

distributed 1 / sockets

# last time

## RAID

ordering writes carefully

waste space rather than point to unallocated/wrong

fsck (filesystem check) recovery

redo logging ('writeahead logging')

write intention to log; redo committed parts of logs on reboot

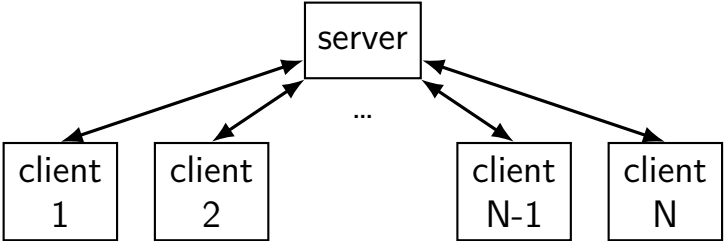
snapshots via copy-on-write

copy only parts that change  
indirection for inode array

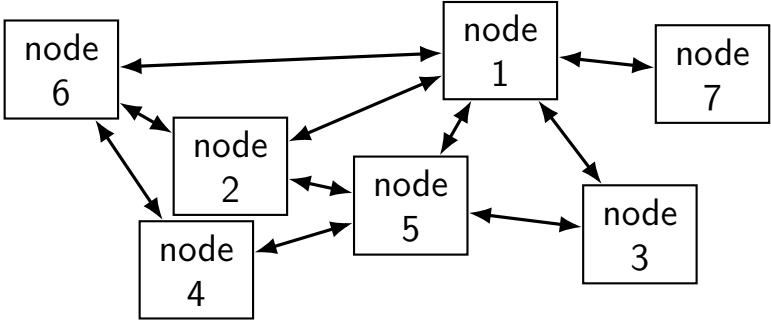
# distributed systems

multiple machines working together to perform a single task  
called a *distributed system*

# some distributed systems models

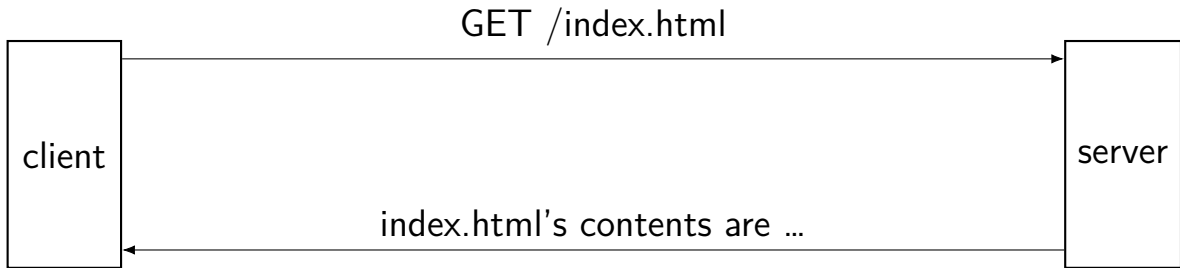


client/server

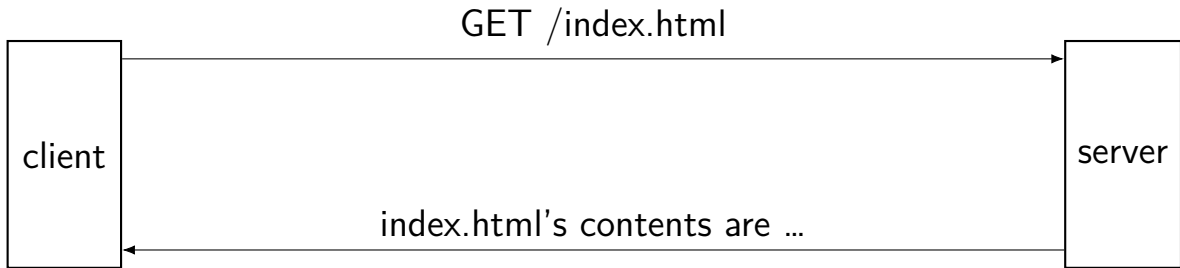


peer-to-peer

# client/server model

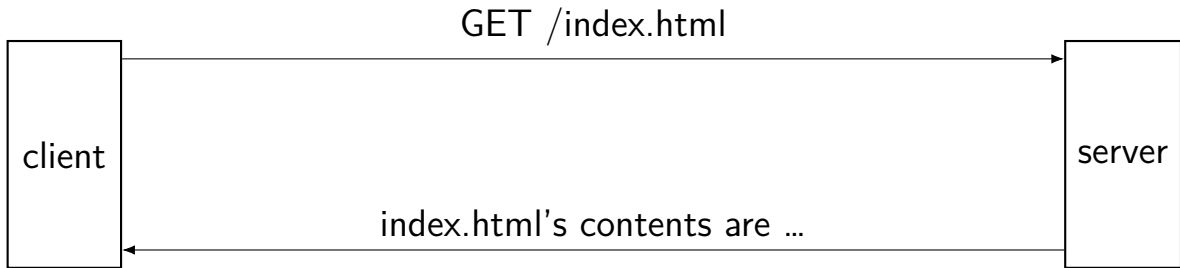


# client/server model



client(s): "sometimes on"  
sends requests to server(s)  
needs to know  
how to contact server

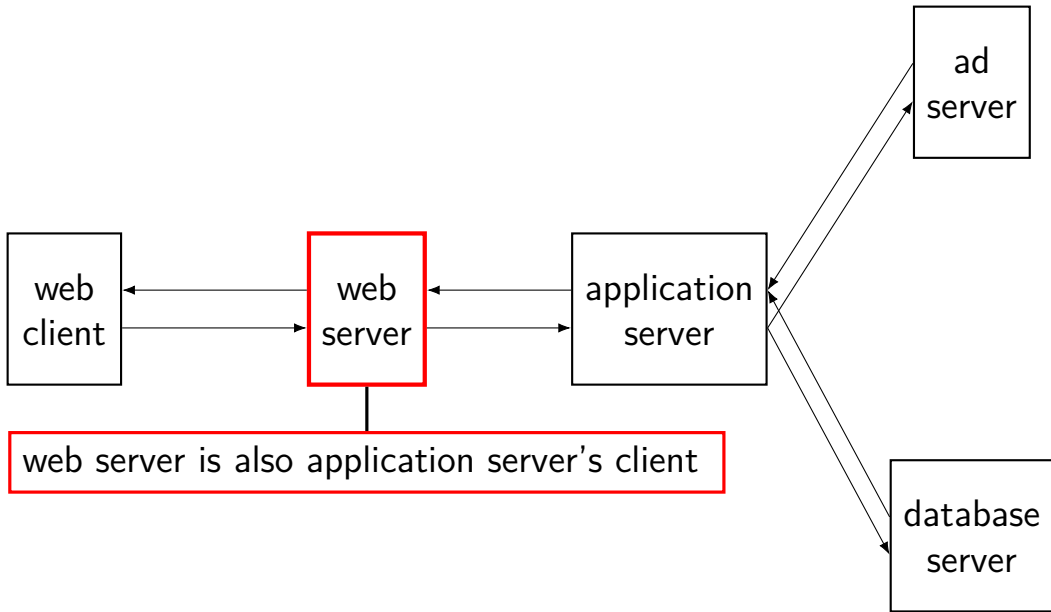
# client/server model



client(s): "sometimes on"  
sends requests to server(s)  
needs to know  
how to contact server

server(s): "always on"  
responds to client requests  
never initiates contact  
with a client

# layers of servers?





# example: Wikipedia architecture

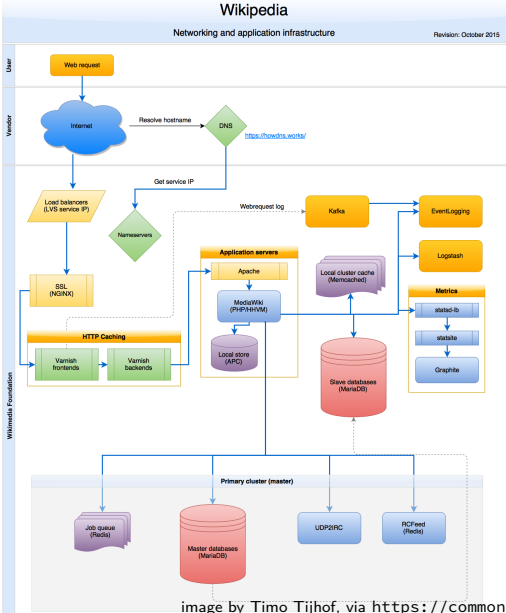
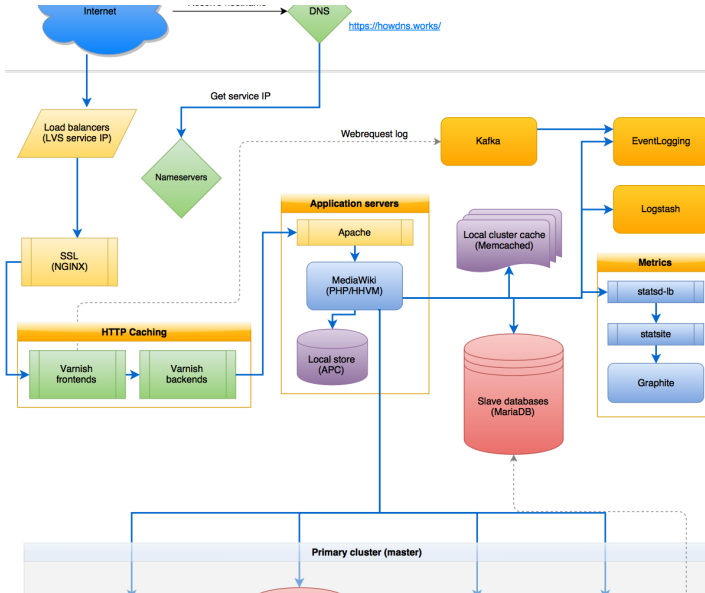


image by Timo Tijhof, via [https://commons.wikimedia.org/wiki/File:Wikipedia\\_webrequest\\_flow\\_2015-10.png](https://commons.wikimedia.org/wiki/File:Wikipedia_webrequest_flow_2015-10.png)

# example: Wikipedia architecture (zoom)



# peer-to-peer

no always-on server everyone knows about

hopefully, no one bottleneck — “scalability”

any machine can contact any other machine

every machine plays an approx. equal role?

set of machines may change over time

# why distributed?

multiple machine owners **collaborating**

**delegation** of responsibility to other entity  
put (part of) service “in the cloud”

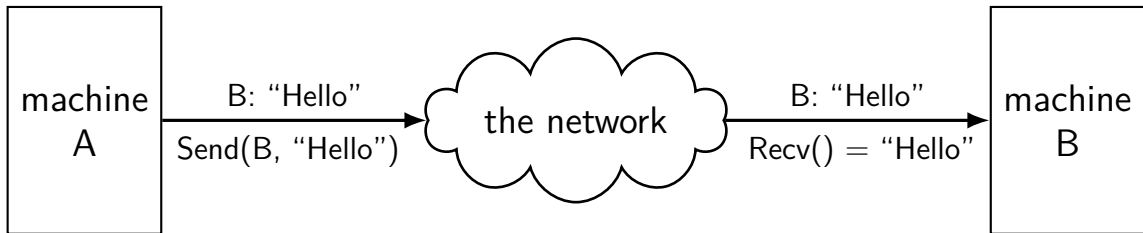
**combine** many **cheap machines** to replace expensive machine

easier to **add incrementally**

**redundancy** — one machine can fail and *system* still works?

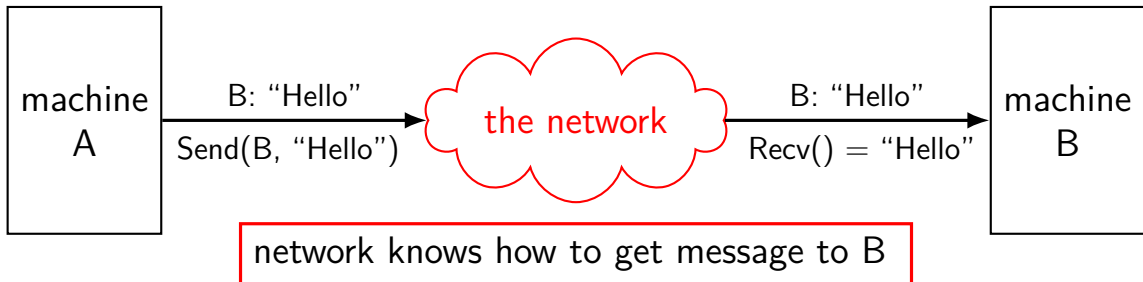
# mailbox model

*mailbox* abstraction: send/receive messages



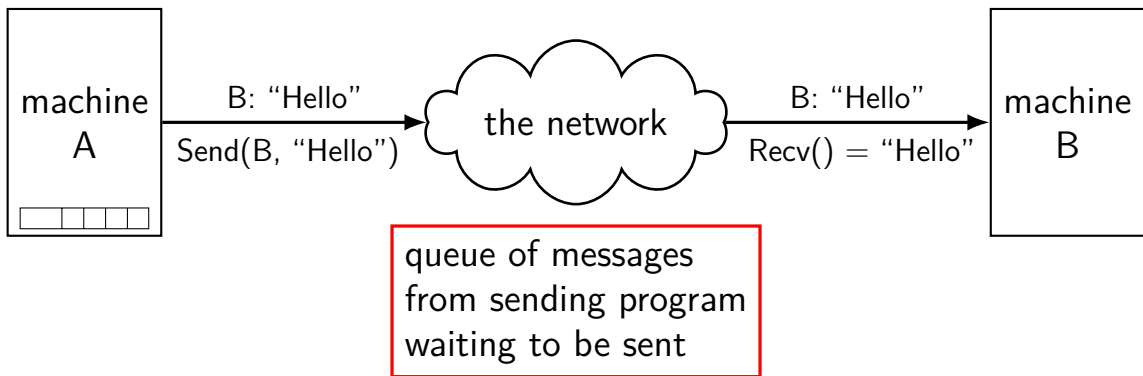
# mailbox model

*mailbox* abstraction: send/receive messages



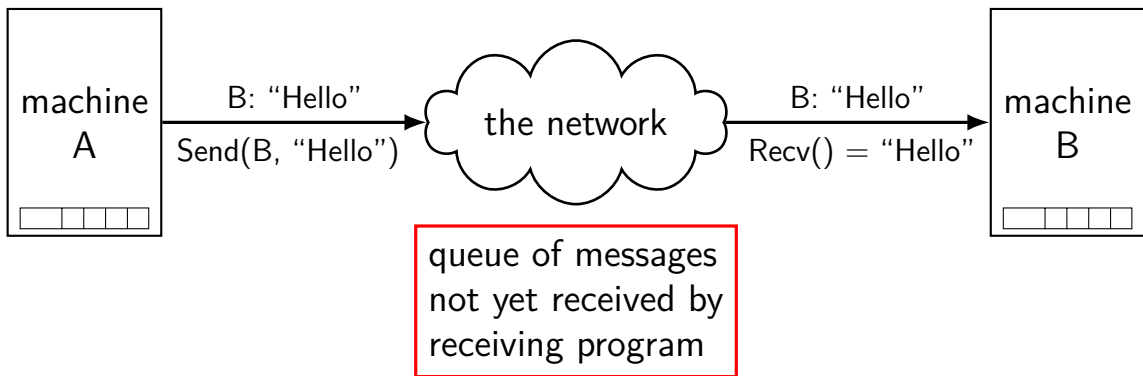
# mailbox model

*mailbox* abstraction: send/receive messages



# mailbox model

*mailbox* abstraction: send/receive messages





# what about servers?

client/server model: server wants to reply to clients

might want to send/receive multiple messages

# what about servers?

client/server model: server wants to reply to clients

might want to send/receive multiple messages

can build this with mailbox idea

- send a 'return address'

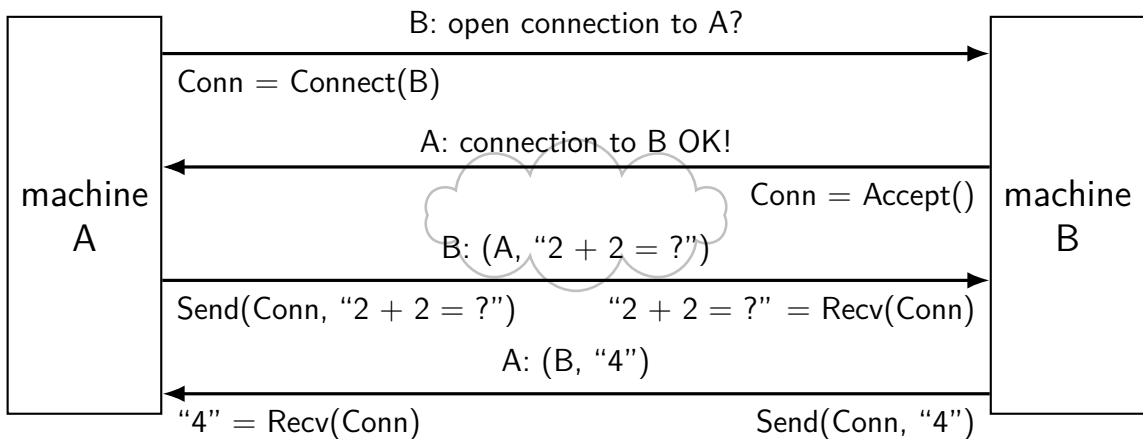
- need to track related messages

common abstraction that does this: the connection

## extension: connections

*connections*: two-way channel for messages

extra operations: connect, accept



# connections over mailboxes

real Internet: mailbox-style communication

connections implemented on top of this

including handling errors, transmitting more data than fits in message, ...

full details: take networking (CS/ECE 4457)

## connections versus pipes

connections look kinda like two-direction pipes

in fact, in POSIX will have the same API:

each end gets file descriptor representing connection

can use `read()` and `write()`

# connection missing pieces?

how to specify the machine?

multiple programs on one machine? who gets the message?

# names and addresses

**name**

logical identifier

hostname `www.virginia.edu`

hostname `mail.google.com`

hostname `mail.google.com`

filename `/home/cr4bd/NOTES.txt`

variable `counter`

service name `https`

**address**

location/how to locate

IPv4 address `128.143.22.36`

IPv4 address `216.58.217.69`

IPv6 address `2607:f8b0:4004:80b::2005`

inode# `120800873`

and device `0x2eh/0x46d`

memory address `0x7FFF9430`

port number `443`

# hostnames

typically use *domain name system* (DNS) to find machine names

maps logical names like `www.virginia.edu`

- chosen for humans

- hierarchy of names

...to *addresses* the network can use to move messages

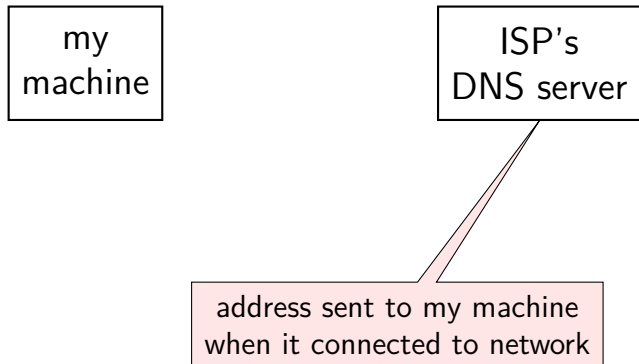
- numbers

- ranges of numbers assigned to different parts of the network

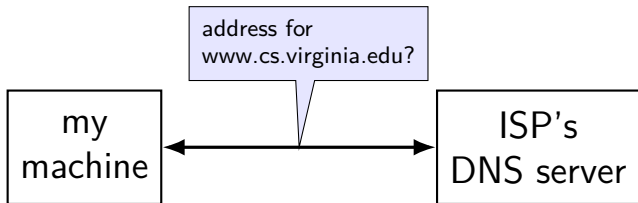
- network *routers* knows “send this range of numbers goes this way”



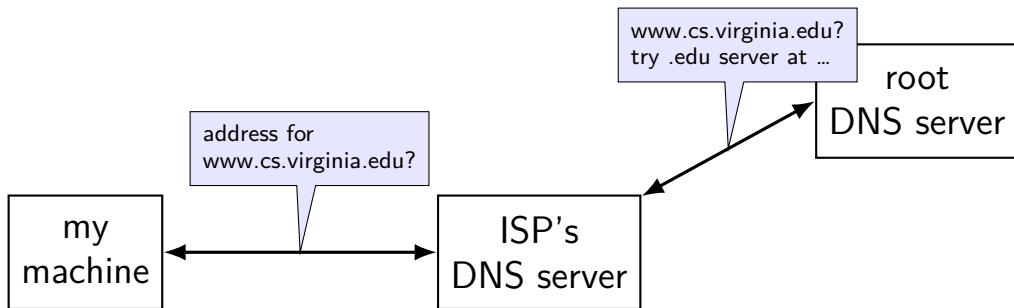
# DNS: distributed database



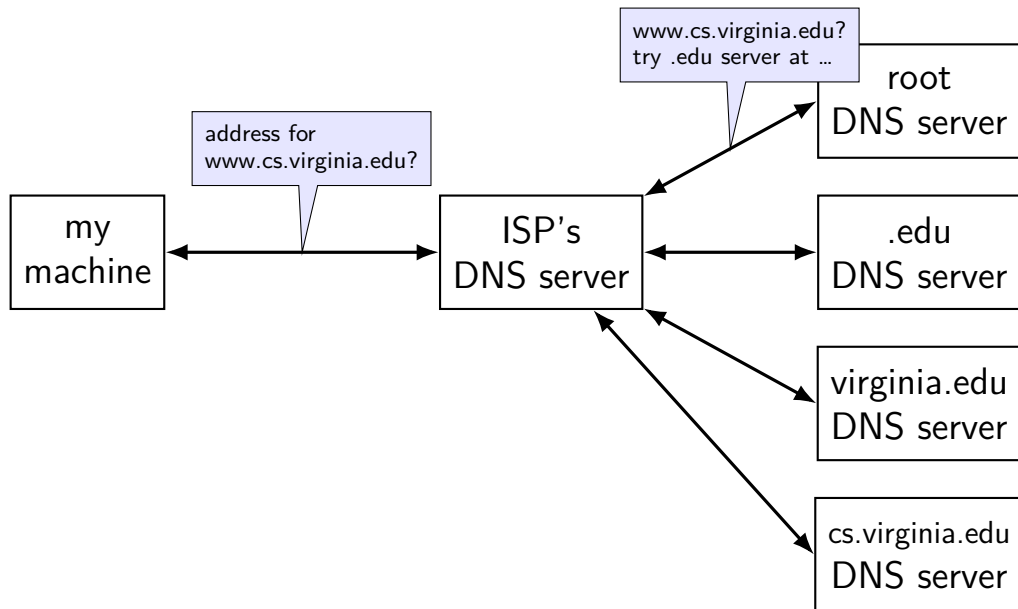
# DNS: distributed database



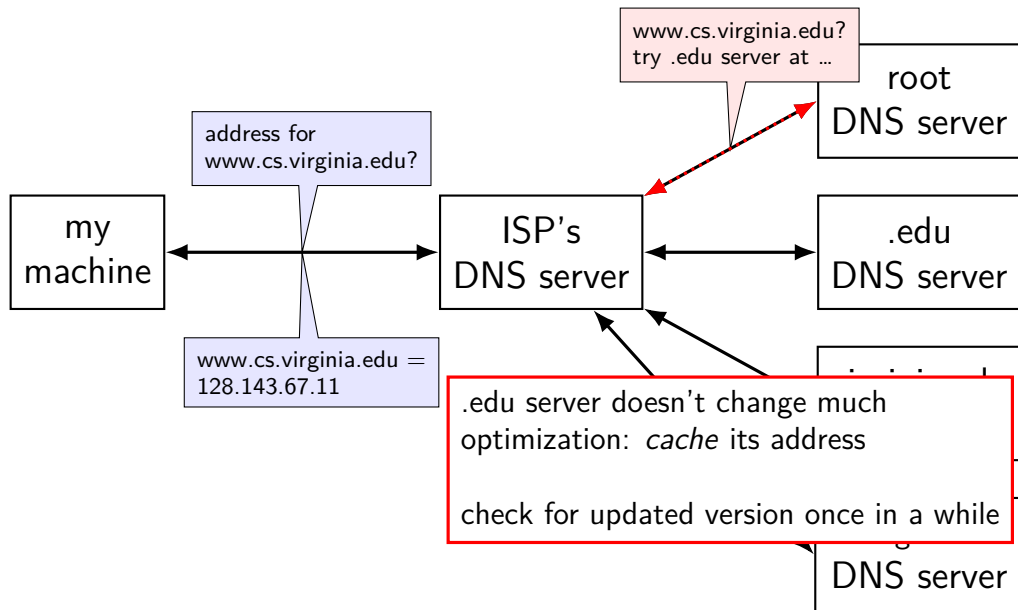
# DNS: distributed database



# DNS: distributed database



# DNS: distributed database



# IPv4 addresses

32-bit numbers

typically written like 128.143.67.11

four 8-bit decimal values separated by dots

first part is most significant

same as  $128 \cdot 256^3 + 143 \cdot 256^2 + 67 \cdot 256 + 11 = 2\,156\,782\,459$

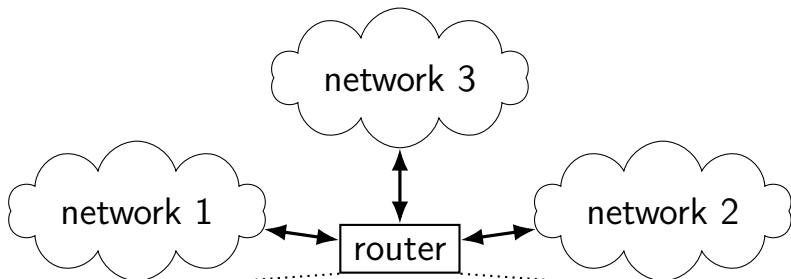
organizations get blocks of IPs

e.g. UVA has 128.143.0.0–128.143.255.255

e.g. Google has 216.58.192.0–216.58.223.255 and

74.125.0.0–74.125.255.255 and 35.192.0.0–35.207.255.255

# IPv4 addresses and routing tables



| if I receive data for...      | send it to... |
|-------------------------------|---------------|
| 128.143.0.0—128.143.255.255   | network 1     |
| 192.107.102.0—192.107.102.255 | network 1     |
| ...                           | ...           |
| 4.0.0.0—7.255.255.255         | network 2     |
| 64.8.0.0—64.15.255.255        | network 2     |
| ...                           | ...           |
| anything else                 | network 3     |

# selected special IPv4 addresses

127.0.0.0 — 127.255.255.255 — localhost

AKA loopback

the machine we're on

typically only 127.0.0.1 is used

192.168.0.0–192.168.255.255 and

10.0.0.0–10.255.255.255 and

172.16.0.0–172.31.255.255

“private” IP addresses

not used on the Internet

commonly connected to Internet with **network address translation**

also 100.64.0.0–100.127.255.255 (but with restrictions)

169.254.0.0–169.254.255.255

link-local addresses — ‘never’ forwarded by routers



# network address translation

IPv4 addresses are kinda scarce

solution: *convert* many private addrs. to one public addr.

locally: use private IP addresses for machines

outside: private IP addresses become a single public one

commonly how home networks work (and some ISPs)

# IPv6 addresses

IPv6 like IPv4, but with 128-bit numbers

written in hex, 16-bit parts, separated by colons ( : )

strings of 0s represented by double-colons ( :: )

typically given to users in blocks of  $2^{80}$  or  $2^{64}$  addresses  
no need for address translation?

2607:f8b0:400d:c00::6a =

2607:f8b0:400d:0c00:0000:0000:0000:006a

2607f8b0400d0c0000000000000000006a<sub>SIXTEEN</sub>

# selected special IPv6 addresses

`::1` = localhost

anything starting with `fe80` = link-local addresses  
never forwarded by routers

# port numbers

we run multiple programs on a machine

IP addresses identifying machine — not enough

# port numbers

we run multiple programs on a machine

IP addresses identifying machine — not enough

so, add 16-bit *port numbers*

think: multiple PO boxes at address

# port numbers

we run multiple programs on a machine

IP addresses identifying machine — not enough

so, add 16-bit *port numbers*

think: multiple PO boxes at address

0–49151: typically assigned for particular services

80 = http, 443 = https, 22 = ssh, ...

49152–65535: allocated on demand

default “return address” for client connecting to server

# protocols

protocol = agreement on how to communicate

sytnax (format of messages, etc.)

semantics (meaning of messages — actions to take, etc.)

# human protocol: telephone

|                             |                                |
|-----------------------------|--------------------------------|
| caller: pick up phone       |                                |
| caller: check for service   |                                |
| caller: dial                |                                |
| caller: wait for ringing    |                                |
|                             | callee: "Hello?"               |
| caller: "Hi, it's Casey..." |                                |
|                             | callee: "Hi, so how about ..." |
| caller: "Sure, ..."         |                                |
| ...                         | ...                            |
|                             | callee: "Bye!"                 |
| caller: "Bye!"              |                                |
| hang up                     | hang up                        |



# layered protocols

IP: protocol for sending data by IP addresses

- mailbox model

- limited message size

UDP: send *datagrams* built on IP

- still mailbox model, but *with port numbers*

TCP: reliable connections built on IP

- adds port numbers

- adds resending data if error occurs

- splits big amounts of data into many messages

HTTP: protocol for sending files, etc. built on TCP

## other notable protocols (transport layer)

TLS: Transport Layer Security — built on TCP  
like TCP, but adds encryption + authentication

SSH: secure shell (remote login) — built on TCP

SCP/SFTP: secure copy/secure file transfer — built on SSH

HTTPS: HTTP, but over TLS instead of TCP

FTP: file transfer protocol

...

## other notable protocols (transport layer)

TLS: Transport Layer Security — built on TCP  
like TCP, but adds encryption + authentication

SSH: secure shell (remote login) — built on TCP

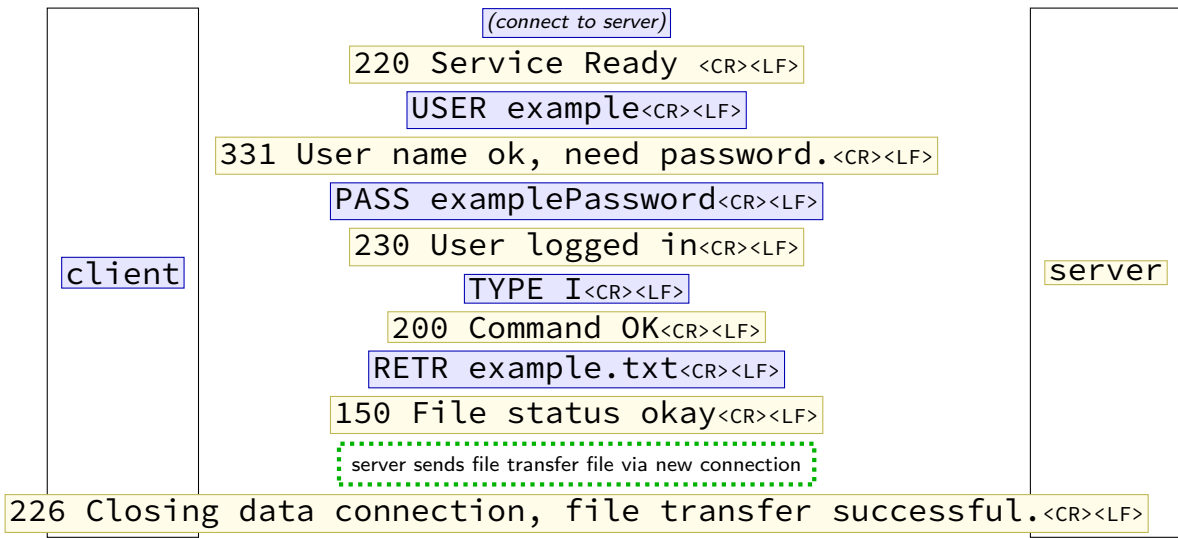
SCP/SFTP: secure copy/secure file transfer — built on SSH

HTTPS: HTTP, but over TLS instead of TCP

FTP: file transfer protocol

...

# FTP protocol (simplified)



# notable things about FTP

FTP is **stateful** — previous commands change future ones

- logging in for whole connection

- change current directory

- set image file type (binary, not text)

FTP uses **separate connections for transferring data**

- PASV: client connects separately to server

- PORT: client specifies where server connects

- (+ very rarely used default: connect back to port 20)

status codes for every command

# sockets

socket: POSIX abstraction of network I/O queue

- any kind of network

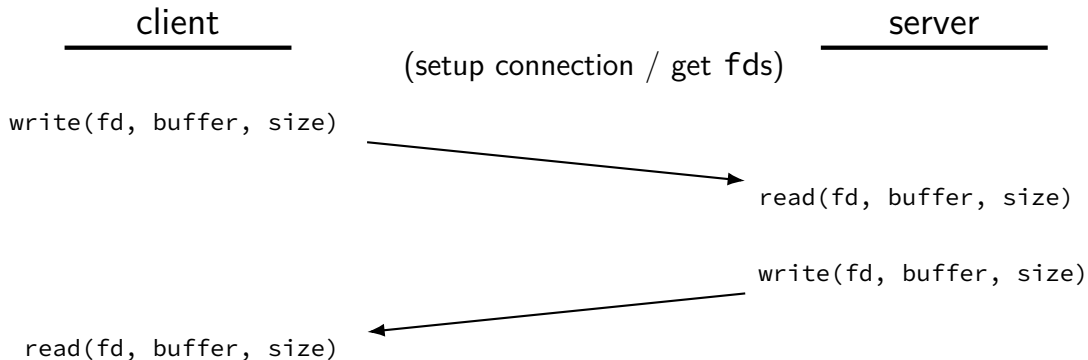
- can also be used between processes on same machine

a kind of **file descriptor**

# connected sockets

sockets can represent a connection

act like **bidirectional pipe**



# echo client/server

```
void client_for_connection(int socket_fd) {
    int n; char send_buf[MAX_SIZE]; char recv_buf[MAX_SIZE];
    while (prompt_for_input(send_buf, MAX_SIZE)) {
        n = write(socket_fd, send_buf, strlen(send_buf));
        if (n != strlen(send_buf)) {...error?...}
        n = read(socket_fd, recv_buf, MAX_SIZE);
        if (n <= 0) return; // error or EOF
        write(STDOUT_FILENO, recv_buf, n);
    }
}
```

---

```
void server_for_connection(int socket_fd) {
    int read_count, write_count; char request_buf[MAX_SIZE];
    while (1) {
        read_count = read(socket_fd, request_buf, MAXSIZE);
        if (read_count <= 0) return; // error or EOF
        write_count = write(socket_fd, request_buf, read_count);
        if (read_count != write_count) {...error?...}
    }
}
```



# echo client/server

```
void client_for_connection(int socket_fd) {
    int n; char send_buf[MAX_SIZE]; char recv_buf[MAX_SIZE];
    while (prompt_for_input(send_buf, MAX_SIZE)) {
        n = write(socket_fd, send_buf, strlen(send_buf));
        if (n != strlen(send_buf)) {...error?...}
        n = read(socket_fd, recv_buf, MAX_SIZE);
        if (n <= 0) return; // error or EOF
        write(STDOUT_FILENO, recv_buf, n);
    }
}
```

---

```
void server_for_connection(int socket_fd) {
    int read_count, write_count; char request_buf[MAX_SIZE];
    while (1) {
        read_count = read(socket_fd, request_buf, MAXSIZE);
        if (read_count <= 0) return; // error or EOF
        write_count = write(socket_fd, request_buf, read_count);
        if (read_count != write_count) {...error?...}
    }
}
```

# echo client/server

```
void client_for_connection(int socket_fd) {
    int n; char send_buf[MAX_SIZE]; char recv_buf[MAX_SIZE];
    while (prompt_for_input(send_buf, MAX_SIZE)) {
        n = write(socket_fd, send_buf, strlen(send_buf));
        if (n != strlen(send_buf)) {...error?...}
        n = read(socket_fd, recv_buf, MAX_SIZE);
        if (n <= 0) return; // error or EOF
        write(STDOUT_FILENO, recv_buf, n);
    }
}
```

---

```
void server_for_connection(int socket_fd) {
    int read_count, write_count; char request_buf[MAX_SIZE];
    while (1) {
        read_count = read(socket_fd, request_buf, MAXSIZE);
        if (read_count <= 0) return; // error or EOF
        write_count = write(socket_fd, request_buf, read_count);
        if (read_count != write_count) {...error?...}
    }
}
```

## aside: send/recv

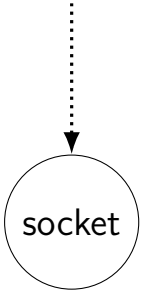
sockets have some alternate read/write-like functions:

recv, recvfrom, recvmsg  
send, sendmsg

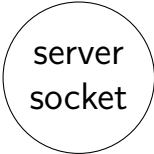
have some additional options we won't need in this class

# sockets and server sockets

```
client:  
fd = socket(...)
```



client

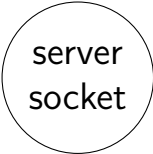
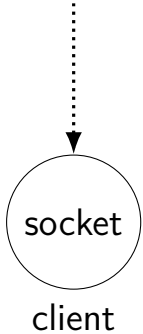


```
server:  
ss_fd = socket(...)  
...  
listen(ss_fd, ...)
```

server

# sockets and server sockets

```
client:  
fd = socket(...)
```



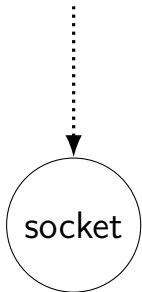
```
server:  
ss_fd = socket(...)  
...  
listen(ss_fd, ...)
```

socket() function — create socket fd

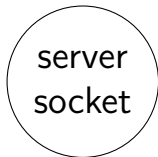
server

# sockets and server sockets

```
client:  
fd = socket(...)
```



client

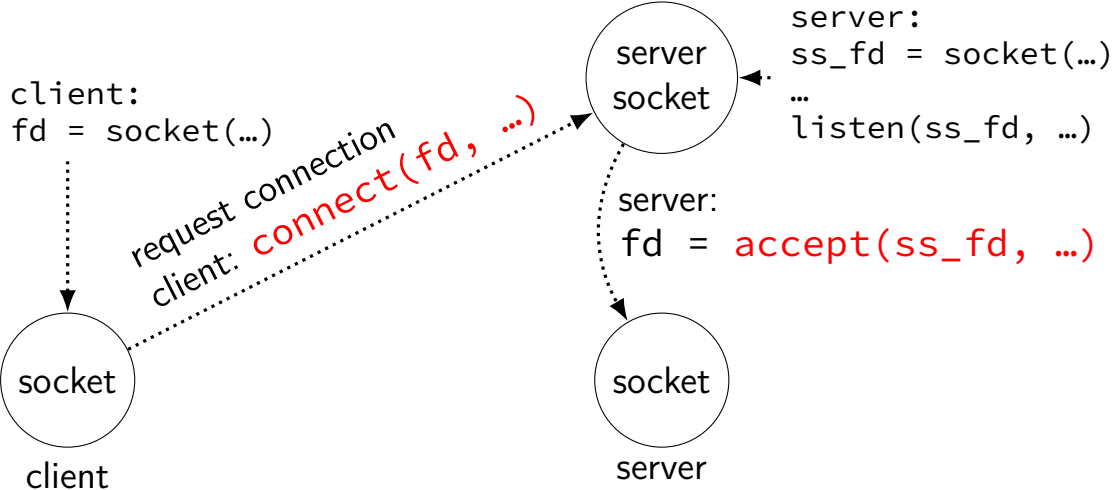


```
server:  
ss_fd = socket(...)  
...  
listen(ss_fd, ...)
```

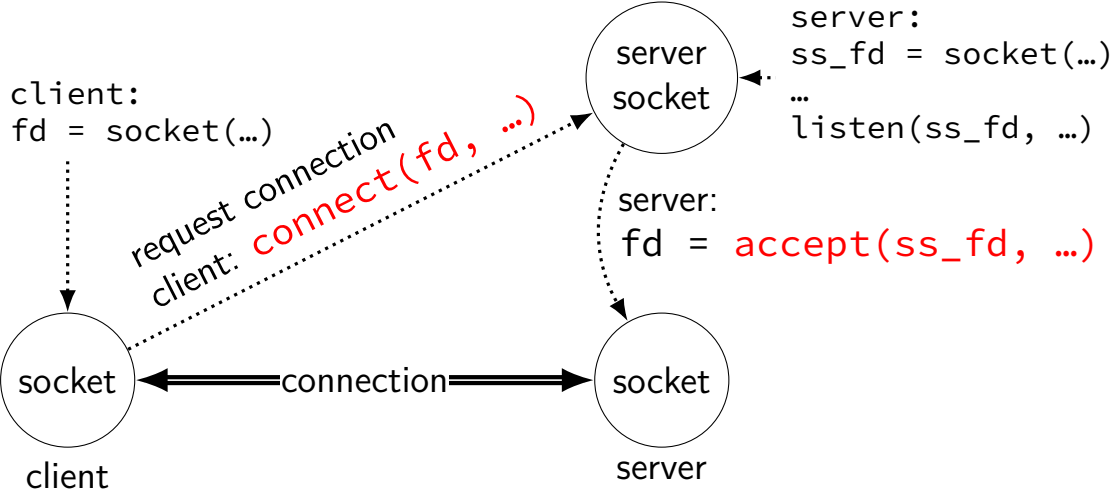
listen() — turn socket into server socket  
still has a file descriptor, but ...  
can only accept() — create normal socket

server

# sockets and server sockets



# sockets and server sockets





# connections in TCP/IP

connection identified by *5-tuple*

used to mark messages sent on network

used by OS to lookup “where is the file descriptor?”

(protocol=TCP, local IP addr., local port, remote IP addr., remote port)

how messages are tagged on the network

(other notable protocol value: UDP)

both ends always have an address+port

what is the IP address, port number? set with `bind()` function

*typically* always done for servers, not done for clients

system will choose default if you don't

# connections on my desktop

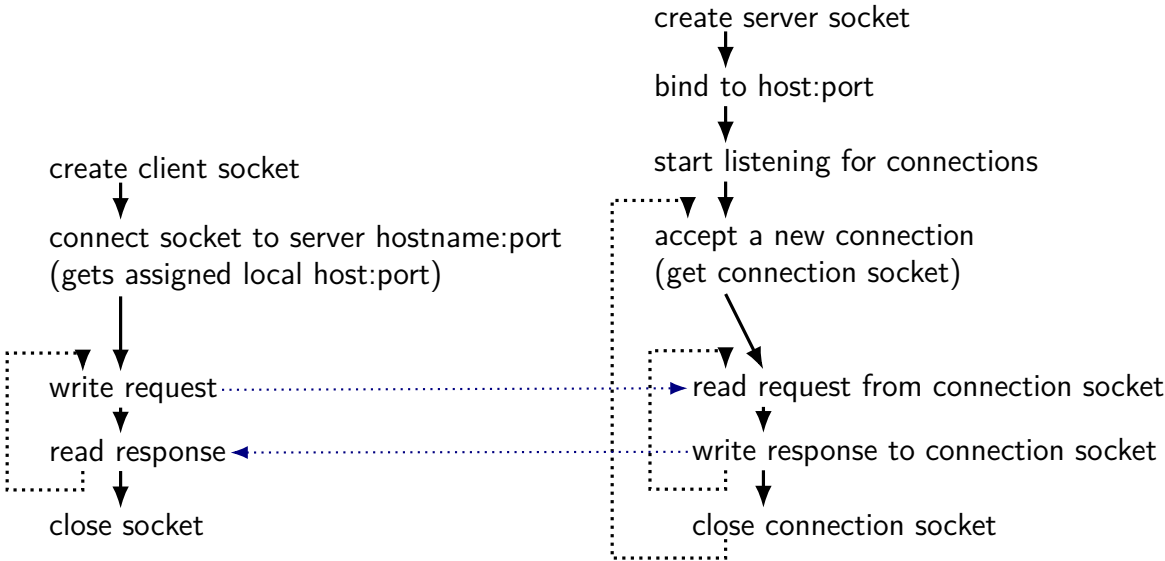
cr4bd@reiss-t3620

: /zf14/cr4bd ; netstat --inet --inet6 --numeric

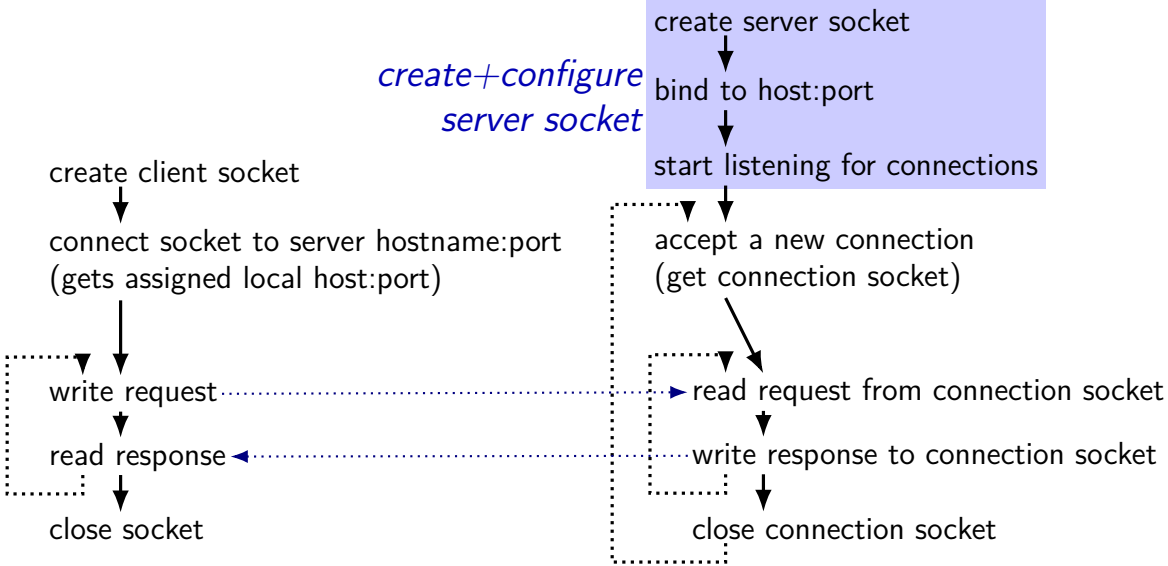
Active Internet connections (w/o servers)

| Proto | Recv-Q | Send-Q | Local Address       | Foreign Address      | State      |
|-------|--------|--------|---------------------|----------------------|------------|
| tcp   | 0      | 0      | 128.143.67.91:49202 | 128.143.63.34:22     | ESTABLISHE |
| tcp   | 0      | 0      | 128.143.67.91:803   | 128.143.67.236:2049  | ESTABLISHE |
| tcp   | 0      | 0      | 128.143.67.91:50292 | 128.143.67.226:22    | TIME_WAIT  |
| tcp   | 0      | 0      | 128.143.67.91:54722 | 128.143.67.236:2049  | TIME_WAIT  |
| tcp   | 0      | 0      | 128.143.67.91:52002 | 128.143.67.236:111   | TIME_WAIT  |
| tcp   | 0      | 0      | 128.143.67.91:732   | 128.143.67.236:63439 | TIME_WAIT  |
| tcp   | 0      | 0      | 128.143.67.91:40664 | 128.143.67.236:2049  | TIME_WAIT  |
| tcp   | 0      | 0      | 128.143.67.91:54098 | 128.143.67.236:111   | TIME_WAIT  |
| tcp   | 0      | 0      | 128.143.67.91:49302 | 128.143.67.236:63439 | TIME_WAIT  |
| tcp   | 0      | 0      | 128.143.67.91:50236 | 128.143.67.236:111   | TIME_WAIT  |
| tcp   | 0      | 0      | 128.143.67.91:22    | 172.27.98.20:49566   | ESTABLISHE |
| tcp   | 0      | 0      | 128.143.67.91:51000 | 128.143.67.236:111   | TIME_WAIT  |
| tcp   | 0      | 0      | 127.0.0.1:50438     | 127.0.0.1:631        | ESTABLISHE |
| tcp   | 0      | 0      | 127.0.0.1:631       | 127.0.0.1:50438      | ESTABLISHE |

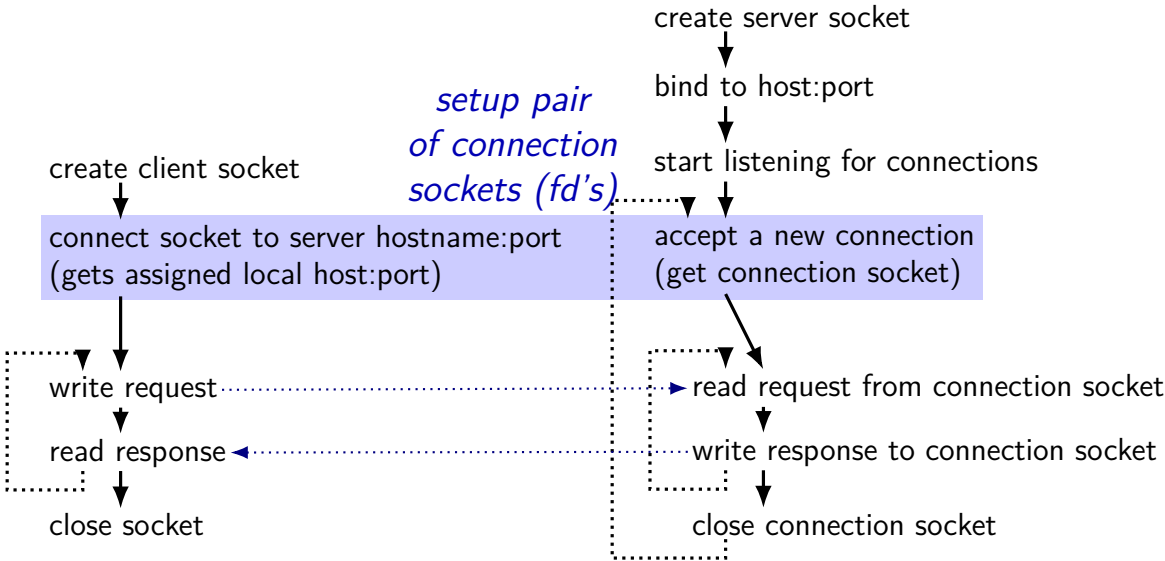
# client/server flow (one connection at a time)



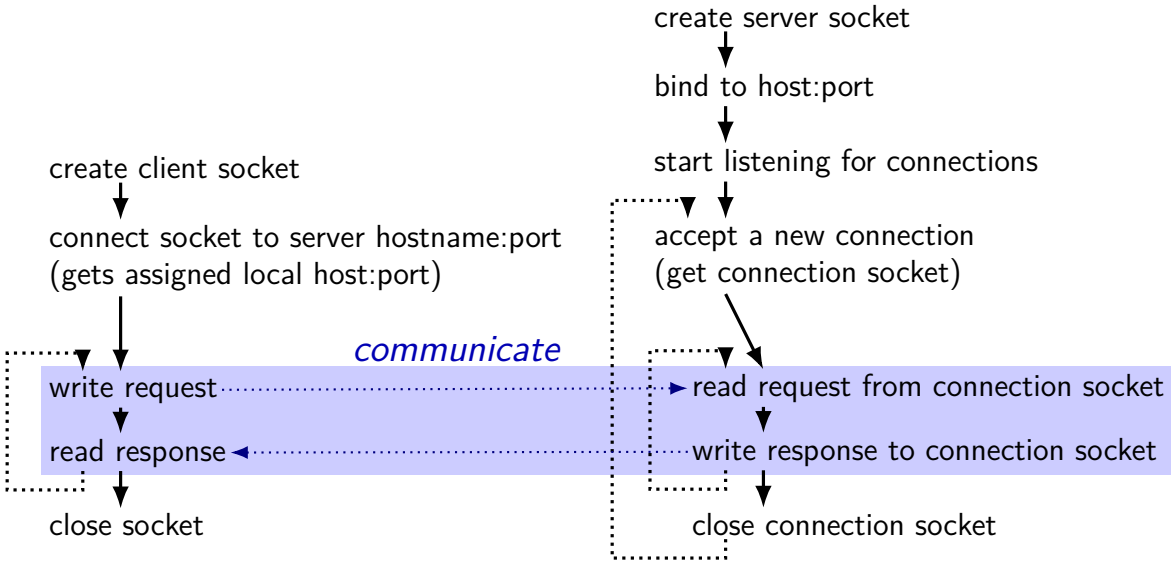
# client/server flow (one connection at a time)



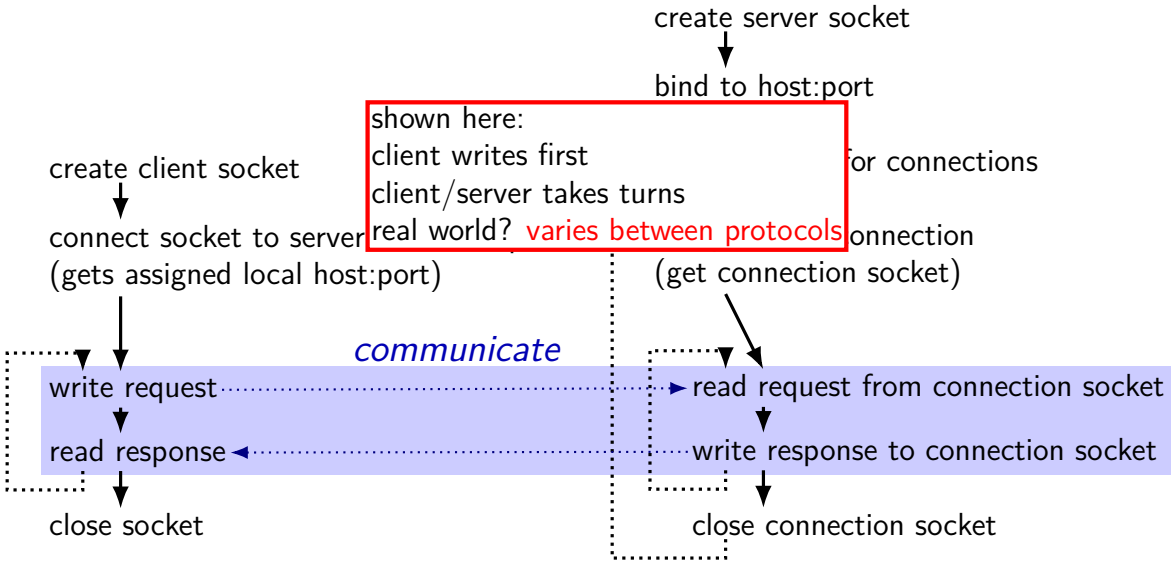
# client/server flow (one connection at a time)



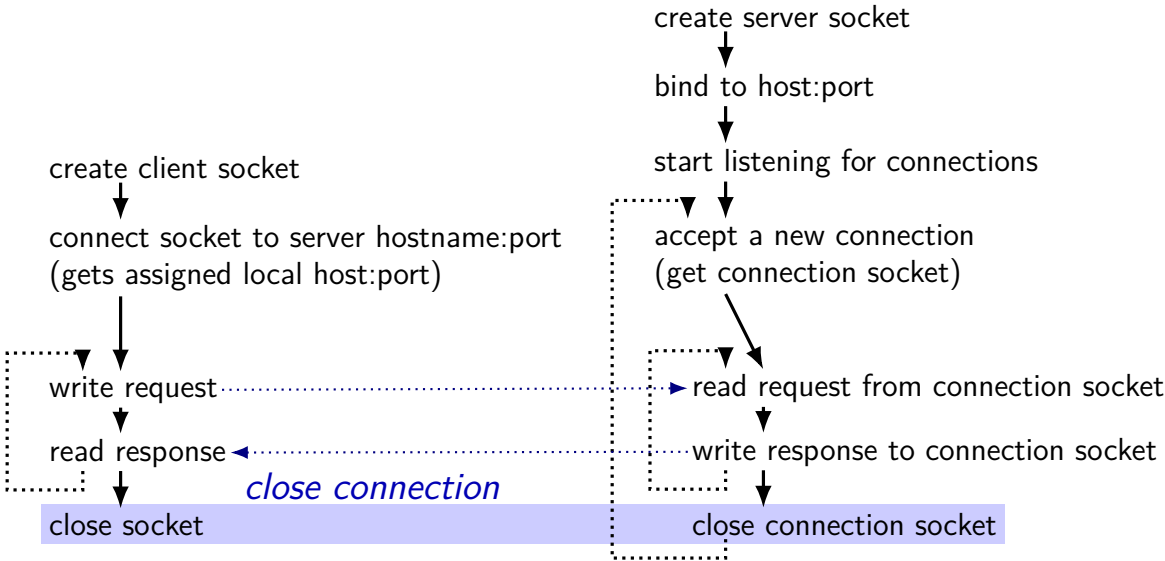
# client/server flow (one connection at a time)



# client/server flow (one connection at a time)

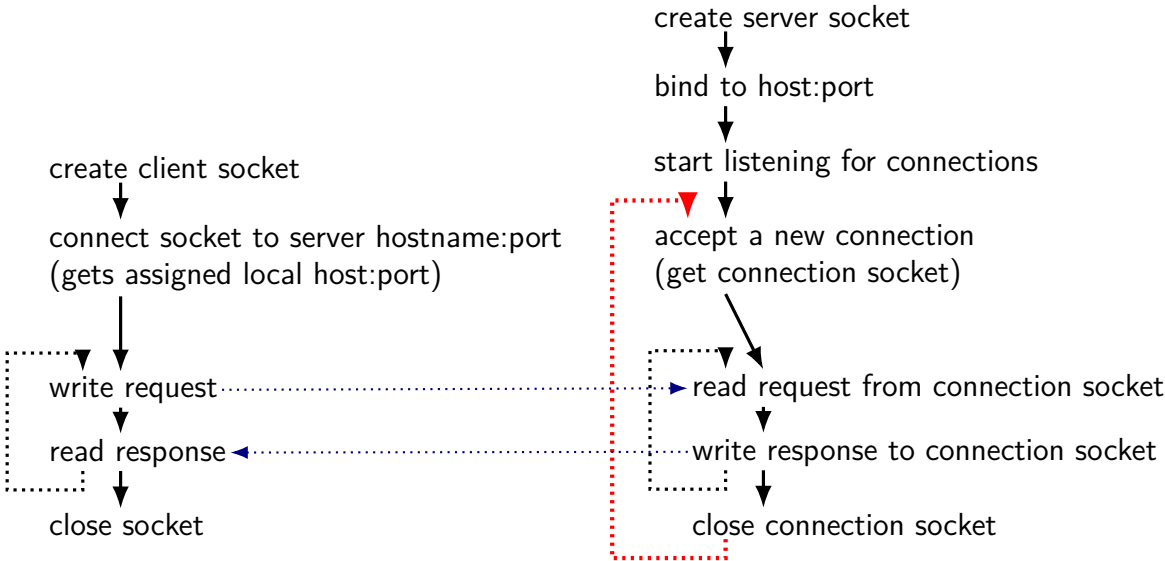


# client/server flow (one connection at a time)





# client/server flow (one connection at a time)



# connection setup: client — manual addresses

```
int sock_fd;

server = /* code on later slide */;
sock_fd = socket(
    AF_INET, /* IPv4 */
    SOCK_STREAM, /* byte-oriented */
    IPPROTO_TCP
);
if (sock_fd < 0) { /* handle error */ }

struct sockaddr_in addr;
addr.sin_family = AF_INET;
addr.sin_addr.s_addr = htonl(2156872459); /* 128.143.67.11 */
addr.sin_port = htons(80); /* port 80 */
if (connect(sock_fd, (struct sockaddr*) &addr, sizeof(addr)) {
    /* handle error */
}
DoClientStuff(sock_fd); /* read and write from sock_fd */
close(sock_fd);
```

# connection setup: client — manual addresses

```
int sock_fd;
```

```
server = /* code on later slide */;
```

```
sock_fd = socket(  
    AF_INET, /* IPv4 */  
    SOCK_STREAM, /* byte-oriented */  
    IPPROTO_TCP  
);
```

if specify IPv4 instead of IPv6 or local-only sockets

specify TCP (byte-oriented) instead of UDP ('datagram' oriented)

```
struct sockaddr_in  
addr.sin_family = AF_INET;
```

```
addr.sin_addr.s_addr = htonl(2156872459); /* 128.143.67.11 */
```

```
addr.sin_port = htons(80); /* port 80 */
```

```
if (connect(sock_fd, (struct sockaddr*) &addr, sizeof(addr)) {
```

```
    /* handle error */
```

```
}
```

```
DoClientStuff(sock_fd); /* read and write from sock_fd */
```

```
close(sock_fd);
```

# connection setup: client — manual addresses

```
int sock_fd;

server = /* code */
sock_fd = socket(
    AF_INET, /* IPv4 */
    SOCK_STREAM, /* byte-oriented */
    IPPROTO_TCP
);
if (sock_fd < 0) { /* handle error */ }

struct sockaddr_in addr;
addr.sin_family = AF_INET;
addr.sin_addr.s_addr = htonl(2156872459); /* 128.143.67.11 */
addr.sin_port = htons(80); /* port 80 */
if (connect(sock_fd, (struct sockaddr*) &addr, sizeof(addr)) {
    /* handle error */
}
DoClientStuff(sock_fd); /* read and write from sock_fd */
close(sock_fd);
```

htonl/s = host-to-network long/short  
network byte order = big endian

## connection setup: client — manual addresses

```
int sock_fd;
server = /* struct representing IPv4 address + port number
          declared in <netinet/in.h>
sock_fd = s see man 7 ip on Linux for docs
            AF_INET
            SOCK_STREAM, /* byte-oriented */
            IPPROTO_TCP
);
if (sock_fd < 0) { /* handle error */ }

struct sockaddr_in addr;
addr.sin_family = AF_INET;
addr.sin_addr.s_addr = htonl(2156872459); /* 128.143.67.11 */
addr.sin_port = htons(80); /* port 80 */
if (connect(sock_fd, (struct sockaddr*) &addr, sizeof(addr)) {
    /* handle error */
}
DoClientStuff(sock_fd); /* read and write from sock_fd */
close(sock_fd);
```

# sockaddr\_in

```
/* from 'man 7 ip' */  
struct sockaddr_in {  
    sa_family_t    sin_family; /* address family: always AF_INET */  
    in_port_t      sin_port;    /* port in network byte order */  
    struct in_addr sin_addr;    /* internet address */  
};
```

```
/* Internet address. */  
struct in_addr {  
    uint32_t        s_addr;     /* address in network byte order */  
};
```

# sockaddr\_in

```
/* from 'man 7 ip' */
struct sockaddr_in {
    sa_family_t    sin_family; /* address family: always AF_INET */
    in_port_t      sin_port;   /* port in network byte order */
    struct in_addr sin_addr;   /* internet address */
};
```

```
/* Internet address. */
struct in_addr {
    uint32_t       s_addr;    /* address in network byte order */
};
```

# sockaddr\_in

```
/* from 'man 7 ip' */
struct sockaddr_in {
    sa_family_t    sin_family; /* address family: always AF_INET */
    in_port_t      sin_port;    /* port in network byte order */
    struct in_addr sin_addr;    /* internet address */
};
```

```
/* Internet address. */
struct in_addr {
    uint32_t      s_addr;      /* address in network byte order */
};
```

trick: multiple versions of address struct  
each have “type” information in same spot  
OS/library checks before using



# sockaddr\_in6

```
/* from 'man 7 ipv6' */
struct sockaddr_in6 {
    sa_family_t      sin6_family;    /* always AF_INET6 */
    in_port_t        sin6_port;      /* port number */
    uint32_t          sin6_flowinfo; /* IPv6 flow information */
    struct in6_addr   sin6_addr;      /* IPv6 address */
    uint32_t          sin6_scope_id;  /* Scope ID (new in 2.4) */
};

struct in6_addr {
    unsigned char     s6_addr[16];    /* IPv6 address */
};
```

# sockaddr\_in6

```
/* from 'man 7 ipv6' */
struct sockaddr_in6 {
    sa_family_t    sin6_family;    /* always AF_INET6 */
    in_port_t      sin6_port;      /* port number */
    uint32_t       sin6_flowinfo;  /* IPv6 flow information */
    struct in6_addr sin6_addr;      /* IPv6 address */
    uint32_t       sin6_scope_id;   /* Scope ID (new in 2.4) */
};

struct in6_addr {
    unsigned char  s6_addr[16];    /* IPv6 address */
};
```

## connection setup: client, using addrinfo

```
int sock_fd;
struct addrinfo *server = /* code on next slide */;

sock_fd = socket(
    server->ai_family,
    // ai_family = AF_INET (IPv4) or AF_INET6 (IPv6) or ...
    server->ai_socktype,
    // ai_socktype = SOCK_STREAM (bytes) or ...
    server->ai_protocol,
    // ai_protocol = IPPROTO_TCP or ...
);
if (sock_fd < 0) { /* handle error */ }
if (connect(sock_fd, server->ai_addr, server->ai_addrlen) < 0) {
    /* handle error */
}
freeaddrinfo(server);
DoClientStuff(sock_fd); /* read and write from sock_fd */
close(sock_fd);
```

# connection setup: client, using addrinfo

```
int sock_fd;
struct addrinfo *server = /* code on next slide */;

sock_fd = socket(
    server->ai_family,
    // ai_family = AF_INET (IPv4) or AF_INET6 (IPv6) or ...
    server->ai_socktype,
    // server->ai_protocol (TCP or UDP)
    server);
// addrinfo contains all information needed to setup socket
// set by getaddrinfo function (next slide)
// handles IPv4 and IPv6
// handles DNS names, service names
if (socket_error != 0) {
    /* handle error */
}
freeaddrinfo(server);
DoClientStuff(sock_fd); /* read and write from sock_fd */
close(sock_fd);
```

# connection setup: client, using addrinfo

```
int sock_fd;
struct addrinfo *server = /* code on next slide */;

sock_fd = socket(
    server->ai_family,
    // ai_family = AF_INET (IPv4) or AF_INET6 (IPv6) or ...
    server->ai_socktype,
    // ai_socktype = SOCK_STREAM (bytes) or ...
    server->ai_protocol
    // ai_protocol = IPPROTO_TCP or ...
);
if (sock_fd < 0) { /* handle error */ }
if (connect(sock_fd, server->ai_addr, server->ai_addrlen) < 0) {
    /* handle error */
}
freeaddrinfo(server);
DoClientStuff(sock_fd); /* read and write from sock_fd */
close(sock_fd);
```

## connection setup: client, using addrinfo

```
int sock_fd;
struct addrinfo server;

sock_fd = server;
    // ai_family = AF_INET (IPv4) or AF_INET6 (IPv6) or ...
    server->ai_socktype,
    // ai_socktype = SOCK_STREAM (bytes) or ...
    server->ai_protocol
    // ai_protocol = IPPROTO_TCP or ...
);
if (sock_fd < 0) { /* handle error */ }
if (connect(sock_fd, server->ai_addr, server->ai_addrlen) < 0) {
    /* handle error */
}
freeaddrinfo(server);
DoClientStuff(sock_fd); /* read and write from sock_fd */
close(sock_fd);
```

## connection setup: client, using addrinfo

```
int  
str  
sock_fd = socket(  
    server->ai_family,  
    // ai_family = AF_INET (IPv4) or AF_INET6 (IPv6) or ...  
    server->ai_socktype,  
    // ai_socktype = SOCK_STREAM (bytes) or ...  
    server->ai_protocol  
    // ai_protocol = IPPROTO_TCP or ...  
);  
if (sock_fd < 0) { /* handle error */ }  
if (connect(sock_fd, server->ai_addr, server->ai_addrlen) < 0) {  
    /* handle error */  
}  
freeaddrinfo(server);  
DoClientStuff(sock_fd); /* read and write from sock_fd */  
close(sock_fd);
```

## connection setup: lookup address

```
/* example hostname, portname = "www.cs.virginia.edu", "443" */
const char *hostname; const char *portname;
...
struct addrinfo *server;
struct addrinfo hints;
int rv;
memset(&hints, 0, sizeof(hints));
hints.ai_family = AF_UNSPEC; /* for IPv4 OR IPv6 */
// hints.ai_family = AF_INET4; /* for IPv4 only */

hints.ai_socktype = SOCK_STREAM; /* byte-oriented --- TCP */
rv = getaddrinfo(hostname, portname, &hints, &server);
if (rv != 0) { /* handle error */ }

/* eventually freeaddrinfo(result) */
```



## connection setup: lookup address

```
/* example hostname, portname = "www.cs.virginia.edu", "443" */
const char *hostname; const char *portname;
...
struct addrinfo *server;
struct addrinfo hints;
int rv;
memset(&hints, 0, sizeof(hints));
hints.ai_family = AF_INET;
// hints.ai_flags = AI_NUMERICSERV;
NB: pass pointer to pointer to addrinfo to fill in
hints.ai_socktype = SOCK_STREAM; /* byte-oriented --- TCP */
rv = getaddrinfo(hostname, portname, &hints, &server);
if (rv != 0) { /* handle error */ }

/* eventually freeaddrinfo(result) */
```

## connection setup: lookup address

```
/* exam AF_UNSPEC: choose between IPv4 and IPv6 for me */  
const d AF_INET, AF_INET6: choose IPv4 or IPV6 respectively  
...  
struct addrinfo *server;  
struct addrinfo hints;  
int rv;  
memset(&hints, 0, sizeof(hints));  
hints.ai_family = AF_UNSPEC; /* for IPv4 OR IPv6 */  
// hints.ai_family = AF_INET4; /* for IPv4 only */  
  
hints.ai_socktype = SOCK_STREAM; /* byte-oriented --- TCP */  
rv = getaddrinfo(hostname, portname, &hints, &server);  
if (rv != 0) { /* handle error */ }  
  
/* eventually freeaddrinfo(result) */
```

## connection setup: multiple server addresses

```
struct addrinfo *server;
...
rv = getaddrinfo(hostname, portname, &hints, &server);
if (rv != 0) { /* handle error */ }

for (struct addrinfo *current = server; current != NULL;
     current = current->ai_next) {
    sock_fd = socket(current->ai_family, current->ai_socktype, current->ai_protocol);
    if (sock_fd < 0) continue;
    if (connect(sock_fd, current->ai_addr, current->ai_addrlen) == 0)
        break;
}
close(sock_fd); // connect failed
}
freeaddrinfo(server);
DoClientStuff(sock_fd);
close(sock_fd);
```

# connection setup: multiple server addresses

```
struct addrinfo *server;
...
rv = getaddrinfo(hostname, portname, &hints, &server);
if (rv != 0) { /* handle error */ }

for (struct addrinfo *current = server; current != NULL;
     current = current->ai_next) {
    sock_fd = socket(current->ai_family, current->ai_socktype, current->ai_protocol);
    if (sock_fd < 0) continue;
    if (connect(sock_fd, current->ai_addr, current->ai_addrlen) == 0)
        break;
}
close(sock_fd);
}
freeaddrinfo(server);
DoClient();
close(sock_fd);
```

addrinfo is a linked list

name can correspond to multiple addresses

example: redundant copies of web server

example: an IPv4 address and IPv6 address

example: wired + wireless connection on one machine

## connection setup: old lookup function

```
/* example hostname, portnum= "www.cs.virginia.edu", 443*/
const char *hostname; int portnum;
...
struct hostent *server_ip;
server_ip = gethostbyname(hostname);

if (server_ip == NULL) { /* handle error */ }

struct sockaddr_in addr;
addr.s_addr = *(struct in_addr*) server_ip->h_addr_list[0];
addr.sin_port = htons(portnum);
sock_fd = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP);
connect(sock_fd, &addr, sizeof(addr));
...
```

## connection setup: server, address setup

```
/* example (hostname, portname) = ("127.0.0.1", "443") */
const char *hostname; const char *portname;
...
struct addrinfo *server;
struct addrinfo hints;
int rv;

memset(&hints, 0, sizeof(hints));
hints.ai_family = AF_INET; /* for IPv4 */
/* or: */ hints.ai_family = AF_INET6; /* for IPv6 */
/* or: */ hints.ai_family = AF_UNSPEC; /* I don't care */
hints.ai_flags = AI_PASSIVE;

rv = getaddrinfo(hostname, portname, &hints, &server);
if (rv != 0) { /* handle error */ }
```

# connection setup: server, address setup

```
/* example (hostname, portname) = ("127.0.0.1", "443") */
const char *hostname; const char *portname;
...
struct addrinfo *server;
struct addrinfo hints;
int rv;

memset(&hints, 0, sizeof(hints));
hints.ai_family = AF_INET; /* for IPv4 */
/* or: */ hints.ai_family = AF_INET6; /* for IPv6 */
/* or: */ hints.ai_family = AF_UNSPEC; /* care */
hints.ai_flags = AI_PASSIVE; /* hostname could also be NULL
                               means "use all possible addresses"
                               only makes sense for servers */
rv = getaddrinfo(hostname, portname, &hints, &server);
if (rv != 0) { /* handle error */ }
```

# connection setup: server, address setup

```
/* example (hostname, portname) = ("127.0.0.1", "443") */
const char *hostname; const char *portname;
...
struct addrinfo *server;
struct addrinfo hints;
int rv;

memset(&hints, 0, sizeof(hints));
hints.ai_family = AF_INET; /* for IPv4 */
/* or: */ hints.ai_family = AF_INET6; /* for IPv6 */
/* or: */ hints.ai_family = AF_UNSPEC; /* either are */
hints.ai_flags = AI_NUMERICSERV; /* portname could also be NULL
means "choose a port number for me"
only makes sense for servers */
rv = getaddrinfo(hostname, portname, &hints, NULL, &server, &rv);
if (rv != 0) { /* handle error */ }
```



## connection setup: server, address setup

```
/* example (hostname, portname, &hints, &server) */
const char *hostname = "localhost";
...
struct addrinfo *server;
struct addrinfo hints;
int rv;

memset(&hints, 0, sizeof(hints));
hints.ai_family = AF_INET; /* for IPv4 */
/* or: */ hints.ai_family = AF_INET6; /* for IPv6 */
/* or: */ hints.ai_family = AF_UNSPEC; /* I don't care */
hints.ai_flags = AI_PASSIVE;

rv = getaddrinfo(hostname, portname, &hints, &server);
if (rv != 0) { /* handle error */ }
```

## connection setup: server, addrinfo

```
struct addrinfo *server;
... getaddrinfo(...) ...

int server_socket_fd = socket(
    server->ai_family,
    server->ai_socktype,
    server->ai_protocol
);

if (bind(server_socket_fd, ai->ai_addr, ai->ai_addr_len)) < 0) {
    /* handle error */
}
listen(server_socket_fd, MAX_NUM_WAITING);
...
int socket_fd = accept(server_socket_fd, NULL);
```

## aside: on server port numbers

Unix convention: must be `root` to use ports 0–1023

`root` = superuser = 'administrator user' = what `sudo` does

so, for testing: probably ports  $> 1023$

## connection setup: server, manual

```
int server_socket_fd = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP);
struct sockaddr_in addr;
addr.sin_family = AF_INET;
addr.sin_addr.s_addr = INADDR_ANY; /* "any address I can use" */
    /* or: addr.s_addr.in_addr = INADDR_LOOPBACK (127.0.0.1) */
    /* or: addr.s_addr.in_addr = htonl(...); */
addr.sin_port = htons(9999); /* port number 9999 */

if (bind(server_socket_fd, &addr, sizeof(addr)) < 0) {
    /* handle error */
}
listen(server_socket_fd, MAX_NUM_WAITING);
...
int socket_fd = accept(server_socket_fd, NULL);
```

## connection setup: server, manual

```
int server_socket_fd = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP);
struct sockaddr_in addr;
addr.sin_family = AF_INET;
addr.sin_addr.s_addr = INADDR_ANY; /* "any address I can use" */
/* or: addr.s_addr.in_addr = INADDR_LOOPBACK (127.0.0.1) */
/* or: addr.s_addr.in_addr = htonl(...); */
addr.sin_port = htons(9999); /* port number 9999 */

if (bind(server_socket_fd, &addr, sizeof(addr)) < 0) {
    /* handle error */
}
listen
...
int socket_fd = accept(server_socket_fd, NULL);
```

INADDR\_ANY: accept connections for any address I can!  
alternative: specify specific address

## connection setup: server, manual

```
int server_socket_fd = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP);
struct sockaddr_in addr;
addr.sin_family = AF_INET;
addr.sin_addr.s_addr = INADDR_ANY; /* "any address I can use" */
/* or: addr.s_addr.in_addr = INADDR_LOOPBACK (127.0.0.1) */
/* or: addr.s_addr.in_addr = htonl(...); */
addr.sin_port = htons(9999); /* port number 9999 */

if (bind(server_socket_fd, &addr, sizeof(addr)) < 0) {
    /* handle error */
}
list
...
int socket_fd = accept(server_socket_fd, NULL);
```

bind to 127.0.0.1? only accept connections from same machine  
what we recommend for FTP server assignment

## connection setup: server, manual

```
int server_socket_fd = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP);
struct sockaddr_in addr;
addr.sin_family = AF_INET;
addr.sin_addr.s_addr = INADDR_ANY; /* "any address I can use" */
    /* or: addr.s_addr.in_addr = INADDR_LOOPBACK (127.0.0.1) */
    /* or: addr.s_addr.in_addr = htonl(...); */
addr.sin_port = htons(9999); /* port number 9999 */

if (bind(server_socket_fd, &addr, sizeof(addr)) < 0) {
    /* handle error */
}
listen(server_socket_fd, MAX_NUM_WAITING),
...
int socket_fd = accept(server_socket_fd, NULL);
```

choose the number of unaccepted connections

# client/server flow (multiple connections)

