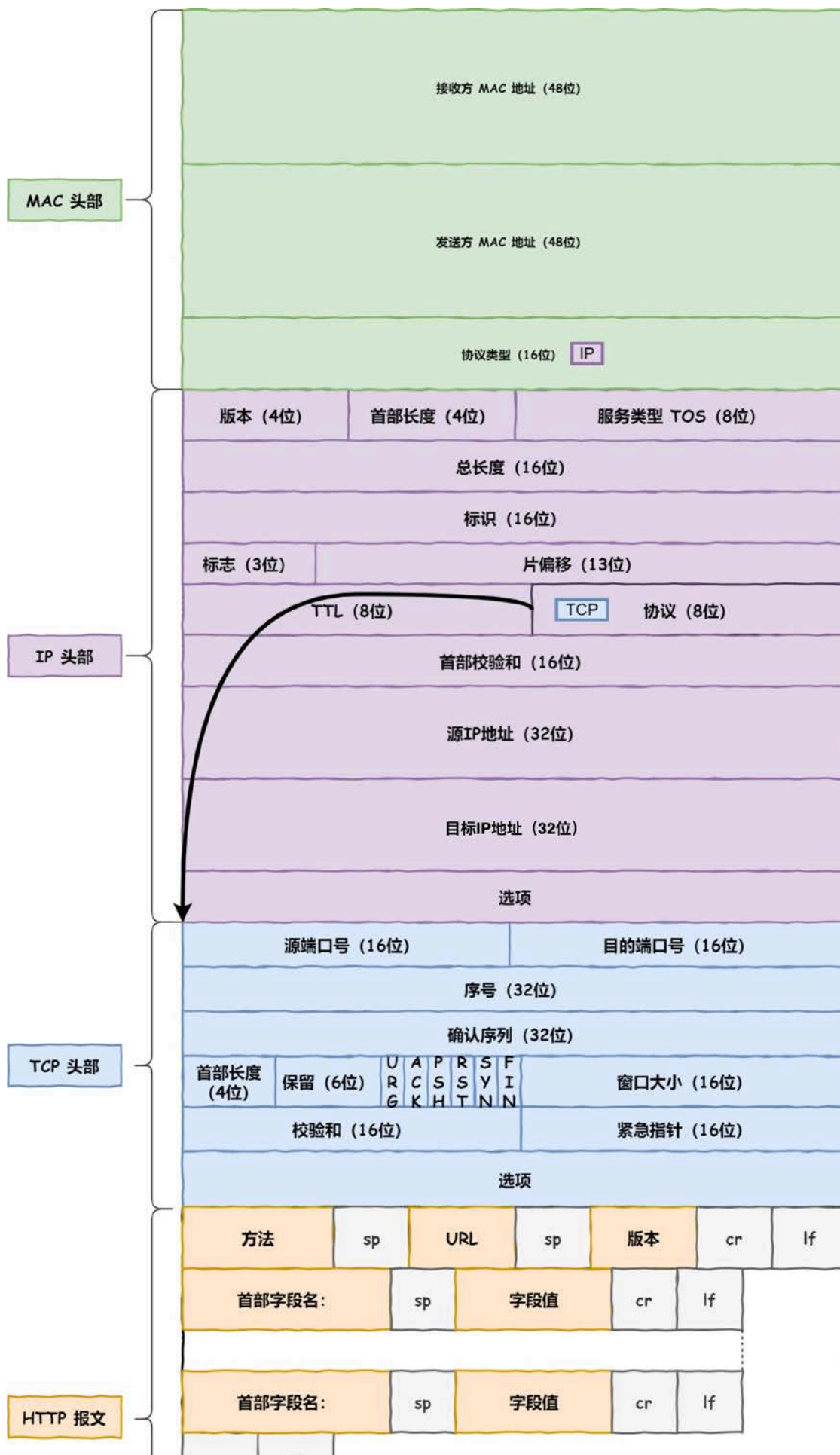
```
[root@lincoding ~]# arp -a  
? (192.168.3.20) at f0:76:1c:58:f4:bc [ether] on eth2
```

IP地址                  MAC 地址                  网口名称

MAC 报文生成

至此，网络包的报文如下图。







此时，加上了 MAC 头部的数据包万分感谢，说道：“感谢 MAC 大佬，我知道我下一步要去哪了！我现在有很多头部兄弟，相信我可以到达最终的目的地！”。

带着众多头部兄弟的数据包，终于准备要出门了。

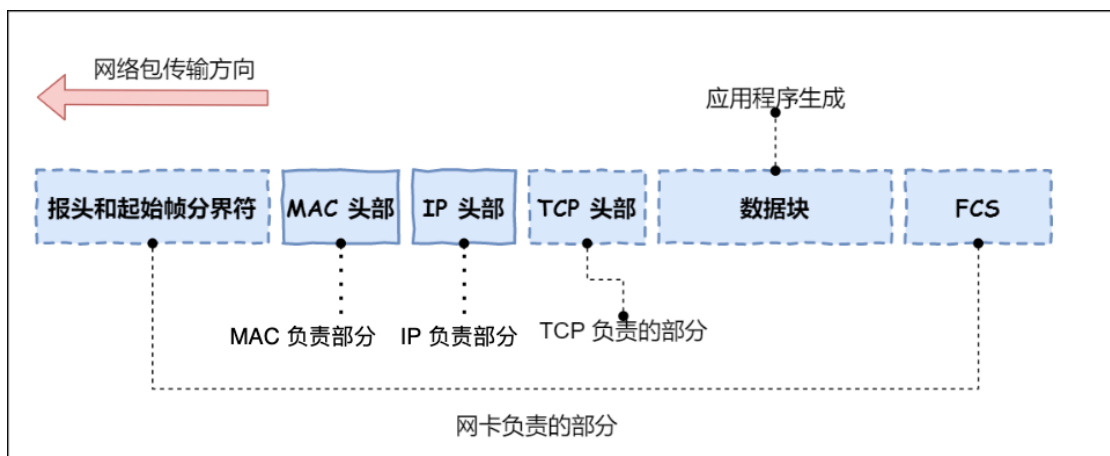
---

## 出口 —— 网卡

网络包只是存放在内存中的一串二进制数字信息，没有办法直接发送给对方。因此，我们需要将**数字信息转换为电信号**，才能在网线上传输，也就是说，这才是真正的数据发送过程。

负责执行这一操作的是**网卡**，要控制网卡还需要靠**网卡驱动程序**。

网卡驱动获取网络包之后，会将其**复制**到网卡内的缓存区中，接着会在其**开头加上报头和起始帧分界符**，在**末尾加上用于检测错误的帧校验序列**。



- 起始帧分界符是一个用来表示包起始位置的标记
- 末尾的 FCS（帧校验序列）用来检查包传输过程是否有损坏

最后网卡会将包转为电信号，通过网线发送出去。

唉，真是不容易，发一个包，真是历经千辛万苦。至此，一个带有许多头部的数据终于踏上寻找目的地的征途了！

---

## 送别者 —— 交换机

下面来看一下包是如何通过交换机的。交换机的设计是将网络包**原样**转发到目的地。交换机工作在 MAC 层，也称为**二层网络设备**。

交换机的包接收操作

首先，电信号到达网线接口，交换机里的模块进行接收，接下来交换机里的模块将电信号转换为数字信号。

然后通过包末尾的 FCS 校验错误，如果没问题则放到缓冲区。这部分操作基本和计算机的网卡相同，但交换机的工作方式和网卡不同。

计算机的网卡本身具有 MAC 地址，并通过核对收到的包的接收方 MAC 地址判断是不是发给自己的，如果不是发给自己的则丢弃；相对地，交换机的端口不核对接收方 MAC 地址，而是直接接收所有的包并存放到缓冲区中。因此，和网卡不同，**交换机的端口不具有 MAC 地址。**

将包存入缓冲区后，接下来需要查询一下这个包的接收方 MAC 地址是否已经在 MAC 地址表中有记录了。

交换机的 MAC 地址表主要包含两个信息：

- 一个是设备的 MAC 地址，
- 另一个是该设备连接在交换机的哪个端口上。

交换机内部有一张 MAC 地址与网线端口的映射表。  
当接收到包时，会将相应的端口号和发送 MAC 地址写入表中，  
这样就可以根据地址判断出该设备连接在哪个端口上了。  
交换机就是根据这些信息判断，应该把包转发到哪儿的。

MAC 地址表	端口	控制信息
00-60-97-A5-43-3C	1	...
00-00-C0-16-AE-FD	2	...
00-02-B3-1C-9C-F9	3	...
....	...	...

交换机

举个例子，如果收到的包的接收方 MAC 地址为 00-02-B3-1C-9C-F9，则与图中表中的第 3 行匹配，根据端口列的信息，可知这个地址位于 3 号端口上，然后就可以通过交换电路将包发送到相应的端口了。

所以，交换机根据 MAC 地址表查找 MAC 地址，然后将信号发送到相应的端口。

当 MAC 地址表找不到指定的 MAC 地址会怎么样？

地址表中找不到指定的 MAC 地址。这可能是因为有该地址的设备还没有向交换机发送过包，或者这个设备一段时间没有工作导致地址被从地址表中删除了。

这种情况下，交换机无法判断应该把包转发到哪个端口，只能将包转发到除了源端口之外的所有端口上，无论该设备连接在哪个端口上都能收到这个包。

这样做不会产生什么问题，因为以太网的设计本来就是将包发送到整个网络的，然后只有相应的接收者才接收包，而其他设备则会忽略这个包。

有人会说：“这样做会发送多余的包，会不会造成网络拥塞呢？”

其实完全不用过于担心，因为发送了包之后目标设备会作出响应，只要返回了响应包，交换机就可以将它的地址写入 MAC 地址表，下次也就不需要把包发到所有端口了。

局域网中每秒可以传输上千个包，多出一两个包并无大碍。

此外，如果接收方 MAC 地址是一个**广播地址**，那么交换机会将包发送到除源端口之外的所有端口。

以下两个属于广播地址：

- MAC 地址中的 FF:FF:FF:FF:FF:FF
- IP 地址中的 255.255.255.255

数据包通过交换机转发抵达了路由器，准备要离开土生土长的子网了。此时，数据包和交换机离别时说道：“感谢交换机兄弟，帮我转发到出境的大门，我要出远门啦！”

---

## 出境大门 —— 路由器

### 路由器与交换机的区别

网络包经过交换机之后，现在到达了**路由器**，并在此被转发到下一个路由器或目标设备。

这一步转发的工作原理和交换机类似，也是通过查表判断包转发的目标。

不过在具体的操作过程上，路由器和交换机是有区别的。

- 因为**路由器**是基于 IP 设计的，俗称**三层**网络设备，路由器的各个端口都具有 MAC 地址和 IP 地址；
- 而**交换机**是基于以太网设计的，俗称**二层**网络设备，交换机的端口不具有 MAC 地址。

### 路由器基本原理

路由器的端口具有 MAC 地址，因此它能够成为以太网的发送方和接收方；同时还具有 IP 地址，从这个意义上来说，它和计算机的网卡是一样的。

当转发包时，首先路由器端口会接收发给自己的以太网包，然后路由表查询转发目标，再由相应的端口作为发送方将以太网包发送出去。

路由器的包接收操作

首先，电信号到达网线接口部分，路由器中的模块会将电信号转成数字信号，然后通过包末尾的 FCS 进行错误校验。

如果没问题则检查 MAC 头部中的接收方 MAC 地址，看看是不是发给自己的包，如果是就放到接收缓冲区中，否则就丢弃这个包。

总的来说，路由器的端口都具有 MAC 地址，只接收与自身地址匹配的包，遇到不匹配的包则直接丢弃。

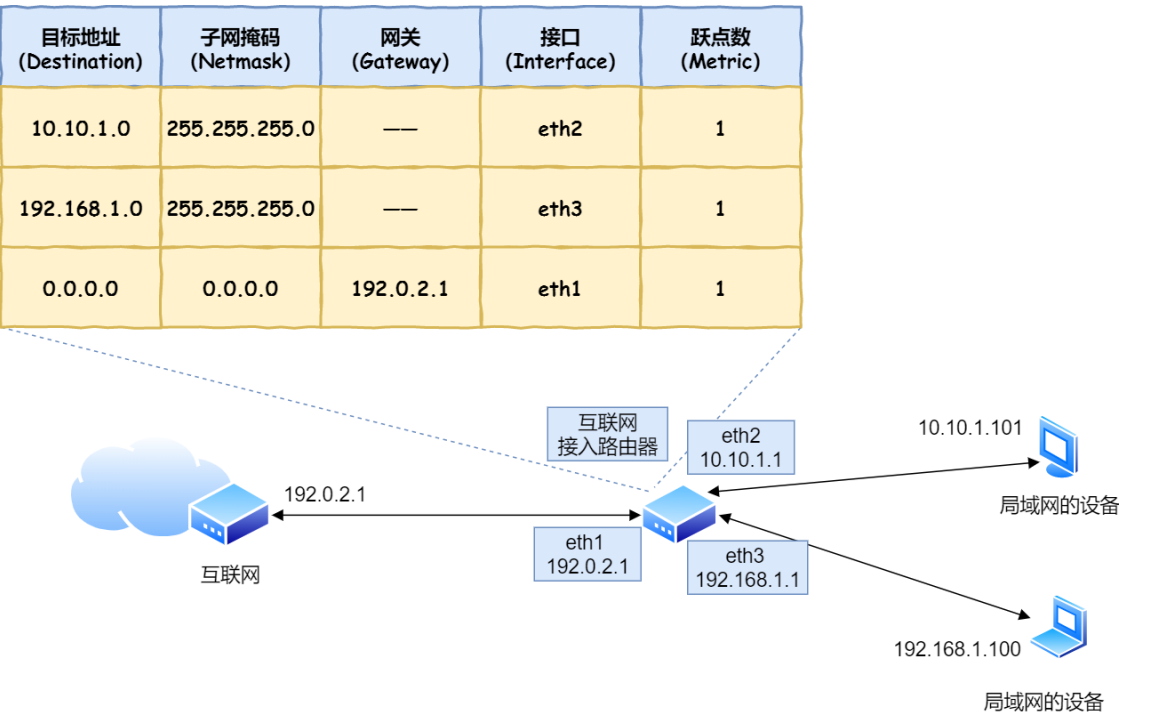
查询路由表确定输出端口

完成包接收操作之后，路由器就会去掉包开头的 MAC 头部。

MAC 头部的作用就是将包送达路由器，其中的接收方 MAC 地址就是路由器端口的 MAC 地址。因此，当包到达路由器之后，MAC 头部的任务就完成了，于是 MAC 头部就会被丢弃。

接下来，路由器会根据 MAC 头部后方的 IP 头部中的内容进行包的转发操作。

转发操作分为几个阶段，首先是查询路由表判断转发目标。



具体的工作流程根据上图，举个例子。

假设地址为 **10.10.1.101** 的计算机要向地址为 **192.168.1.100** 的服务器发送一个包，这个包先到达图中的路由器。

判断转发目标的第一步，就是根据包的接收方 IP 地址查询路由表中的目标地址栏，以找到相匹配的记录。

路由匹配和前面讲的一样，每个条目的子网掩码和 **192.168.1.100** IP 做 **& 与运算** 后，得到的结果与对应条目的目标地址进行匹配，如果匹配就会作为候选转发目标，如果不匹配就继续与下个条目进行路由匹配。

如第二条目的子网掩码 **255.255.255.0** 与 **192.168.1.100** IP 做 **& 与运算** 后，得到结果是 **192.168.1.0**，这与第二条目的目标地址 **192.168.1.0** 匹配，该第二条目记录就会被作为转发目标。

实在找不到匹配路由时，就会选择**默认路由**，路由表中子网掩码为 **0.0.0.0** 的记录表示「默认路由」。

#### 路由器的发送操作

接下来就会进入包的**发送操作**。

首先，我们需要根据**路由表的网关列**判断对方的地址。

- 如果网关是一个 IP 地址，则这个 IP 地址就是我们要转发到的目标地址，**还未抵达终点**，还需继续需要路由器转发。
- 如果网关为空，则 IP 头部中的接收方 IP 地址就是要转发到的目标地址，也是就终于找到 IP 包头里的目标地址了，说明**已抵达终点**。

知道对方的 IP 地址之后，接下来需要通过 **ARP** 协议根据 IP 地址查询 **MAC** 地址，并将查询的结果作为接收方 **MAC** 地址。

路由器也有 **ARP** 缓存，因此首先会在 **ARP** 缓存中查询，如果找不到则发送 **ARP** 查询请求。

接下来是发送方 **MAC** 地址字段，这里填写输出端口的 **MAC** 地址。还有一个以太类型字段，填写 **0800**（十六进制）表示 IP 协议。

网络包完成后，接下来会将其转换成电信号并通过端口发送出去。这一步的工作过程和计算机也是相同的。

发送出去的网络包会通过**交换机**到达下一个路由器。由于接收方 **MAC** 地址就是下一个路由器的地址，所以交换机会根据这一地址将包传输到下一个路由器。

接下来，下一个路由器会将包转发给再下一个路由器，经过层层转发之后，网络包就到达了最终的目的地。

不知你发现了没有，在网络包传输的过程中，**源 IP 和目标 IP 始终是不会变的**，一直变化的是 **MAC 地址**，因为需要 MAC 地址在以太网内进行**两个设备**之间的包传输。

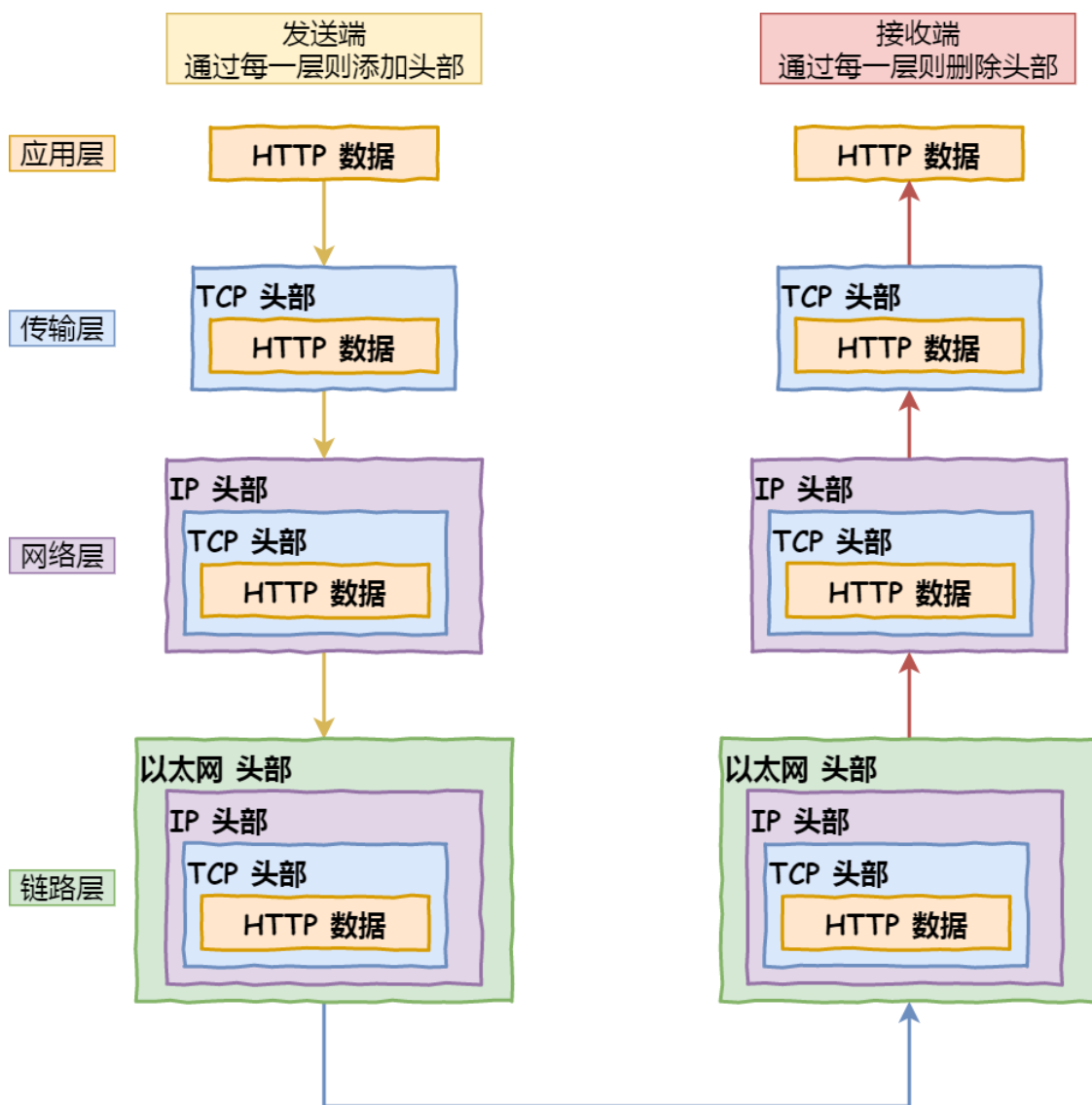
数据包通过多个路由器道友的帮助，在网络世界途经了很多路程，最终抵达了目的地的城门！城门值守的路由器，发现了这个小兄弟数据包原来是找城内的人，于是它就将数据包送进了城内，再经由城内的交换机帮助下，最终转发到了目的地了。数据包感慨万千的说道：“多谢这一路上，各路大侠的相助！”

---

## 互相扒皮 —— 服务器 与 客户端

数据包抵达了服务器，服务器肯定高兴呀，正所谓有朋自远方来，不亦乐乎？

服务器高兴的不得了，于是开始扒数据包的皮！就好像你收到快递，能不兴奋吗？



数据包抵达服务器后，服务器会先扒开数据包的 MAC 头部，查看是否和服务器自己的 MAC 地址符合，符合就将包收起来。

接着继续扒开数据包的 IP 头，发现 IP 地址符合，根据 IP 头中协议项，知道自己上层是 TCP 协议。

于是，扒开 TCP 的头，里面有序列号，需要看一看这个序列包是不是我想要的，如果是就放入缓存中然后返回一个 ACK，如果不是就丢弃。TCP 头部里面还有端口号，HTTP 的服务器正在监听这个端口号。

于是，服务器自然就知道是 HTTP 进程想要这个包，于是就将包发给 HTTP 进程。

服务器的 HTTP 进程看到，原来这个请求是要访问一个页面，于是就把这个网页封装在 HTTP 响应报文里。

HTTP 响应报文也需要穿上 TCP、IP、MAC 头部，不过这次是源地址是服务器 IP 地址，目的地址是客户端 IP 地址。

穿好头部衣服后，从网卡出去，交由交换机转发到出城的路由器，路由器就把响应数据包发到了下一个路由器，就这样跳啊跳。

最后跳到了客户端的城门把守的路由器，路由器扒开 IP 头部发现是要找城内的人，于是又把包发给了城内的交换机，再由交换机转发到客户端。

客户端收到了服务器的响应数据包后，同样也非常的高兴，客户能拆快递了！

于是，客户端开始扒皮，把收到的数据包的皮扒剩 HTTP 响应报文后，交给浏览器去渲染页面，一份特别的数据包快递，就这样显示出来了！

最后，客户端要离开了，向服务器发起了 TCP 四次挥手，至此双方的连接就断开了。

---

## 一个数据包臭不要脸的感受

下面内容的「我」，代表「臭美的数据包角色」。注：（括号的内容）代表我的吐槽，三连呸！

我一开始我虽然孤单、不知所措，但没有停滞不前。我依然满怀信心和勇气开始了征途。（你当然有勇气，你是应用层数据，后面有底层兄弟当靠山，我呸！）

我很庆幸遇到了各路神通广大的大佬，有可靠传输的 TCP、有远程定位功能的 IP、有指明下一站位置的 MAC 等（你当然会遇到，因为都被计算机安排好的，我呸！）。

这些大佬都给我前面加上了头部，使得我能在交换机和路由器的转发下，抵达了目的地！（哎，你也不容易，不吐槽了，放过你！）

这一路上的经历，让我认识到了网络世界中各路大侠协作的重要性，是他们维护了网络世界的秩序，感谢他们！（我呸，你应该感谢众多计算机科学家！）

---

## 参考资料

[1] 户根勤。网络是怎么连接的。人民邮电出版社。

[2] 刘超。趣谈网络协议。极客时间。。

---



## 读者问答

读者问：“笔记本的是自带交换机的吗？交换机现在我还不知道是什么”

笔记本不是交换机，交换机通常是 2 个网口以上。

现在家里的路由器其实有了交换机的功能了。交换机可以简单理解成一个设备，三台电脑网线接到这个设备，这三台电脑就可以互相通信了，交换机嘛，交换数据这么理解就可以。

读者问：“如果知道你电脑的 Mac 地址，我可以直接给你发消息吗？”

Mac 地址只能是两个设备之间传递时使用的，如果你要从大老远给我发消息，是离不开 IP 的。

读者问：“请问公网服务器的 Mac 地址是在什么时机通过什么方式获取到的？我看 ARP 获取 Mac 地址只能获取到内网机器的 Mac 地址吧？”

在发送数据包时，如果目标主机不是本地局域网，填入的 Mac 地址是路由器，也就是把数据包转发给路由器，路由器一直转发下一个路由器，直到转发到目标主机的路由器，发现 IP 地址是自己局域网内的主机，就会 ARP 请求获取目标主机的 Mac 地址，从而转发到这个服务器主机。

转发的过程中，源 IP 地址和目标 IP 地址是不会变的（前提：没有使用 NAT 网络的），源 MAC 地址和目标 MAC 地址是会变化的。

---

哈喽，我是小林，就爱图解计算机基础，如果觉得文章对你有帮助，别忘记关注我哦！



扫一扫，关注「小林coding」公众号

图解计算机基础  
认准**小林coding**

每一张图都包含小林的认真  
只为帮助大家能更好的理解

- ① 关注公众号回复「**图解**」  
获取图解系列 PDF
- ② 关注公众号回复「**加群**」  
拉你进百人技术交流群

## 2.3 Linux 系统是如何收发网络包的？

这次，就围绕一个问题来说。

### Linux 系统是如何收发网络包的？

#### 网络模型

为了使得多种设备能通过网络相互通信，和为了解决各种不同设备在网络互联中的兼容性问题，国际标准化组织制定了开放式系统互联通信参考模型（*Open System Interconnection Reference Model*），也就是 OSI 网络模型，该模型主要有 7 层，分别是应用层、表示层、会话层、传输层、网络层、数据链路层以及物理层。

每一层负责的职能都不同，如下：

- 应用层，负责给应用程序提供统一的接口；
- 表示层，负责把数据转换成兼容另一个系统能识别的格式；
- 会话层，负责建立、管理和终止表示层实体之间的通信会话；
- 传输层，负责端到端的数据传输；
- 网络层，负责数据的路由、转发、分片；
- 数据链路层，负责数据的封装和差错检测，以及 MAC 寻址；
- 物理层，负责在物理网络中传输数据帧；

由于 OSI 模型实在太复杂，提出的也只是概念理论上的分层，并没有提供具体的实现方案。

事实上，我们比较常见，也比较实用的是四层模型，即 TCP/IP 网络模型，Linux 系统正是按照这套网络模型来实现网络协议栈的。

TCP/IP 网络模型共有 4 层，分别是应用层、传输层、网络层和网络接口层，每一层负责的职能如下：

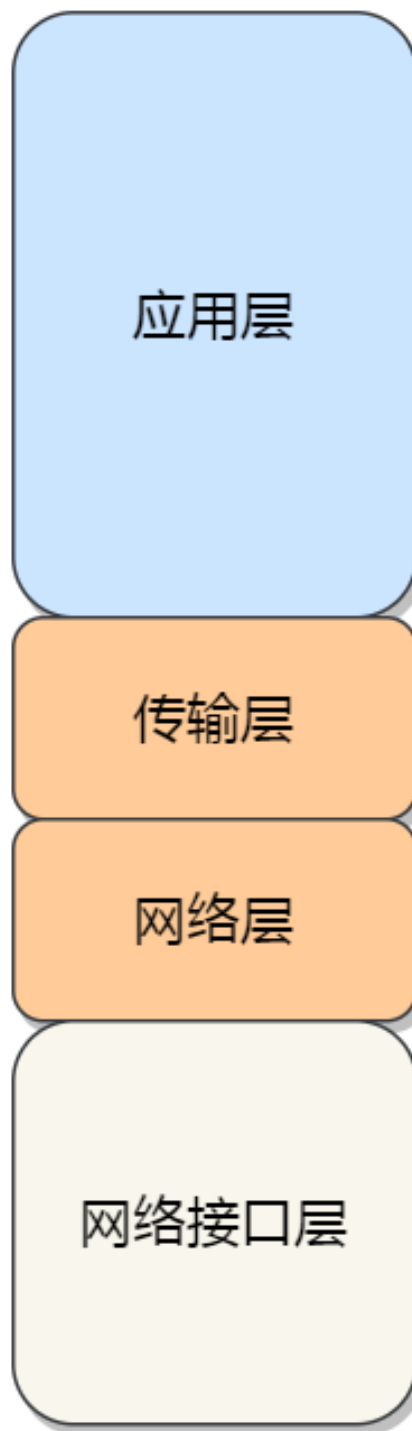
- 应用层，负责向用户提供一组应用程序，比如 HTTP、DNS、FTP 等；
- 传输层，负责端到端的通信，比如 TCP、UDP 等；
- 网络层，负责网络包的封装、分片、路由、转发，比如 IP、ICMP 等；
- 网络接口层，负责网络包在物理网络中的传输，比如网络包的封装、MAC 寻址、差错检测，以及通过网卡传输网络帧等；

TCP/IP 网络模型相比 OSI 网络模型简化了不少，也更加易记，它们之间的关系如下图所示：

## OSI 参考模式



## TCP/IP 模型



不过，我们常说的七层和四层负载均衡，是用 OSI 网络模型来描述的，七层对应的是应用层，四层对应的是传输层。

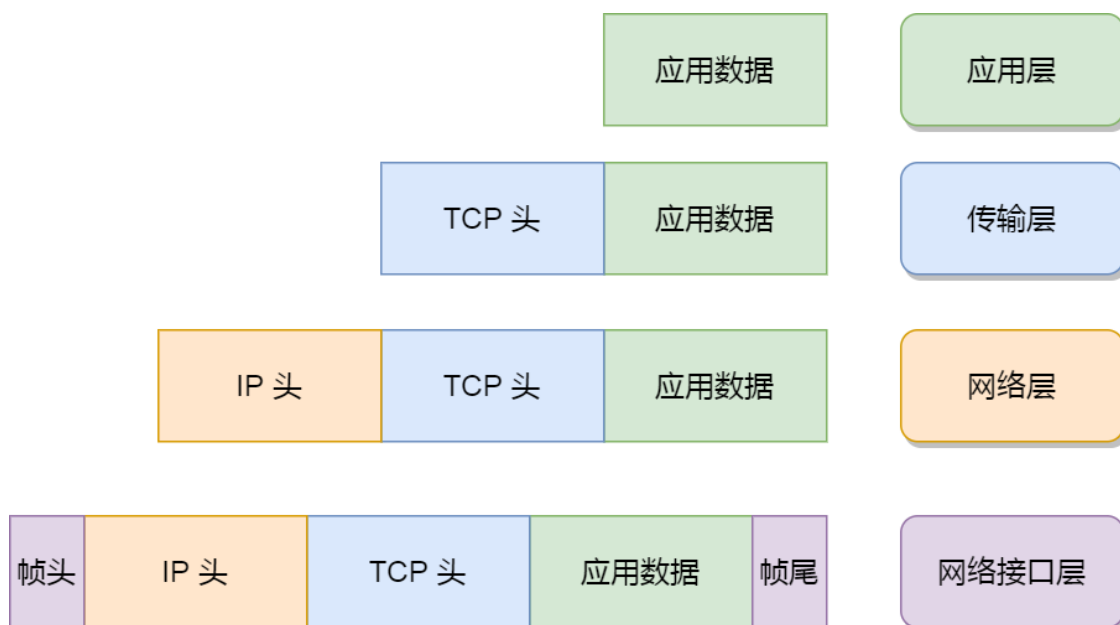
---

## Linux 网络协议栈

我们可以把自己的身体比作应用层中的数据，打底衣服比作传输层中的 TCP 头，外套比作网络层中 IP 头，帽子和鞋子分别比作网络接口层的帧头和帧尾。

在冬天这个季节，当我们要从家里出去玩的时候，自然要先穿个打底衣服，再套上保暖外套，最后穿上帽子和鞋子才出门，这个过程就好像我们把 TCP 协议通信的网络包发出去的时候，会把应用层的数据按照网络协议栈层层封装和处理。

你从下面这张图可以看到，应用层数据在每一层的封装格式。



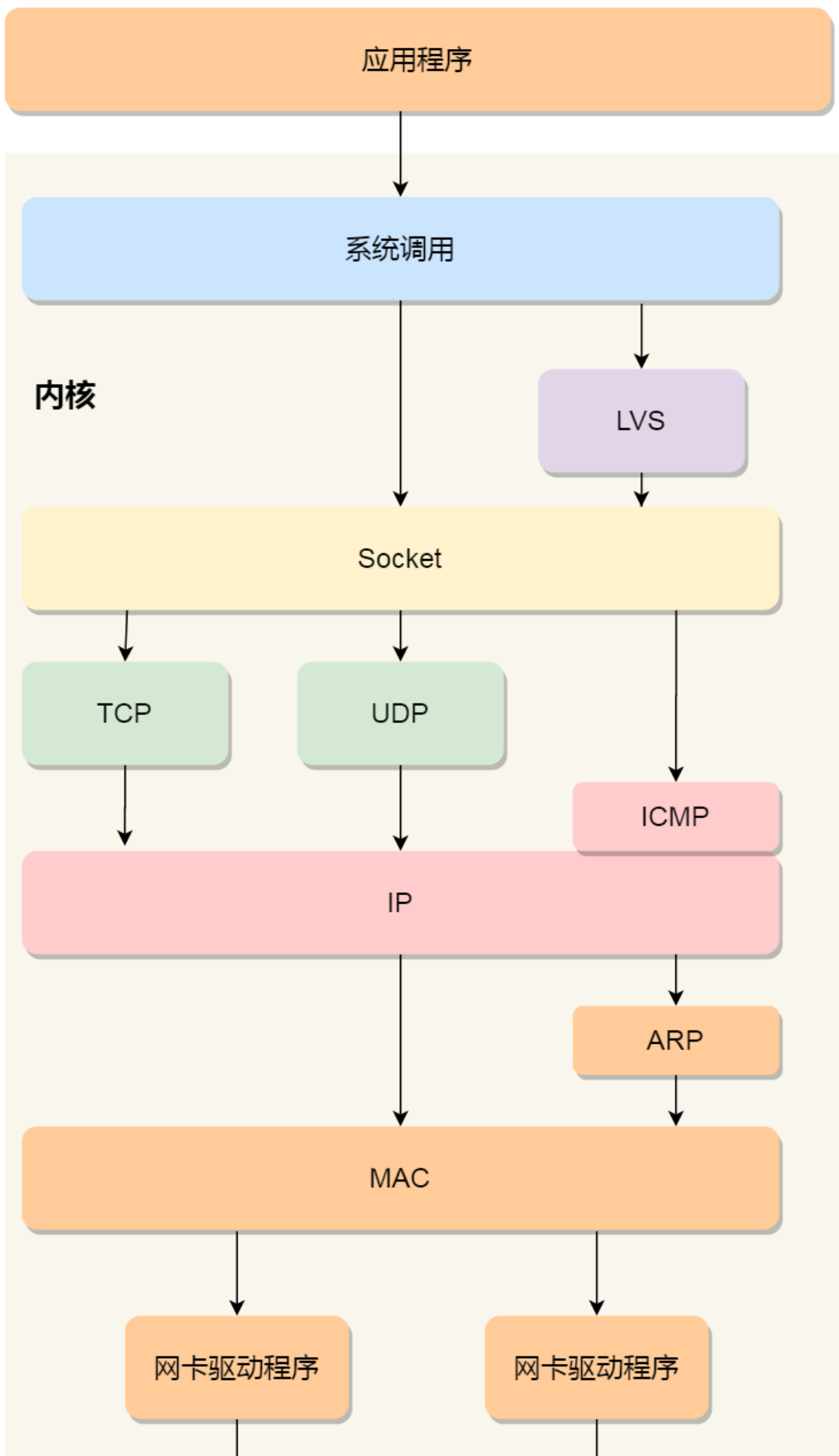
其中：

- 传输层，给应用数据前面增加了 TCP 头；
- 网络层，给 TCP 数据包前面增加了 IP 头；
- 网络接口层，给 IP 数据包前后分别增加了帧头和帧尾；

这些新增的头部和尾部，都有各自的作用，也都是按照特定的协议格式填充，这每一层都增加了各自的协议头，那自然网络包的大小就增大了，但物理链路并不能传输任意大小的数据包，所以在以太网中，规定了最大传输单元（MTU）是 1500 字节，也就是规定了单次传输的最大 IP 包大小。

当网络包超过 MTU 的大小，就会在网络层分片，以确保分片后的 IP 包不会超过 MTU 大小，如果 MTU 越小，需要的分包就越多，那么网络吞吐能力就越差，相反的，如果 MTU 越大，需要的分包就越少，那么网络吞吐能力就越好。

知道了 TCP/IP 网络模型，以及网络包的封装原理后，那么 Linux 网络协议栈的样子，你想必猜到了大概，它其实就类似于 TCP/IP 的四层结构：



从上图的网络协议栈，你可以看到：

- 应用程序需要通过系统调用，来跟 **Socket** 层进行数据交互；
- **Socket** 层的下面就是传输层、网络层和网络接口层；
- 最下面的一层，则是网卡驱动程序和硬件网卡设备；

## Linux 接收网络包的流程

网卡是计算机里的一个硬件，专门负责接收和发送网络包，当网卡接收到一个网络包后，会通过 **DMA** 技术，将网络包写入到指定的内存地址，也就是写入到 **Ring Buffer**，这个是一个环形缓冲区，接着就会告诉操作系统这个网络包已经到达。

那应该怎么告诉操作系统这个网络包已经到达了呢？

最简单的一种方式就是触发中断，也就是每当网卡收到一个网络包，就触发一个中断告诉操作系统。

但是，这存在一个问题，在高性能网络场景下，网络包的数量会非常多，那么就会触发非常多的中断，要知道当 **CPU** 收到了中断，就会停下手里的事情，而去处理这些网络包，处理完毕后，才会回去继续其他事情，那么频繁地触发中断，则会导致 **CPU** 一直没完没了的处理中断，而导致其他任务可能无法继续前进，从而影响系统的整体效率。

所以为了解决频繁中断带来的性能开销，Linux 内核在 2.6 版本中引入了 **NAPI 机制**，它是混合「中断和轮询」的方式来接收网络包，它的核心概念就是**不采用中断的方式读取数据**，而是首先采用中断唤醒数据接收的服务程序，然后 **poll** 的方法来轮询数据。

因此，当有网络包到达时，会通过 **DMA** 技术，将网络包写入到指定的内存地址，接着网卡向 **CPU** 发起硬件中断，当 **CPU** 收到硬件中断请求后，根据中断表，调用已经注册的中断处理函数。

硬件中断处理函数会做如下的事情：

- 需要先「暂时屏蔽中断」，表示已经知道内存中有数据了，告诉网卡下次再收到数据包直接写内存就可以了，不要再通知 **CPU** 了，这样可以提高效率，避免 **CPU** 不停的被中断。
- 接着，发起「软中断」，然后恢复刚才屏蔽的中断。

至此，硬件中断处理函数的工作就已经完成。

硬件中断处理函数做的事情很少，主要耗时的工作都交给软中断处理函数了。

软中断的处理

内核中的 `ksoftirqd` 线程专门负责软中断的处理，当 `ksoftirqd` 内核线程收到软中断后，就会来轮询处理数据。

`ksoftirqd` 线程会从 Ring Buffer 中获取一个数据帧，用 `sk_buff` 表示，从而可以作为一个网络包交给网络协议栈进行逐层处理。

### 网络协议栈

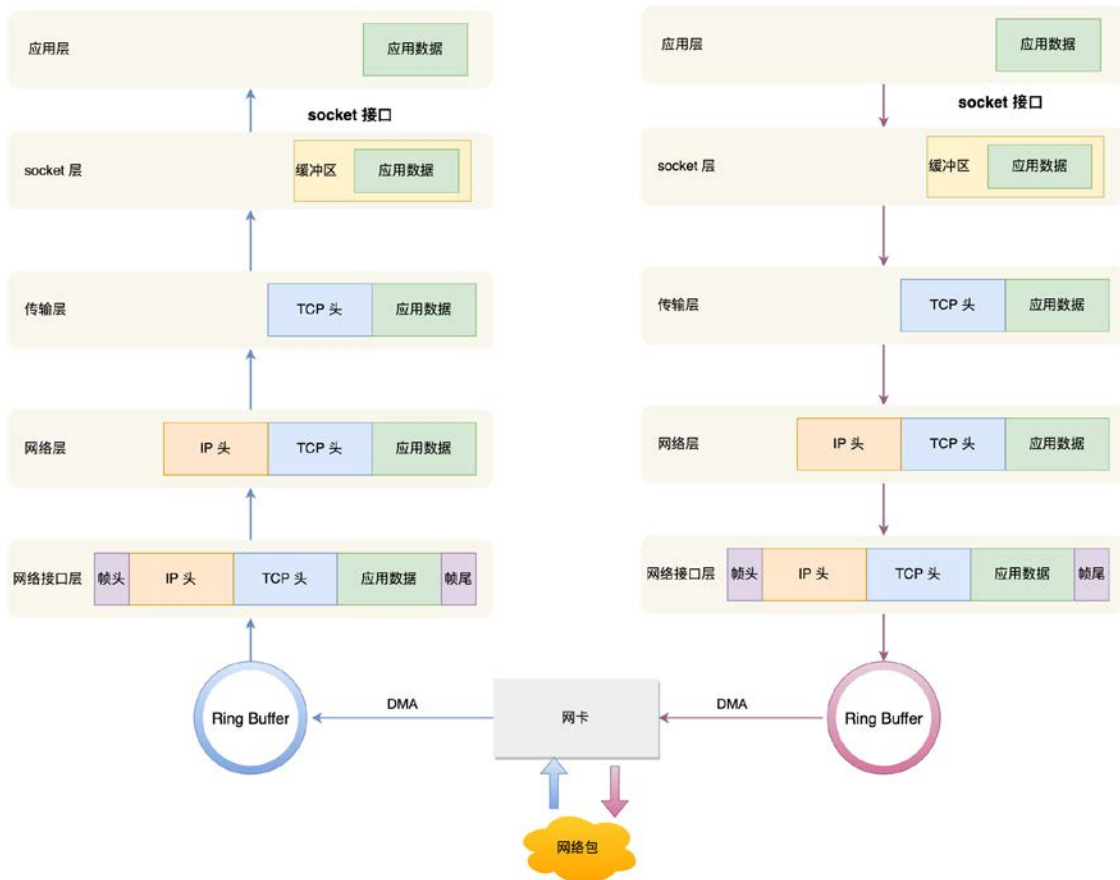
首先，会先进入到网络接口层，在这一层会检查报文的合法性，如果不合法则丢弃，合法则会找出该网络包的上层协议的类型，比如是 IPv4，还是 IPv6，接着再去掉帧头和帧尾，然后交给网络层。

到了网络层，则取出 IP 包，判断网络包下一步的走向，比如是交给上层处理还是转发出去。当确认这个网络包要发送给本机后，就会从 IP 头里看看上一层协议的类型是 TCP 还是 UDP，接着去掉 IP 头，然后交给传输层。

传输层取出 TCP 头或 UDP 头，根据四元组「源 IP、源端口、目的 IP、目的端口」作为标识，找出对应的 Socket，并把数据放到 Socket 的接收缓冲区。

最后，应用层程序调用 Socket 接口，将内核的 Socket 接收缓冲区的数据「拷贝」到应用层的缓冲区，然后唤醒用户进程。

至此，一个网络包的接收过程就已经结束了，你也可以从下图左边部分看到网络包接收的流程，右边部分刚好反过来，它是网络包发送的流程。



## Linux 发送网络包的流程

如上图的右半部分，发送网络包的流程正好和接收流程相反。

首先，应用程序会调用 Socket 发送数据包的接口，由于这个是系统调用，所以会从用户态陷入到内核态中的 Socket 层，内核会申请一个内核态的 `sk_buff` 内存，将用户待发送的数据拷贝到 `sk_buff` 内存，并将其加入到发送缓冲区。

接下来，网络协议栈从 Socket 发送缓冲区中取出 `sk_buff`，并按照 TCP/IP 协议栈从上到下逐层处理。

如果使用的是 TCP 传输协议发送数据，那么先拷贝一个新的 `sk_buff` 副本，这是因为 `sk_buff` 后续在调用网络层，最后到达网卡发送完成的时候，这个 `sk_buff` 会被释放掉。而 TCP 协议是支持丢失重传的，在收到对方的 ACK 之前，这个 `sk_buff` 不能被删除。所以内核的做法就是每次调用网卡发送的时候，实际上传递出去的是 `sk_buff` 的一个拷贝，等收到 ACK 再真正删除。

接着，对 `sk_buff` 填充 TCP 头。这里提一下，`sk_buff` 可以表示各个层的数据包，在应用层数据包叫 `data`，在 TCP 层我们称为 `segment`，在 IP 层我们叫 `packet`，在数据链路层称为 `frame`。

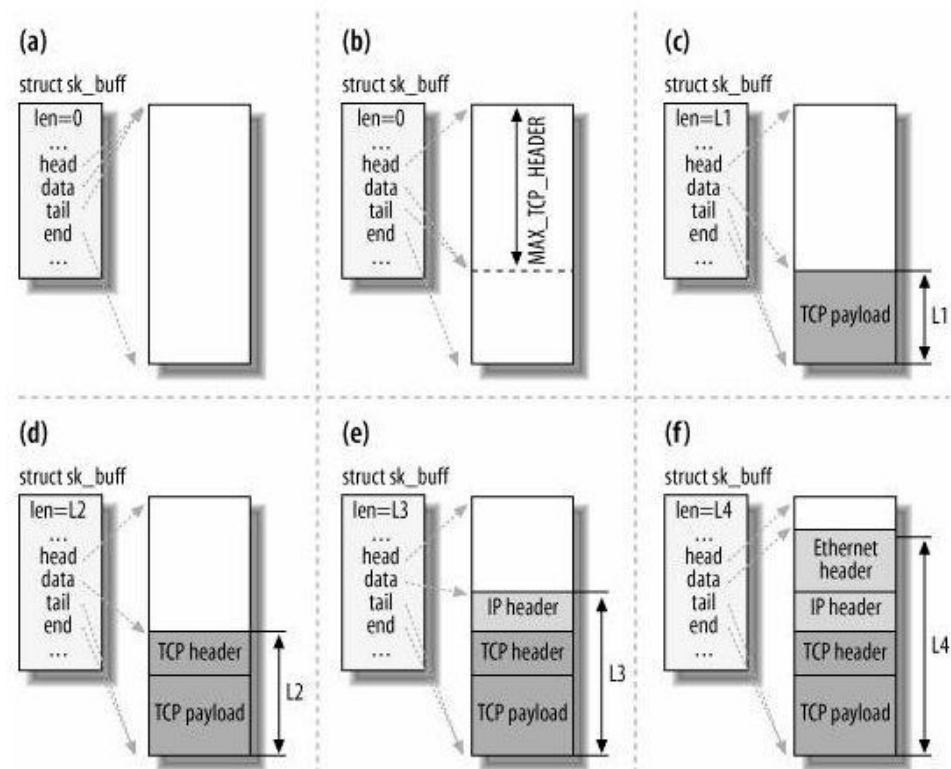


你可能会好奇，为什么全部数据包只用一个结构体来描述呢？协议栈采用的是分层结构，上层向下层传递数据时需要增加包头，下层向上层数据时又需要去掉包头，如果每一层都用一个结构体，那在层之间传递数据的时候，就要发生多次拷贝，这将大大降低 CPU 效率。

于是，为了在层级之间传递数据时，不发生拷贝，只用 `sk_buff` 一个结构体来描述所有的网络包，那它是如何做到的呢？是通过调整 `sk_buff` 中 `data` 的指针，比如：

- 当接收报文时，从网卡驱动开始，通过协议栈层层往上传送数据报，通过增加 `skb->data` 的值，来逐步剥离协议首部。
- 当要发送报文时，创建 `sk_buff` 结构体，数据缓存区的头部预留足够的空间，用来填充各层首部，在经过各下层协议时，通过减少 `skb->data` 的值来增加协议首部。

你可以从下面这张图看到，当发送报文时，`data` 指针的移动过程。



至此，传输层的工作也就都完成了。

然后交给网络层，在网络层里会做这些工作：选取路由（确认下一跳的 IP）、填充 IP 头、`netfilter` 过滤、对超过 MTU 大小的数据包进行分片。处理完这些工作后会交给网络接口层处理。

网络接口层会通过 ARP 协议获得下一跳的 MAC 地址，然后对 `sk_buff` 填充帧头和帧尾，接着将 `sk_buff` 放到网卡的发送队列中。

这一些工作准备好后，会触发「软中断」告诉网卡驱动程序，这里有新的网络包需要发送，驱动程序会从发送队列中读取 `sk_buff`，将这个 `sk_buff` 挂到 `RingBuffer` 中，接着将 `sk_buff` 数据映射到网卡可访问的内存 DMA 区域，最后触发真实的发送。

当数据发送完成以后，其实工作并没有结束，因为内存还没有清理。当发送完成的时候，网卡设备会触发一个硬中断来释放内存，主要是释放 `sk_buff` 内存和清理 `RingBuffer` 内存。

最后，当收到这个 TCP 报文的 ACK 应答时，传输层就会释放原始的 `sk_buff`。

发送网络数据的时候，涉及几次内存拷贝操作？

第一次，调用发送数据的系统调用的时候，内核会申请一个内核态的 `sk_buff` 内存，将用户待发送的数据拷贝到 `sk_buff` 内存，并将其加入到发送缓冲区。

第二次，在使用 TCP 传输协议的情况下，从传输层进入网络层的时候，每一个 `sk_buff` 都会被克隆一个新的副本出来。副本 `sk_buff` 会被送往网络层，等它发送完的时候就会释放掉，然后原始的 `sk_buff` 还保留在传输层，目的是为了实现在 TCP 的可靠传输，等收到这个数据包的 ACK 时，才会释放原始的 `sk_buff`。

第三次，当 IP 层发现 `sk_buff` 大于 MTU 时才需要进行。会再申请额外的 `sk_buff`，并将原来的 `sk_buff` 拷贝为多个小的 `sk_buff`。

## 总结

电脑与电脑之间通常都是通过网卡、交换机、路由器等网络设备连接到一起，那由于网络设备的异构性，国际标准化组织定义了一个七层的 OSI 网络模型，但是这个模型由于比较复杂，实际应用中并没有采用，而是采用了更为简化的 TCP/IP 模型，Linux 网络协议栈就是按照了该模型来实现的。

TCP/IP 模型主要分为应用层、传输层、网络层、网络接口层四层，每一层负责的职责都不同，这也是 Linux 网络协议栈主要构成部分。

当应用程序通过 `Socket` 接口发送数据包，数据包会被网络协议栈从上到下进行逐层处理后，才会被送到网卡队列中，随后由网卡将网络包发送出去。

而在接收网络包时，同样也要先经过网络协议栈从下到上的逐层处理，最后才会被送到应用程序。

---

参考资料：

- Linux 网络包发送过程：<https://mp.weixin.qq.com/s/wThfD9th9e-YGHJJ3HXNQ>
- Linux 网络数据接收流程（TCP）- NAPI：<https://wenfh2020.com/2021/12/29/kernel-tcp-receive/>

- Linux 网络 - 数据包接收过程:  
[https://blog.csdn.net/frank\\_jb/article/details/115841622](https://blog.csdn.net/frank_jb/article/details/115841622)

哈喽，我是小林，就爱图解计算机基础，如果觉得文章对你有帮助，别忘记关注我哦！



扫一扫，关注「小林coding」公众号

图解计算机基础  
认准**小林coding**

每一张图都包含小林的认真  
只为帮助大家能更好的理解

① 关注公众号回复「**图解**」  
获取图解系列 PDF

② 关注公众号回复「**加群**」  
拉你进百人技术交流群