

6. Slightly Advanced Techniques

These aren't *really* advanced, but they're getting out of the more basic levels we've already covered. In fact, if you've gotten this far, you should consider yourself fairly accomplished in the basics of Unix network programming! Congratulations!

So here we go into the brave new world of some of the more esoteric things you might want to learn about sockets. Have at it!

6.1. Blocking

Blocking. You've heard about it--now what the heck is it? In a nutshell, "block" is techie jargon for "sleep". You probably noticed that when you run **listener**, above, it just sits there until a packet arrives. What happened is that it called `recvfrom()`, there was no data, and so `recvfrom()` is said to "block" (that is, sleep there) until some data arrives.

Lots of functions block. `accept()` blocks. All the `recv()` functions block. The reason they can do this is because they're allowed to. When you first create the socket descriptor with `socket()`, the kernel sets it to blocking. If you don't want a socket to be blocking, you have to make a call to `fcntl()`:

```
#include <unistd.h>
#include <fcntl.h>
.
.
.
sockfd = socket(PF_INET, SOCK_STREAM, 0);
fcntl(sockfd, F_SETFL, O_NONBLOCK);
.
.
.
```

By setting a socket to non-blocking, you can effectively "poll" the socket for information. If you try to read from a non-blocking socket and there's no data there, it's not allowed to block--it will return **-1** and *errno* will be set to **EWOULDBLOCK**.

Generally speaking, however, this type of polling is a bad idea. If you put your program in a busy-wait looking for data on the socket, you'll suck up CPU time like it was going out of style. A more elegant solution for checking to see if there's data waiting to be read comes in the following section on `select()`.

6.2. `select()`--Synchronous I/O Multiplexing

This function is somewhat strange, but it's very useful. Take the following situation: you are a server and you want to listen for incoming connections as well as keep reading from the connections you already have.

No problem, you say, just an `accept()` and a couple of `recv()`s. Not so fast, buster! What if you're blocking on an `accept()` call? How are you going to `recv()` data at the same time? "Use non-blocking sockets!" No way! You don't want to be a CPU hog. What, then?

`select()` gives you the power to monitor several sockets at the same time. It'll tell you which ones are ready for reading, which are ready for writing, and which sockets have raised exceptions, if you really want to know that.

Without any further ado, I'll offer the synopsis of `select()`:

```
#include <sys/time.h>
#include <sys/types.h>
#include <unistd.h>

int select(int numfds, fd_set *readfds, fd_set *writefds,
          fd_set *exceptfds, struct timeval *timeout);
```

The function monitors "sets" of file descriptors; in particular *readfds*, *writefds*, and *exceptfds*. If you want to see if you can read from standard input and some socket descriptor, *sockfd*, just add the file descriptors `0` and *sockfd* to the set *readfds*. The parameter *numfds* should be set to the values of the highest file descriptor plus one. In this example, it should be set to *sockfd*+1, since it is assuredly higher than standard input (`0`).

When `select()` returns, *readfds* will be modified to reflect which of the file descriptors you selected which is ready for reading. You can test them with the macro `FD_ISSET()`, below.

Before progressing much further, I'll talk about how to manipulate these sets. Each set is of the type `fd_set`. The following macros operate on this type:

- `FD_ZERO(fd_set *set)` -- clears a file descriptor set
- `FD_SET(int fd, fd_set *set)` -- adds *fd* to the set
- `FD_CLR(int fd, fd_set *set)` -- removes *fd* from the set
- `FD_ISSET(int fd, fd_set *set)` -- tests to see if *fd* is in the set

Finally, what is this weirded out `struct timeval`? Well, sometimes you don't want to wait forever for someone to send you some data. Maybe every 96 seconds you want to print "Still Going..." to the terminal even though nothing has happened. This time structure allows you to specify a timeout period. If the time is exceeded and `select()` still hasn't found any ready file descriptors, it'll return so you can continue processing.

The `struct timeval` has the follow fields:

```
struct timeval {
    int tv_sec;    // seconds
    int tv_usec;   // microseconds
};
```

Just set `tv_sec` to the number of seconds to wait, and set `tv_usec` to the number of microseconds to wait. Yes, that's *microseconds*, not milliseconds. There are 1,000 microseconds in a millisecond, and 1,000 milliseconds in a second. Thus, there are 1,000,000 microseconds in a second. Why is it "usec"? The "u" is supposed to look like the Greek letter μ (Mu) that we use for "micro". Also, when the function returns, *timeout might* be updated to show the time still remaining. This depends on what flavor of Unix you're running.

Yay! We have a microsecond resolution timer! Well, don't count on it. Standard Unix timeslice is around 100 milliseconds, so you might have to wait that long no matter how small you set your `struct timeval`.

Other things of interest: If you set the fields in your `struct timeval` to 0, `select()` will timeout immediately, effectively polling all the file descriptors in your sets. If you set the parameter *timeout* to NULL, it will never timeout, and will wait until the first file descriptor is ready. Finally, if you don't care about waiting for a certain set, you can just set it to NULL in the call to `select()`.

[The following code snippet](#) waits 2.5 seconds for something to appear on standard input:

```
/*
** select.c -- a select() demo
*/

#include <stdio.h>
#include <sys/time.h>
#include <sys/types.h>
#include <unistd.h>

#define STDIN 0 // file descriptor for standard input

int main(void)
{
    struct timeval tv;
    fd_set readfds;

    tv.tv_sec = 2;
    tv.tv_usec = 500000;

    FD_ZERO(&readfds);
    FD_SET(STDIN, &readfds);

    // don't care about writefds and exceptfds:
    select(STDIN+1, &readfds, NULL, NULL, &tv);

    if (FD_ISSET(STDIN, &readfds))
```

```

        printf("A key was pressed!\n");
    else
        printf("Timed out.\n");

    return 0;
}

```

If you're on a line buffered terminal, the key you hit should be RETURN or it will time out anyway.

Now, some of you might think this is a great way to wait for data on a datagram socket--and you are right: it *might* be. Some Unices can use select in this manner, and some can't. You should see what your local man page says on the matter if you want to attempt it.

Some Unices update the time in your struct `timeval` to reflect the amount of time still remaining before a timeout. But others do not. Don't rely on that occurring if you want to be portable. (Use `gettimeofday()` if you need to track time elapsed. It's a bummer, I know, but that's the way it is.)

What happens if a socket in the read set closes the connection? Well, in that case, `select()` returns with that socket descriptor set as "ready to read". When you actually do `recv()` from it, `recv()` will return `0`. That's how you know the client has closed the connection.

One more note of interest about `select()`: if you have a socket that is `listen()`ing, you can check to see if there is a new connection by putting that socket's file descriptor in the *readfds* set.

And that, my friends, is a quick overview of the almighty `select()` function.

But, by popular demand, here is an in-depth example. Unfortunately, the difference between the dirt-simple example, above, and this one here is significant. But have a look, then read the description that follows it.

[This program](#) acts like a simple multi-user chat server. Start it running in one window, then **telnet** to it ("**telnet hostname 9034**") from multiple other windows. When you type something in one **telnet** session, it should appear in all the others.

```

/*
** selectserver.c -- a cheezy multiperson chat server
*/

#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>

#define PORT 9034    // port we're listening on

```

```
int main(void)
{
    fd_set master;    // master file descriptor list
    fd_set read_fds;  // temp file descriptor list for select()
    struct sockaddr_in myaddr;    // server address
    struct sockaddr_in remoteaddr; // client address
    int fdmax;    // maximum file descriptor number
    int listener; // listening socket descriptor
    int newfd;    // newly accept()ed socket descriptor
    char buf[256]; // buffer for client data
    int nbytes;
    int yes=1;    // for setsockopt() SO_REUSEADDR, below
    socklen_t addrlen;
    int i, j;

    FD_ZERO(&master);    // clear the master and temp sets
    FD_ZERO(&read_fds);

    // get the listener
    if ((listener = socket(PF_INET, SOCK_STREAM, 0)) == -1) {
        perror("socket");
        exit(1);
    }

    // lose the pesky "address already in use" error message
    if (setsockopt(listener, SOL_SOCKET, SO_REUSEADDR, &yes,
        sizeof(int)) == -1) {
        perror("setsockopt");
        exit(1);
    }

    // bind
    myaddr.sin_family = AF_INET;
    myaddr.sin_addr.s_addr = INADDR_ANY;
    myaddr.sin_port = htons(PORT);
    memset(&(myaddr.sin_zero), '\0', 8);
    if (bind(listener, (struct sockaddr *)&myaddr, sizeof(myaddr)) == -1) {
        perror("bind");
        exit(1);
    }

    // listen
    if (listen(listener, 10) == -1) {
        perror("listen");
        exit(1);
    }

    // add the listener to the master set
    FD_SET(listener, &master);
```

```

// keep track of the biggest file descriptor
fdmax = listener; // so far, it's this one

// main loop
for(;;) {
    read_fds = master; // copy it
    if (select(fdmax+1, &read_fds, NULL, NULL, NULL) == -1) {
        perror("select");
        exit(1);
    }

    // run through the existing connections looking for data to read
    for(i = 0; i <= fdmax; i++) {
        if (FD_ISSET(i, &read_fds)) { // we got one!!
            if (i == listener) {
                // handle new connections
                addrlen = sizeof(remoteaddr);
                if ((newfd = accept(listener, (struct sockaddr *)&remoteaddr,
                                    &addrlen)) == -1) {
                    perror("accept");
                } else {
                    FD_SET(newfd, &master); // add to master set
                    if (newfd > fdmax) { // keep track of the maximum
                        fdmax = newfd;
                    }
                    printf("selectserver: new connection from %s on "
                           "socket %d\n", inet_ntoa(remoteaddr.sin_addr), newfd);
                }
            } else {
                // handle data from a client
                if ((nbytes = recv(i, buf, sizeof(buf), 0)) <= 0) {
                    // got error or connection closed by client
                    if (nbytes == 0) {
                        // connection closed
                        printf("selectserver: socket %d hung up\n", i);
                    } else {
                        perror("recv");
                    }
                    close(i); // bye!
                    FD_CLR(i, &master); // remove from master set
                } else {
                    // we got some data from a client
                    for(j = 0; j <= fdmax; j++) {
                        // send to everyone!
                        if (FD_ISSET(j, &master)) {
                            // except the listener and ourselves
                            if (j != listener && j != i) {
                                if (send(j, buf, nbytes, 0) == -1) {

```

```

    }
    }
    }
    }
    }
    } // it's SO UGLY!
    }
    }
    }
    return 0;
}
perror("send");
}

```

Notice I have two file descriptor sets in the code: *master* and *read_fds*. The first, *master*, holds all the socket descriptors that are currently connected, as well as the socket descriptor that is listening for new connections.

The reason I have the *master* set is that `select()` actually *changes* the set you pass into it to reflect which sockets are ready to read. Since I have to keep track of the connections from one call of `select()` to the next, I must store these safely away somewhere. At the last minute, I copy the *master* into the *read_fds*, and then call `select()`.

But doesn't this mean that every time I get a new connection, I have to add it to the *master* set? Yup! And every time a connection closes, I have to remove it from the *master* set? Yes, it does.

Notice I check to see when the *listener* socket is ready to read. When it is, it means I have a new connection pending, and I `accept()` it and add it to the *master* set. Similarly, when a client connection is ready to read, and `recv()` returns `0`, I know the client has closed the connection, and I must remove it from the *master* set.

If the client `recv()` returns non-zero, though, I know some data has been received. So I get it, and then go through the *master* list and send that data to all the rest of the connected clients.

And that, my friends, is a less-than-simple overview of the almighty `select()` function.

6.3. Handling Partial `send()`s

Remember back in the [section about `send\(\)`](#), above, when I said that `send()` might not send all the bytes you asked it to? That is, you want it to send 512 bytes, but it returns 412. What happened to the remaining 100 bytes?

Well, they're still in your little buffer waiting to be sent out. Due to circumstances beyond your control, the kernel decided not to send all the data out in one chunk, and now, my friend, it's up to you to get the data out there.

You could write a function like this to do it, too:

```
#include <sys/types.h>
#include <sys/socket.h>

int sendall(int s, char *buf, int *len)
{
    int total = 0;          // how many bytes we've sent
    int bytesleft = *len;   // how many we have left to send
    int n;

    while(total < *len) {
        n = send(s, buf+total, bytesleft, 0);
        if (n == -1) { break; }
        total += n;
        bytesleft -= n;
    }

    *len = total; // return number actually sent here

    return n==-1?-1:0; // return -1 on failure, 0 on success
}
```

In this example, *s* is the socket you want to send the data to, *buf* is the buffer containing the data, and *len* is a pointer to an `int` containing the number of bytes in the buffer.

The function returns **-1** on error (and *errno* is still set from the call to `send()`.) Also, the number of bytes actually sent is returned in *len*. This will be the same number of bytes you asked it to send, unless there was an error. `sendall()` will do its best, huffing and puffing, to send the data out, but if there's an error, it gets back to you right away.

For completeness, here's a sample call to the function:

```
char buf[10] = "Beej!";
int len;

len = strlen(buf);
if (sendall(s, buf, &len) == -1) {
    perror("sendall");
    printf("We only sent %d bytes because of the error!\n", len);
}
```

What happens on the receiver's end when part of a packet arrives? If the packets are variable length, how does the receiver know when one packet ends and another begins? Yes, real-world scenarios are a royal pain in the donkeys. You probably have to *encapsulate* (remember that from the [data encapsulation section](#) way back there at the beginning?) Read on for details!

6.4. Son of Data Encapsulation

What does it really mean to encapsulate data, anyway? In the simplest case, it means you'll stick a header on there with either some identifying information or a packet length, or both.

What should your header look like? Well, it's just some binary data that represents whatever you feel is necessary to complete your project.

Wow. That's vague.

Okay. For instance, let's say you have a multi-user chat program that uses **SOCK_STREAMs**. When a user types ("says") something, two pieces of information need to be transmitted to the server: what was said and who said it.

So far so good? "What's the problem?" you're asking.

The problem is that the messages can be of varying lengths. One person named "tom" might say, "Hi", and another person named "Benjamin" might say, "Hey guys what is up?"

So you send () all this stuff to the clients as it comes in. Your outgoing data stream looks like this:

```
t o m H i B e n j a m i n H e y g u y s w h a t i s u p ?
```

And so on. How does the client know when one message starts and another stops? You could, if you wanted, make all messages the same length and just call the `sendall()` we implemented, [above](#). But that wastes bandwidth! We don't want to send () 1024 bytes just so "tom" can say "Hi".

So we *encapsulate* the data in a tiny header and packet structure. Both the client and server know how to pack and unpack (sometimes referred to as "marshal" and "unmarshal") this data. Don't look now, but we're starting to define a *protocol* that describes how a client and server communicate!

In this case, let's assume the user name is a fixed length of 8 characters, padded with `'\0'`. And then let's assume the data is variable length, up to a maximum of 128 characters. Let's have a look at a sample packet structure that we might use in this situation:

1. `len` (1 byte, unsigned) -- The total length of the packet, counting the 8-byte user name and chat data.
2. `name` (8 bytes) -- The user's name, NUL-padded if necessary.
3. `chatdata` (*n*-bytes) -- The data itself, no more than 128 bytes. The length of the packet should be calculated as the length of this data plus 8 (the length of the name field, above).

Why did I choose the 8-byte and 128-byte limits for the fields? I pulled them out of the air, assuming they'd be long enough. Maybe, though, 8 bytes is too restrictive for your needs, and you can have a 30-byte name field, or whatever. The choice is up to you.

Using the above packet definition, the first packet would consist of the following information (in hex and ASCII):

```

    0A      74 6F 6D 00 00 00 00 00      48 69
(length)  T o m      (padding)      H i

```

And the second is similar:

```

    14      42 65 6E 6A 61 6D 69 6E      48 65 79 20 67 75 79 73 20 77 ...
(length)  B e n j a m i n      H e y      g u y s      w ...

```

(The length is stored in Network Byte Order, of course. In this case, it's only one byte so it doesn't matter, but generally speaking you'll want all your binary integers to be stored in Network Byte Order in your packets.)

When you're sending this data, you should be safe and use a command similar to [sendall\(\)](#), above, so you know all the data is sent, even if it takes multiple calls to `send()` to get it all out.

Likewise, when you're receiving this data, you need to do a bit of extra work. To be safe, you should assume that you might receive a partial packet (like maybe we receive "14 42 65 6E" from Benjamin, above, but that's all we get in this call to `recv()`). We need to call `recv()` over and over again until the packet is completely received.

But how? Well, we know the number of bytes we need to receive in total for the packet to be complete, since that number is tacked on the front of the packet. We also know the maximum packet size is $1+8+128$, or 137 bytes (because that's how we defined the packet.)

What you can do is declare an array big enough for two packets. This is your work array where you will reconstruct packets as they arrive.

Every time you `recv()` data, you'll feed it into the work buffer and check to see if the packet is complete. That is, the number of bytes in the buffer is greater than or equal to the length specified in the header (+1, because the length in the header doesn't include the byte for the length itself.) If the number of bytes in the buffer is less than 1, the packet is not complete, obviously. You have to make a special case for this, though, since the first byte is garbage and you can't rely on it for the correct packet length.

Once the packet is complete, you can do with it what you will. Use it, and remove it from your work buffer.

Whew! Are you juggling that in your head yet? Well, here's the second of the one-two punch: you might have read past the end of one packet and onto the next in a single `recv()` call. That is, you have a work buffer with one complete packet, and an incomplete part of the next packet! Bloody heck. (But this is why you made your work buffer large enough to hold two packets--in case this happened!)

Since you know the length of the first packet from the header, and you've been keeping track of the number of bytes in the work buffer, you can subtract and calculate how many of the bytes in the work buffer belong to the second (incomplete) packet. When you've handled the first one, you can clear it out of the work buffer and move the partial second packet down the to front of the buffer so it's all ready to go for the next `recv()`.

(Some of you readers will note that actually moving the partial second packet to the beginning of the work buffer takes time, and the program can be coded to not require this by using a circular buffer. Unfortunately for the rest of you, a discussion on circular buffers is beyond the scope of this

article. If you're still curious, grab a data structures book and go from there.)

I never said it was easy. Ok, I did say it was easy. And it is; you just need practice and pretty soon it'll come to you naturally. By Excalibur I swear it!

6.5. Broadcast Packets--Hello, World!

So far, this guide has talked about sending data from one host to one other host. But it is possible, I insist, that you can, with the proper authority, send data to multiple hosts *at the same time*!

With UDP (only UDP, not TCP) and standard IPv4, this is done through a mechanism called *broadcasting*. With IPv6 (not appearing in this guide...yet), broadcasting isn't supported, and you have to resort to the often superior technique of *multicasting*. But enough of the starry-eyed future--we're stuck in the 32-bit present.

But wait! You can't just run off and start broadcasting willy-nilly; You have to set the socket option **SO_BROADCAST** before you can send a broadcast packet out on the network. It's like a one of those little plastic covers they put over the missile launch switch! That's just how much power you hold in your hands!

But seriously, though, there is a danger to using broadcast packets, and that is: every system that receives a broadcast packet must undo all the onion-skin layers of data encapsulation until it finds out what port the data is destined to. And then it hands the data over or discards it. In either case, it's a lot of work for each machine that receives the broadcast packet, and since it is all of them on the local network, that could be a lot of machines doing a lot of unnecessary work. When the game Doom first came out, this was a complaint about its network code.

Yes, I said the local network. There is more than one way to skin a cat... wait a minute. Is there really more than one way to skin a cat? What kind of expression is that? Uh, and likewise, there is more than one way to send a broadcast packet, but the broadcast packets will usually be restricted to your local network no matter how you send them.

So now to the meat and potatoes of the whole thing: how do you specify the destination address for a broadcast message? There are two common ways.

1. Send the data to your broadcast address. This is your network number with all one-bits set for the host portion of the address. For instance, at home my network is 192.168.1.0, my netmask is 255.255.255.0, so the last byte of the address is my host number (because the first three bytes, according to the netmask, are the network number). So my broadcast address is 192.168.1.255. Under Unix, the **ifconfig** command will actually give you all this data. (If you're curious, the bitwise logic to get your broadcast address is **network_number** OR (NOT **netmask**).)
2. Send the data to the "global" broadcast address. This is 255.255.255.255, aka **INADDR_BROADCAST**. Many machines will automatically bitwise AND this with your network number to convert it to a network broadcast address, but some won't. It varies.

So what happens if you try to send data on the broadcast address without first setting the **SO_BROADCAST** socket option? Well, let's fire up good old [talker and listener](#) and see what happens.

```
$ talker 192.168.1.2 foo
sent 3 bytes to 192.168.1.2
$ talker 192.168.1.255 foo
sendto: Permission denied
$ talker 255.255.255.255 foo
sendto: Permission denied
```

Yes, it's not happy at all...because we didn't set the **SO_BROADCAST** socket option. Do that, and now you can `sendto()` anywhere you want!

In fact, that's the *only difference* between a UDP application that can broadcast and one that can't. So let's take the old **talker** application and add one section that sets the **SO_BROADCAST** socket option. We'll call this program [broadcaster.c](#):

```
/*
** broadcaster.c -- a datagram "client" like talker.c, except
**                  this one can broadcast
*/

#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>
#include <errno.h>
#include <string.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>

#define SERVERPORT 4950    // the port users will be connecting to

int main(int argc, char *argv[])
{
    int sockfd;
    struct sockaddr_in their_addr; // connector's address information
    struct hostent *he;
    int numbytes;
    int broadcast = 1;
    //char broadcast = '1'; // if that doesn't work, try this

    if (argc != 3) {
        fprintf(stderr, "usage: broadcaster hostname message\n");
        exit(1);
    }

    if ((he=gethostbyname(argv[1])) == NULL) { // get the host info
        perror("gethostbyname");
```

```

        exit(1);
    }

    if ((sockfd = socket(AF_INET, SOCK_DGRAM, 0)) == -1) {
        perror("socket");
        exit(1);
    }

    // this call is the difference between this program and talker.c:
    if (setsockopt(sockfd, SOL_SOCKET, SO_BROADCAST, &broadcast,
        sizeof(broadcast)) == -1) {
        perror("setsockopt (SO_BROADCAST)");
        exit(1);
    }

    their_addr.sin_family = AF_INET;        // host byte order
    their_addr.sin_port = htons(SERVERPORT); // short, network byte order
    their_addr.sin_addr = *((struct in_addr *)he->h_addr);
    memset(&(their_addr.sin_zero), '\0', 8); // zero the rest of the struct

    if ((numbytes=sendto(sockfd, argv[2], strlen(argv[2]), 0,
        (struct sockaddr *)&their_addr, sizeof(struct sockaddr))) == -1) {
        perror("sendto");
        exit(1);
    }

    printf("sent %d bytes to %s\n", numbytes, inet_ntoa(their_addr.sin_addr));

    close(sockfd);

    return 0;
}

```

What's different between this and a "normal" UDP client/server situation? Nothing! (With the exception of the client being allowed to send broadcast packets in this case.) As such, go ahead and run the old UDP [listener](#) program in one window, and **broadcaster** in another. You should be now be able to do all those sends that failed, above.

```

$ talker 192.168.1.2 foo
sent 3 bytes to 192.168.1.2
$ talker 192.168.1.255 foo
sent 3 bytes to 192.168.1.255
$ talker 255.255.255.255 foo
sent 3 bytes to 255.255.255.255

```

And you should see **listener** responding that it got the packets.

Well, that's kind of exciting. But now fire up **listener** on another machine next to you on the same network so that you have two copies going, one on each machine, and run **broadcaster** again with your broadcast address... Hey! Both **listeners** get the packet even though you only called `sendto()` once! Cool!

If the **listener** gets data you send directly to it, but not data on the broadcast address, it could be that you have a firewall on your local machine that is blocking the packets. (Yes, Pat and Bapper, thank you for realizing before I did that this is why my sample code wasn't working. I told you I'd mention you in the guide, and here you are. So *nyah*.)

Again, be careful with broadcast packets. Since every machine on the LAN will be forced to deal with the packet whether it `recvfrom()`s it or not, it can present quite a load to the entire computing network. They are definitely to be used sparingly and appropriately.

[Prev](#)[Contents](#)[Next](#)