

上海大学

SHANGHAI UNIVERSITY

毕业设计（论文）

UNDERGRADUATE PROJECT (THESIS)

题目：基于 socket 的文件传输的设计与实现

学院 计算机工程与科学

专业 计算机科学与技术

学号 10122060

学生姓名 吕伟彬

指导老师 雷咏梅

起讫日期 2014.02.17—2014.05.30

【目录】

基于 socket 的文件传输设计与实现

摘要

本文从先讨论了 socket 网络编程的方法, 然后讨论了文件传输协议的设计与实现, 主要包括 HTTP、FTP 和 BitTorrent 协议。涉及到以上几种文件传输协议的协议规范, 以及开源实现的实现细节, 包括 libcurl (HTTP 和 FTP) 以及 libtransmission (BitTorrent)。同时描述了我本人实现的文件传输管理器 wdl 以及一个简单的 HTTP 服务器 wdlS 的设计与实现。wdl 支持 FTP、HTTP 和 BitTorrent 协议, 多任务多线程下载, 支持断点续传; 可以从互联网上下载文件。wdls 是一个非常简单的 HTTP 服务器, 实现 HTTP200, HTTP404 和 HTTP206, 支持断点续传, 可以与 wdl 很好地配合。wdl 和 wdlS 都是在 Linux 完全采用 C 语言实现, wdlS 不依赖任何第三方库。

关键词: Socket,linux,文件传输协议, HTTP 服务器

ABSTRACT

This Document describes the design and implement of file transfer protocol using Socket, including HTTP, FTP and BitTorrent Protocol. The specifications and open source implements of those protocols are involved. Libcurl is introduced as HTTP and FTP implement, while libtransmission as BitTorrent implement. Also, my file transfer manager named wdl and simple HTTP server named wdlS are described in the document. Wdl supports FTP, HTTP and BitTorrent Protocol, multi-task, multi-thread and breakpoint-resume. It's OK to download files from the Internet using wdl. And wdlS is a very simple HTTP server, which only implements HTTP200, HTTP206 and HTTP404, and breakpoint-resume too. Wdl and wdlS work well together. Wdl and wdlS are both written in C language in Linux System. In addition, wdlS does not depend on any third-party library.

Keywords:Socket,linux,file transfer protocol,HTTP server

绪论

本课题从科研与实际应用出发,从 Socket 网络编程出发,旨在提高学生对网络数据传输的理解与认识。并能在实际生产中加以运用。

Socket 最早作为 BSD UNIX 的进程通信机制,随着网络的出现,它也开始支持网络通信;事实上,在 UNIX 系统下,网络通信被认为是一种特殊的进程间通信。网络的快速发展,使得 Socket 网络编程越来越流行,以至于现在很多人都忽略了 Socket 作为本机进程间通信的功能。本文中所谈到的 Socket,除非特别说明,都是只指网络套接字。而且 Windows 最早的 IP/TCP 实现就来源于 BSD UNIX 的代码,所以 Windows 下的网络编程接口也是 Socket 兼容的。

Socket 实际上就是操作系统内核(主要是 IP/TCP 协议栈)对外提供的网络编程接口,它属于系统调用层面,而不是特定语言的标准库。因为大部分操作系统都是 C 语言实现的,因此,Socket 接口一般都是以 C 函数的形式提供。而像 C++, Java, Python 这些语言的网路支持,要么是对 C 语言的封装,如 C++;要么就是通过虚拟机(如 Java)或者解释器(如 Python)间接调用 Socket 接口;总之,它们最底层还是要调用操作系统提供的 Socket 接口。因此,了解熟悉 Socket 接口对理解网络编程,无论哪个层面的,都非常有益。

既然 Socket 是作为操作系统网络协议栈对外的接口,那么了解相应的网络协议(主要是 TCP/IP)就很有必要了。本文的前一部分会简要描述 TCP/IP 协议的规范,但不会很深入,对这方面感兴趣的读者可以参考 IP/TCP 详解^[1]。接着,我们讨论 Socket 网络编程接口,从系统调用的层面展示网络编程的方法,这要求读者对 UNIX 系统下的 C 语言编程有一定基础,但这些并不在本文的讨论范围内;对这方面感兴趣的读者可以参考 UNIX 环境高级编程^[2]。

Socket 套接字其实不仅仅局限于网络编程,它有多种协议类型,最典型的就 AF_INET 和 AF_INET6,分别对应的是 IPv4 和 IPv6 套接字。AF_UNIX 或者 AF_LOCAL 是 UNIX 域套接字,用于本地进程间通信。还有其他协议特定的类型,比如 AF_IPX 是 IPX 协议套接字,AF_X25 是 ITU-T X.25 / ISO-8208 协议套接字,以及内核的 netlink 套接字 AF_NETLINK。本文只讨论 AF_INET,IPv6 套接字与 IPv4 类似,主要是地址结构不同。

而 AF_INET 协议下又可以指定 SOCK_STREAM 和 SOCK_DGRAM,分别对应 TCP 和 UDP 协议;甚至可以指定 SOCK_RAW 使用原始套接字,原始套接字可以接受到 IP 的数据报文,由程序员手动构造和解析 TCP 或者 UDP 协议首部字段。本文只讨论 SOCK_STREAM。也就是说只讨论 TCP 套接字,同时也只描述 TCP 和 IP 协议规范。

本文还讨论基于 Socket 的网络文件传输协议,FTP、HTTP 和 BitTorrent 协议。FTP 协议是流行的文件传输协议,主要用在托管大量文件的服务器上;HTTP 最初为了互联网浏览网页而设计,虽然协议中有很多与网页服务相关的字段,但本质就是文件传输的协议;BitTorrent 近年来也越来越流行,主要用在大文件的传输上。

最后结合上述的内容来描述我所完成的文件下载器 wdl 和一个简单的 HTTP 服务器 wdl_s,主要从网络编程的角度描述,而不是从软件工程的角度。

Socket 网络编程是比较底层的一种网络编程模型,国外,主要是美国有很好的研究与实现,主流的网络协议 TCP/IP、FTP、HTTP 等等都来源于国外;而国内对计算机网络底层的研究比较少,也没有什么成熟的产品,要改变现状,首先要从根本上提高计算机行业从业人员的专业素质,不能仅仅满足于开发一些上层应用,而对底层的具体实现知之甚少。

第一章 TCP/IP

引言

TCP/IP 起源于 60 年代末美国政府资助的一个分组交换网络研究项目，到 90 年代已经发展成为计算机之间最常用的联网方式。因为其协议规范以及很多实现都不需要花钱或者只要花很少的钱就可以获得，它被认为是一个开放的系统，是全球互联网的基础。Socket 本质上就是操作系统对外提供的操作 TCP/IP 协议的编程接口，了解 TCP/IP 协议对使用 Socket 进行网络编程有重要意义。本章简要描述了 TCP/IP 协议的规范。

我们有时说 TCP/IP 协议族，起始包括了很多其他协议，比如 ICMP、IGMP、UDP 等，只是 IP 和 TCP 是这些协议中用得最多的。本章也只讨论 IP 和 TCP 协议。

第一节 分层

网络协议通常是分层次的，OSI 定义了七层，这里为了简化讨论，我们采用 TCP/IP 的分层模型。

应用层
传输层
网络层
链路层

图 1-1 TCP/IP 协议族的层次

每一层都有不同的功能。最底层的是链路层，也可以称为数据链路层和网络接口层，一般包括了操作系统中的设备驱动程序，主要是网络接口。它们处理与电缆（或者任何其他传输媒介）的物理接口的细节。然后是网络层，这层处理数据分组在网络中的活动，比如路由选路等，这一层主要包括的网络协议有 IP 协议，ICMP 协议和 IGMP 协议。网络层上面是传输层，这层为数据提供端到端的通信，主要是 TCP 和 UDP 协议，也就是说 TCP 协议是基于 IP 协议的。最上面则是应用层，这层不处理任何网络数据传输的细节，而负责应用程序特定的数据细节，这一层的协议有 HTTP、FTP、Telnet 等。

Socket 网络编程主要是完成传输层的连接和数据传输，在此基础上可以实现 FTP、HTTP 等协议，甚至是用户自定义的应用层协议。下面我们具体讨论 IP 协议和 TCP 协议。

第二节 IP：网际协议

IP 是 TCP/IP 协议族中最核心的协议。所有的 TCP、UDP、ICMP 以及 IGMP 数据都是以 IP 数据报格式传输。IP 提供的是不可靠、无连接的数据报传输，所谓的不可靠是指 IP 协议不保证数据真正到达目的地，如果当中的某个地方出错了，比如路由器缓存区用完，那么 IP 数据报会被直接丢弃。无连接则指 IP 数据报之间没有任何关系，它们是互相独立的，IP 协议不维护任何有关数据报状态的信息。可靠性和可连接性都由上层协议，比如 TCP 来完成。

第一条 IP 首部

每个 IP 数据报都有一个 IP 首部，一般其长度是 20 字节，除非包含有选项字段。下图表示了 IP 首部结构，最高位在左边，记作 0bit；最低位在右边，为 31bit。

版本号 (4 位)	首部长度的 (4 位)	服务类型 TOS (8 位)	字节总长度 (16 位)	
标识 (16 位)			标志 (3 位)	偏移量 (13 位)
生存时间 TTL (8 位)		协议号 (8 位)	首部校验和 (16 位)	
源 IP 地址 (32 位)				
目的 IP 地址 (32 位)				
选项 (如果有)				
负载数据 (...)				

图 1-2 IP 首部结构

版本号是 4，也就是 IPv4。首部长度是 4 位二进制数，最大 15，表示的是首部占 32bit 字的数目，因此 IP 首部最长为 60 字节。如果没有选项，IP 首部长度为 20 字节。服务类型（TOS）字段包括一个 3bit 的优先级子字段（现在已被忽略），4bit 的 TOS 字段和 1bit 的保留位。4bit 的 TOS 字段分别表示：最小时延、最大吞吐量、最高可靠性和最小费用。总长度字段表示了整个 IP 数据报的长度，以字节为单位，总长度减去 IP 首部长度就得到了实际负载数据的长度，IP 数据报总长度最大为 65535 字节。标识字段唯一标识每一份 IP 数据报，通常每发送一份 IP 数据报它的值就加 1，最后会回归到 0。TTL 表示 time-to-live，指一个 IP 数据报在网络上的最大生存时间，不过它其实指的不是真实的时间，而是路由的跳数，每经过一个路由，TTL 会减 1，某个路由器发现 IP 数据报的 TTL 为 0，则直接丢弃，然后发送 ICMP 报文通知源主机。这也是 traceroute 的工作方式。

为了加深理解，我们看一下 linux 内核中 IP 首部的结构；

```
struct iphdr {
    #if defined(__LITTLE_ENDIAN_BITFIELD)
        __u8    ihl:4,
               version:4;
    #elif defined (__BIG_ENDIAN_BITFIELD)
        __u8    version:4,
               ihl:4;
    #else
        #error  "Please fix <asm/byteorder.h>"
    #endif
        __u8    tos;
        __be16  tot_len;
        __be16  id;
```

```
__be16 frag_off;

__u8  ttl;

__u8  protocol;

__sum16 check;

__be32 saddr;

__be32 daddr;

/*The options start here. */

};
```

第二条 IP 地址

互联网上的每个接口都必须有一个唯一的地址，也就是 IP 地址。就像在前面 IP 首部中看到的一样，IP 地址长 32bit。IP 地址具有一定的结构，分为五类，A、B、C、D 和 E 类地址。都接口主机具有多个 IP 地址，每个接口都有一个对应的 IP 地址。IP 地址还可以分为三类，单播地址（目的地为单个主机，这也是最多的情况），广播地址（目的地是给定网络上的所有主机）以及多播地址（目的地为同一组内所有主机）。

在 TCP/IP 领域中，域名系统（DNS）是一个分布的数据库，它提供了 IP 地址和主机名之间的映射信息。大部分网络程序都允许指定 IP 地址或者域名。

第三节 TCP：传输控制协议

TCP 在 IP 协议之上，提供了一种面向连接的、可靠的字节流服务。面向连接意味着两个使用 TCP 的应用（一般是一个客户端和一个服务器）在交换数据之前必须先建立一条 TCP 连接。而可靠性 TCP 使用以下方式来实现：

- 1. 数据被分割成 TCP 认为的最合适的大小发送，因为 IP 数据报无法保证接受顺序和发送顺序一致，因此再在接收端将缓存 TCP 数据报直到一个完整 TCP 数据报可以组合。
- 2. 当 TCP 发送一个数据包后，会启动一个定时器，等待接受方确认收到这个数据包。如果超时则认为发送失败。同样的，如果收到来自对方的数据，将返回一个确认。
- 3. TCP 将保证它首部和数据的校验和，如果校验和验证失败，则认为该数据包无效，丢弃。
- 4. IP 数据报会发生重复，TCP 将丢弃重复的数据报
- 5. TCP 会进行流量控制，保证接受的数据不会超过缓存区大小。

第一条 TCP 首部

TCP 数据报被封装在 IP 数据报中，如下图所示，

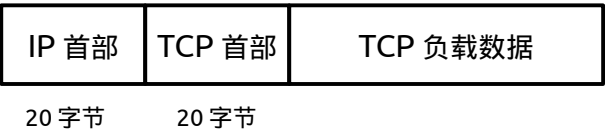


图 1-3 TCP 封装

TCP 首部和 IP 首部一样，如果没有任何选项，也是 20 个字节。这意味着，一个 TCP 数据报一般至少是 40 字节的。下图显示了一个 TCP 协议首部的结构。

源端口号 (16 位)								目的端口号 (16 位)							
序号 (32 位)															
确认序号 (32 位)															
首部长度 (4 位)		保留 (6 位)		U R G	A C K	P S H	R S T	S Y N	F I N	窗口大小 (16 位)					
检验和 (16 位)									紧急指针 (16 位)						
选项 (...)															
负载数据 (...)															

图 1-4 TCP 首部

回想上面的 IP 首部，IP 首部中是没有端口号的，但是有 IP 地址，而 TCP 首部中包含了端口号，用于寻找接受和发送数据的应用进程。源 IP 地址和源端口号确定了一个源主机，目的 IP 地址和目的端口号确定了目的主机，从而唯一确认了一个 TCP 连接。其实，最早的 TCP 规范中，将端口号称为插口 (socket)。序号用来标识从 TCP 数据报的发送字节流，表示在这个报文段中的第一个数据字节。用于接收端进行报文排序。确认序号用于接收端向发送端确认数据已成功接收，其值是上次接受到数据字节序号加 1，只有在 ACK 被设置时才有效。窗口大小用于接收端向发送端告知接收缓冲区的大小，如果接受端的缓冲区过小，发送端一般会暂停发送，直到接受端有了足够的缓冲区。

TCP 首部中有 6 个标志位，它们分别是

- 1. URG 紧急指针有效
- 2. ACK 确认序号有效
- 3. PSH 接受方应该尽快将这个报文段交给应用层
- 4. RST 重建连接
- 5. SYN 同步序号，用于发起连接
- 6. FIN 发送端已经发送结束

下面同样拿出 linux 内核中 TCP 首部的结构，以加深对 TCP 首部的理解。

```
struct tcphdr {
    __be16 source;
    __be16 dest;
    __be32 seq;
    __be32 ack_seq;
```



```
#if defined(__LITTLE_ENDIAN_BITFIELD)
    __u16  res1:4, doff:4,
           fin:1, syn:1, rst:1, psh:1,
           ack:1, urg:1, ece:1, cwr:1;
#elif defined(__BIG_ENDIAN_BITFIELD)
    __u16  doff:4, res1:4,
           cwr:1, ece:1, urg:1, ack:1,
           psh:1, rst:1, syn:1, fin:1;
#else
#error      "Adjust your <asm/byteorder.h> defines"
#endif

    __be16 window;
    __sum16 check;
    __be16 urg_ptr;
};
```

第二条 TCP 连接

TCP 的连接过程是著名的三次握手。首先客户端设置一个初始序号，发送一个 SYN 请求；服务器返回该 SYN 的 ACK 确认的同时也发送一个他自己的 SYN 请求，客户端确认服务器的 SYN 请求后，连接建立。

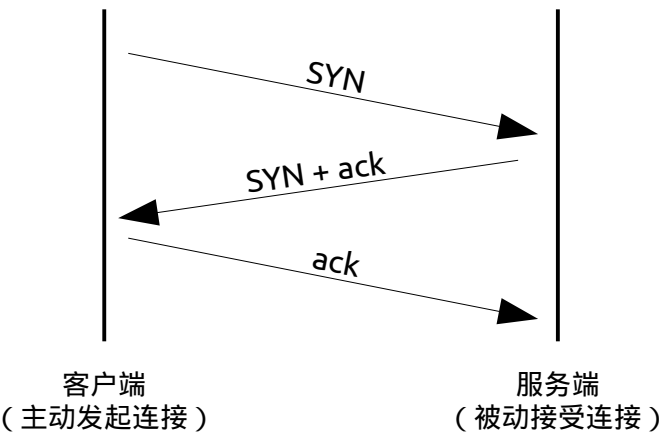


图 1-5 TCP 连接过程

第二章 Socket 网络编程

引言

本章简单地介绍 Socket 套接字 API，主要围绕 IPv4/TCP 描述；同时会用 wdlS 的实现作为具体实例。wdls 是完全采用 C 实现简单 HTTP 服务器，它不依赖任何第三方库，直接使用系统调用和标准 C 库。

第一节 套接字地址结构

大多数套接字函数需要一个套接字地址结构作为参数。不同协议使用不同的地址结构，比如 IPv4 的地址结构为 `sockaddr_in`，IPv6 的地址结构为 `sockaddr_in6`。`sockaddr_in` 结构定义在 `netinet/in.h` 中：

```
/* Structure describing an Internet (IP) socket address. */
#define __SOCK_SIZE__ 16          /* sizeof(struct sockaddr) */
struct sockaddr_in {
    __kernel_sa_family_t  sin_family; /* Address family */
    __be16                 sin_port;   /* Port number */
    struct in_addr          sin_addr;   /* Internet address */

    /* Pad to size of `struct sockaddr'. */
    unsigned char          __pad[__SOCK_SIZE__ - sizeof(short int) -
                                   sizeof(unsigned short int) - sizeof(struct in_addr)];
};
#define sin_zero    __pad          /* for BSD UNIX comp. -FvK */
```

`sin_family` 表示地址协议，可以是 `AF_INET` 或者 `AF_INET6`。`sin_port` 表示端口号，是以网络字节序表示的 16 位无符号数。`sin_addr` 表示真正的 IP 地址，这里当然是二进制表示的，而不是我们熟知的点分十进制。最后的 `__pad` 是填充字段，为了保证 `sockaddr_in` 的长度与 `sockaddr` 结构的长度一致。

`in_addr` 是一个只有一个字段 `s_addr` 的结构，一般引用 `sin_addr` 字段时可以直接 `sin_addr.s_addr`。

`sockaddr` 是一个通用的地址结构，为了统一套接字函数的地址参数而出现。

```
struct sockaddr {
    sa_family_t sa_family; /* address family, AF_xxx */
    char        sa_data[14]; /* 14 bytes of protocol address */
};
```

下面我们看一下 wdlS 是如何处理地址结构的。wdls 作为一个服务器，它需要把套接字和与本机地址绑定，因此要设置自己的 IP 地址结构。下面是 wdlS 处理本机地址结构的实现代码。

```
struct sockaddr_in addr;
memset(&addr, 0, sizeof(addr));

addr.sin_family = AF_INET;

addr.sin_addr.s_addr = htonl(INADDR_ANY);

addr.sin_port = htons(port);
```

先初始化 `sockaddr_in` 结构，所有字段设置为 0。然后设置地址类型为 `AF_INET`，表示 IPv4 地址。这里把 `sin_addr` 设置为 `INADDR_ANY` 表示让操作系统选择一个本机地址，因为我的测试机器是只有一个网络接口，也就是一个 IP 地址，所以这么处理是没有问题的，但是对于那些有多个 IP 地址的多接口主机，让系统选择一个 IP 地址不是一个好办法。最后设置 `sin_port` 端口号。这里的 `htonl` 和 `htons` 都是将本机字节序转化为网络字节序。其含义是 `host to network long` 和 `host to network short`，网络数据采用大端字节序，如果本机字节序也是大端，那么这两个函数不做任何处理，否则转化为大端字节序表示。关于更多字节序的内容，读者可以参考 [ON HOLY WARS AND A PLEA FOR PEACE^{\[3\]}](#)。

同时，系统提供了将点分十进制表示的 IPv4 地址转化为二进制形式的函数。

```
#include <arpa/inet.h>

int inet_aton(const char *strptr, struct in_addr *addrptr)

若转化成功则返回 1，否则返回 0
```

`strptr` 是点分十进制表示的 IP 地址字符串，而 `addrptr` 是用于返回二进制表示的地址结构的指针。如果要指定某个 IP 地址，可以这样设置

```
inet_aton("123.123.123.123",&(addr.sin_addr));
```

还有几个类似的地址转换函数 `inet_addr` 和 `inet_ntoa`，不再赘述。

第二节 TCP 套接字编程

本节描述编写一个完整的 TCP 客户端/服务器程序所需要的基本套接字函数。同样会引入 `wdls` 中的具体实例来讲解。

第一条 Socket 函数

为了执行网络通信，进程要做的第一件事情就是调用 `socket` 函数来创建一个指定的套接字描述符。套接字的类型可以是 IPv4 协议、IPv6 协议等，我们只描述 IPv4 协议套接字的 TCP 协议。要创建一个 TCP 套接字，要像如下调用 `socket` 函数：

```
int sockfd=socket(AF_INET,SOCK_STREAM,0); 或者
int sockfd=socket(AF_INET,SOCK_STREAM,IPPROTO_TCP);
```

第三个参数如果指定为 0，则系统会自动选择一种类型为 `SOCK_STREAM` 的套接字，一般就是 TCP 套接字，当然也可以手动指定 `IPPROTO_TCP` 创建 TCP 套接字。第一种，就是 `wdls` 中创建套接字的方式。

在 UNIX 系统下，网络套接字被抽象成一个特殊的文件，因此，套接字描述符和文件描述符也是统一管理的，也就是说，文件描述符和套接字描述符是不会重复的。如果已经打开了文件 0、1、2、3、4，再创建一个套接字描述符，则为 5。反之亦然。

第二条 connect 函数

connect 函数用于 TCP 客户端建立与 TCP 服务器的连接。它有三个参数，第一个是 socket 函数返回的有效套接字描述符，第二个和第三个参数分别是 TCP 服务器的地址结构和地址结构的长度。

connect 是客户端调用向服务器发起连接的函数，因此服务端 wdlS 并不需要调用。下面代码展示了客户端如何调用 connect 函数向服务器发起 TCP 连接。connect 函数成功返回 0，失败返回-1。

```
sockaddr_in addr;
addr.sin_family = AF_INET;
inet_aton("123.123.123.123",&(addr.sin_addr));    /* 设置地址 */
addr.sin_port = htons(80);                        /* 设置端口号 */
connect(sockfd, (struct sockaddr*)&addr, sizeof(addr));
```

调用 connect 函数后内核会向目标发起 TCP 连接，这里涉及到 TCP 连接过程中的三次握手，因此 connect 函数是会阻塞的。

第三条 bind 函数

bind 函数将一个本地协议地址赋予一个套接字描述符，对于 IPv4，是一个 32 位的地址和一个 16 位的端口号的组合。对于客户端来说，这个函数是可以不调用的，而且大部分客户端都不调用它；不调用 bind 不代表套接字描述符就没有绑定的地址，而是系统内部自动赋予一个本地的地址和端口号。对于客户端来说，由系统自动设置地址和端口号是没有问题的，但是服务器一般都要手动调用 bind 的。因为服务器一般有多个网络接口，更重要的是服务器需要把自己的地址和端口号告诉客户端，如果让系统自动设置，那么这个地址和端口号都是随机的。

回想套接字地址结构一节，我们设置了一个本地的套接字地址结构，然后 wdlS 就可以直接将其与套接字描述符绑定。

```
bind(sockfd, (struct sockaddr *) &addr, sizeof(addr));
```

第四条 listen 函数

用 socket 创建的套接字默认是主动的，所谓主动套接字就是可以执行 connect 主动发起连接的；但是服务器往往是被动等待客户端连接，这就需要调用 listen 函数了，listen 函数有两个作用，第一个就是将主动模式的套接字转化为被动模式，第二个是设置内核应该为套接字排队的最大连接个数。

在 wdlS 中如下调用 listen 函数。

```
listen(sockfd, BACKLOG);
```

其中 BACKLOG 是数字 10 的宏定义。

listen 函数虽然有“监听”的意思，但只是设置套接字的被动模式，并没有任何网络操作，因此也不会阻塞。下面所讲的 accept 函数才真正监听网络连接。

第五条 accept 函数

accept 函数由 TCP 服务器调用，用于从已完成连接队列中返回下一个连接，如果连接队列为空，则进入阻

塞。注意，一个套接字在调用 listen 后就可以接受来自客户端的连接，不过只有服务器调用了 accept 函数后，该连接才真正被服务器处理，在此期间，连接被放在连接完成队列里面。不过一般服务器调用 listen 后就会立马调用 accept，因此，这个时间差基本可以忽略。

accept 函数有三个参数，第一个是套接字描述符，第二个和第三个分别是地址结构的指针和一个长度结构的指针，它们是作为返回值存在。如果有连接已经完成，accept 返回这个连接的套接字描述符，同时在第二个和第三个参数分别返回客户端的地址结构和地址结构的长度。

在 wds 中用一个循环调用 accpet 函数，每次成功都返回一个套接字描述符，然后创建一个线程来处理这个连接，代码如下；

```
while ((accfd = accept(sockfd, (struct sockaddr *) &arg->addr,
                        &arg->addrlen)) > 0) {

    arg->sockfd = accfd;

    http_thread(arg);

    arg = http_thread_arg_new();

}
```

每次成功返回一个连接套接字后，就设置一个线程参数结构，然后 http_thread 中创建子线程来处理这个连接。关于 http_thread 会在后面详细描述其实现，现在我们只关注于套接字接口相关的函数。

第六条 I/O 操作

当 TCP 连接建立以后，就可以像读写文件一样调用 read 和 write 来读写网络连接。但是所表现的行为有些不同，对套接字的 read 和 write 可能实际输入或者输出的数据少于请求的，这不算是错误，原因在于内核中的输出缓冲区已满，或者还没有接受到来自网络的足够多的数据。

第七条 TCP 客户端实例

本章节展示了一个简单的 TCP 客户端，该程序从命令行读取一个 IP 地址，然后发起 HTTP 请求，读取响应，并打印在终端。为了简单起见，没有做任何错误处理，只为了体现 socket 网络编程的方法。

```
#include <stdio.h>

#include <stdlib.h>

#include <string.h>

#include <unistd.h>


#include <sys/socket.h>

#include <arpa/inet.h>


int main(int argc, char *argv[])
{
```

```
    if(argc!=2){
        return -1;
    }
    const char *ip=argv[1];
    struct sockaddr_in addr;
    addr.sin_family=AF_INET;
    inet_aton(ip,&(addr.sin_addr));
    addr.sin_port=htons(80);          /* HTTP 服务端口号 80 */

    int sockfd=socket(AF_INET,SOCK_STREAM,0);
    connect(sockfd,(struct sockaddr*)&addr,sizeof(addr));

    const char request="GET / HTTP/1.1\r\n\r\n";          /* 一个最简单的 HTTP 请求 */

    write(sockfd,request,strlen(request));          /* 将 HTTP 请求发送给服务器 */

    char buff[1024];
    int n=read(sockfd,buff,1024);          /* 只读取一次数据 */

    buff[n]='\0';
    printf("%s",buff);

    close(sockfd);          /* 像关闭文件一样关闭套接字 */
    return 0;
}
```

该程序的运行结果如下：（省略了部分内容）

```
wiky@thunder:test$ ./a.out 202.120.127.189
HTTP/1.1 400 Bad Request
Content-Type: text/html; charset=us-ascii
Server: Microsoft-HTTPAPI/2.0
Date: Tue, 20 May 2014 05:54:25 GMT
.....
```

第三章 应用层协议

引言

前两章简单描述了 IP/TCP 协议规范以及相应的 socket 套接字编程。本章讨论基于 IP/TCP 协议的应用层协议，包括 HTTP、FTP 和 BitTorrent，以及同样基于 socket 的开源实现。其中 HTTP 和 FTP 使用流行广泛的开源网络库 libcurl 来描述，而 BitTorrent 使用 libtransmission 来描述。

本章会引用部分源代码来描述，分别是 curl-7.3.2 和 transmission-2.8.2。在引用源代码时会标出所在文件以及所在行号。比如在讨论 FTP 时，lib/ftp.c:2256 表示 curl-7.3.2 源代码目录下 lib/ftp.c 文件中第 2256 行。

第一节 FTP：文件传输协议

FTP 是一个古老的文件传输协议，具体协议规范读者可以参考 RFC959^[4]。这里只做简单介绍，结合 libcurl 的实现。

尽管 HTTP 协议已经代替了 FTP 大多数功能，FTP 仍然是互联网上传输文件的一种有效途径。FTP 客户端可以向服务器发送命令来下载文件、上传文件、创建甚至改变服务器上的目录，只要有相应的权限。FTP 相对于 HTTP 的优势就在于他有权限管理，不过大部分 FTP 服务器都允许匿名访问。回想第一章第一节中描述的网络分层模型，FTP 处于应用层，在 TCP 协议之上，或者说是基于 TCP 协议的。

FTP 服务一般运行在 20 和 21 两个端口。端口 20 用于在客户端和服务器之间传输数据，称之为数据流；而端口 21 则用于传输控制命令，称之为控制流。在 FTP 中有几个术语需要明确，

- 服务端控制进程：服务器上运行的 FTP 控制流管理进程，用于和客户端交换控制信息。
- 服务端数据进程：服务器上运行的 FTP 数据流管理进程，由控制进程管理，与客户端交换数据。
- 客户端控制进程：客户端上运行的 FTP 控制流管理进程，用于和服务器交换控制信息。
- 客户端数据进程：客户端上运行的 FTP 数据流管理进程，由控制进程管理，与服务器交换数据。

下图描述了 FTP 的通信模型。

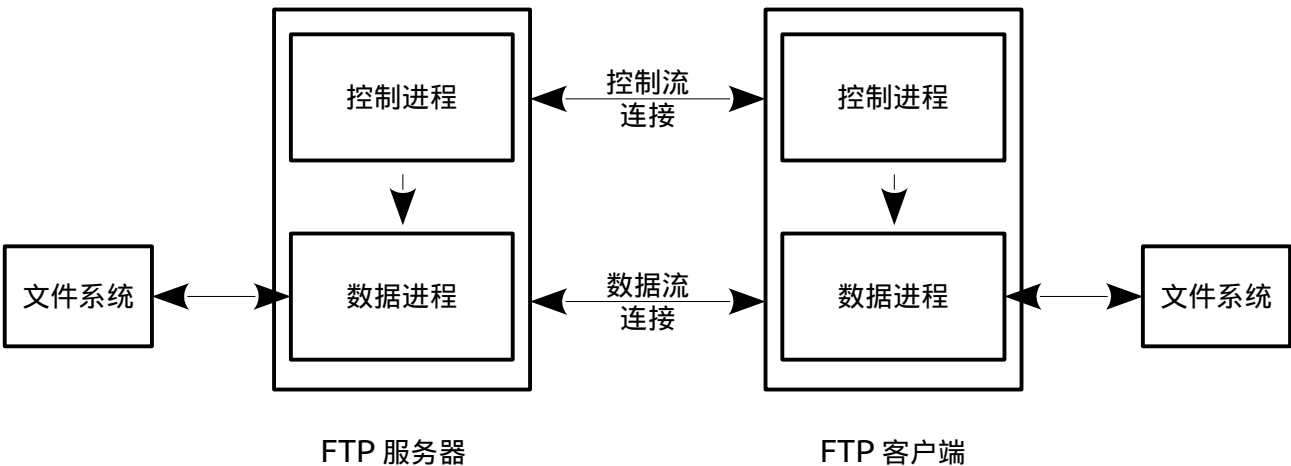


图 3-1 FTP 通信模型

按照 FTP 的设计, 控制进程和数据进程可以不在同一个主机上, 主要保证控制进程可以控制数据进程进行及时的数据传输。当然, 大部分情况下, 控制进程和数据进程都是在同一个系统中, 甚至可能是同一个进程。

FTP 的控制命令以 NVT ASCII 串的格式传输。每个命令以三个或者四个大写的 NVT ASCII 字符开始, 后面带有选项参数和一个 CR/LF 对来表示命令结束。应答由三个 NVT ASCII 数字以及一个选项消息组成。一个长应答也许会有多个消息组成, 第一个消息的三个数字后带有一个破折号, 最后的消息不带有破折号。中间的消息无须携带三个数字, 但是如果带有三个数字, 则也需要破折号。FTP 支持很多控制命令, 这里只列出其中几个,

- RETR: 从远程系统取回一个文件。
- PWD: 显示服务器端的当前工作目录。
- LIST: 显示当前工作目录下的文件列表。

FTP 有两种工作模式, 分别是主动模式和被动模式。在主动模式下, FTP 客户端随机开启一个大于 1024 的端口 N 向服务器的 21 号端口发起连接, 然后开放 N+1 号端口进行监听, 并向服务器发出 PORT N+1 命令。服务器接收到命令后, 会用其本地的 FTP 数据端口 (通常是 20) 来连接客户端指定的端口 N+1, 进行数据传输。

在被动模式下, FTP 客户端随机开启一个大于 1024 的端口 N 向服务器的 21 号端口发起连接, 同时会开启 N+1 号端口。然后向服务器发送 PASV 命令, 通知服务器自己处于被动模式。服务器收到命令后, 会开放一个大于 1024 的端口 P 进行监听, 然后用 PORT P 命令通知客户端, 自己的数据端口是 P。客户端收到命令后, 会通过 N+1 号端口连接服务器的端口 P, 然后在两个端口之间进行数据传输。

总的来说, 主动模式的 FTP 是指服务器主动连接客户端的数据端口, 被动模式的 FTP 是指服务器被动地等待客户端连接自己的数据端口。

被动模式的 FTP 通常用在处于防火墙之后的 FTP 客户访问外界 FTP 服务器的情况, 因为在这种情况下, 防火墙通常配置为不允许外界访问防火墙之后主机, 而只允许由防火墙之后的主机发起的连接请求通过。因此, 在这种情况下不能使用主动模式的 FTP 传输, 而被动模式的 FTP 可以良好的工作。

第二节 libcurl 的 FTP 实现

libcurl 没有实现 FTP 的所有命令, 事实上, 因为 FTP 的一些命令用得很少, 很多 FTP 实现都选择性地忽略了它们。但至少实现了 RETR 命令。RETR 后面跟一个文件名, 向服务器请求相应的文件。我们就用这个命令作为例子讲解 libcurl 是如何实现 RETR 命令的。

在开始 FTP 数据传输之前, 需要先进行 FTP 连接, 而 FTP 连接的基础是 TCP 的连接。

FTP 连接在函数 ftp_connect 中完成。该函数首先将 ftp_statemach_act 函数指针赋值给 connectdata 结构中的 statemach_act 字段, 然后通过 ftp_multi_statemach 函数间接调用 ftp_statemach_act。

ftp_statemach_act 根据 FTP 连接的当前状态不同调用不同的处理函数, 默认情况下采用的是 FTP 被动模式。

```
Lib/ftp.c:3101
```

```
case FTP_PASV:
```

```
result = ftp_state_pasv_resp(conn, ftpcode);
```



```
break;
```

ftp_state_pasv_resp 函数先解析地址，将域名或者字符串形式的 IP 地址转化为二进制表示，然后调用 Curl_connecthost 发起服务器的连接。

Curl_connecthost 内部又调用了 singleipconnect 完成对服务器的连接。

singleipconnect 先调用 Curl_socket 创建一个套接字描述符，然后调用 connect 函数完成连接。

```
Lib/connect.c:955
```

```
res = Curl_socket(conn, ai, &addr, &sockfd);
```

```
Lib/connect.c:1022
```

```
rc = connect(sockfd, &addr.sa_addr, addr.addrlen);
```

```
if(-1 == rc)
```

```
    error = SOCKERRNO;
```

Curl_socket 只是 socket 函数的简单包裹，代码在 lib/connect.c:1272。

完成 TCP 连接后发送 FTP 的 USER 命令进行登录，该操作在函数 ftp_state_user 中完成。

```
Lib/ftp.c:835
```

```
static CURLcode ftp_state_user(struct connectdata *conn){
```

```
    CURLcode result;
```

```
    struct FTP *ftp = conn->data->state.proto.ftp;
```

```
    /* send USER */
```

```
    PPSENDF(&conn->proto.ftpc.pp, "USER %s", ftp->user?ftp->user:"");
```

USER 默认是 anonymous，发送完 USER 后，有些服务器会要求密码，有些则直接返回 2xx 表示登录成功了。如果要求密码，则发送密码，否则登录成功。这些操作在函数 ftp_state_user_resp 中完成。

```
Lib/ftp.c:2626
```

```
if((ftpcode == 331) && (ftpc->state == FTP_USER)) {
```

```
    /* 331 Password required for ...(the server requires to send the user's password too) */
```

```
    PPSENDF(&ftpc->pp, "PASS %s", ftp->passwd?ftp->passwd:"");
```

```
    state(conn, FTP_PASS);
```

```
}else if(ftpcode/100 == 2) {
```

```
    /* 230 User ... Logged in.(the user logged in with or without password) */
```

```
    result = ftp_state_loggedin(conn);
```

```
}
```

```
.....
```

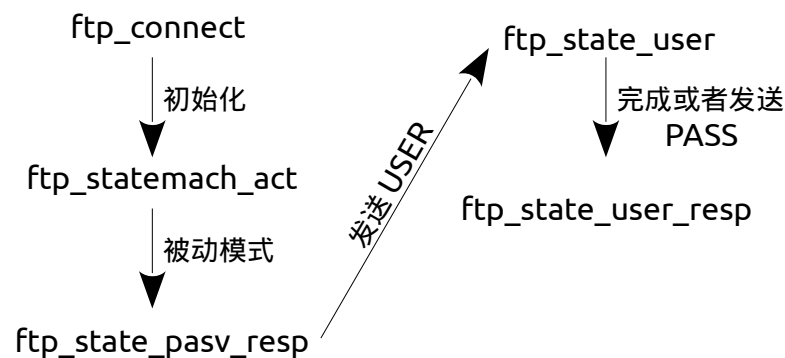


图 3-2 FTP 的连接过程

完成登录后，就可以使用 RETR 命令下载文件了，在函数 ftp_state_retr 中完成。它先判断是否指定了下载位置（断点续传），如果指定了，那么使用 REST 命令从断点开始下载，否则使用 RETR 命令。

```
Lib/ftp.c:2331
PPSENDF(&ftpc->pp, "RETR %s", ftpc->file);
```

PPSENDF 是一个宏定义，展开以后其实最后调用了 Curl_write 发送数据。Curl_write 调用了协议特定的输出函数，其实值是 write 函数的简单封装。

最后我们看看 libcurl 是如何接收来自服务器的响应的。读取响应数据的函数是 ftp_readresp，该函数在 lib/ftp.c:611 实现，内部其实调用了 Curl_pp_readresp，其内部又调用了 Curl_read。Curl_read 是一个支持普通套接字，SSL 套接字和 kerberos 套接字的通用读取函数。在 lib/ftp.c:522 实现。

第三节 HTTP：超文本传输协议

HTTP 是

参考文献

- [1] W.Richard Stevens. TCP/IP Illustrated Volume 1: The Protocol. Addison Wesley, 1994.
 - [2] W.Richard Stevens, Stephen A Rago. 尤晋元, 张亚英, 戚正伟 译。 Advanced Programming in the UNIX Environment. UNIX 环境高级编程。人民邮电出版社, 2006 年 5 月 1 日。
 - [3] Danny Cohen. ON HOLY WARS AND A PLEA FOR PEACE.
<http://www.ietf.org/rfc/ien/ien137.txt>. 1 April 1980.
 - [4] J. Postel, J. Reynolds. FILE TRANSFER PROTOCOL (FTP).
<http://www.ietf.org/rfc/rfc959.txt>. October 1985.
-