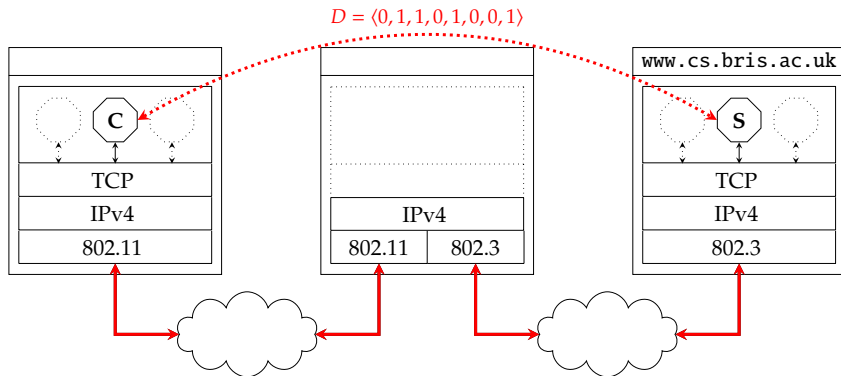


- ▶ **Goal:** *finally* investigate the **application layer**, e.g.,
 - ▶ the (mainly kernel-based) network stack implementation,
 - ▶ the interface between application and network stack, i.e.,
 1. a **raw socket** [3], or
 2. the **POSIX sockets API**,
 and
 - ▶ examples of how *you* can use all this!



$$F = H_{802.11} \parallel H_{IPv4} \parallel H_{TCP} \parallel \langle 0, 1, 1, 0, 1, 0, 0, 1 \rangle \parallel T_{802.11}$$

- **Goal:** *finally* investigate the **application layer**, e.g.,
 - the (mainly kernel-based) network stack implementation,
 - the interface between application and network stack, i.e.,
 1. a **raw socket** [3], or
 2. the **POSIX sockets API**,
 - and
 - examples of how *you* can use all this!

POSIX sockets API (1) – The Interface

| | Function | Description | Blocking? |
|-------|-------------------|--|-----------|
| | socket | Form the data structure used to describe communication end-point | × |
| | bind | Associate socket data structure with (local) address | × |
| | close | Close socket and stop using it | × |
| | shutdown | Close socket and stop using it, with control over how | × |
| | getsockopt | Get or set options for a socket, i.e., control how it functions | × |
| | setsockopt | | |
| UDP { | sendto | Transmit a datagram to (remote) address | ✓ |
| | recvfrom | Receive a datagram from (remote) address | ✓ |
| | listen | Mark socket as passive, i.e., for incoming connections | × |
| TCP { | accept | Wait for a connection to be established | ✓ |
| | connect | Actively establish a connection with (remote) address | ✓ |
| | send | Transmit a segment via connection | ✓ |
| | recv | Receive a segment via connection | ✓ |
| | select | Wait for activity that would allow non-blocking access | ✓ |
| | poll | | |

POSIX sockets API (1) – The Interface

| Function | Description |
|--------------------|--|
| getnameinfo | Convert an internal, machine-readable data structure into a host name |
| getaddrinfo | Convert a host name into an internal, machine-readable data structure |
| inet_aton | Convert a dotted-decimal address into a binary, machine-readable address |
| inet_ntoa | Convert a binary, machine-readable address into a dotted-decimal address |
| htons/htons | Convert a 16/32-bit host order integer into network order |
| ntohs/ntohs | Convert a 16/32-bit network order integer into host order |

POSIX sockets API (2) – An Implementation (in Linux)

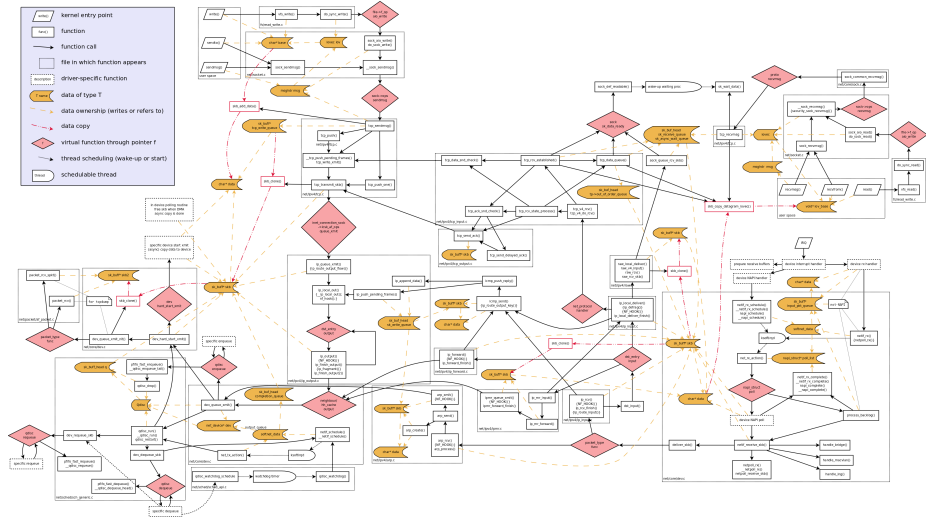
► Some (rough) **design goals** might include

1. offer POSIX-compliant interface,
2. offer RFC-compliant implementation (cf. **Postel's Law**),
3. maximise efficiency (e.g., low-latency, effective use of bandwidth),
4. maximise flexibility (e.g., general- not special-purpose),
5. allow configurability,
6. ...

which lead to some underlying **golden rules**, e.g.,

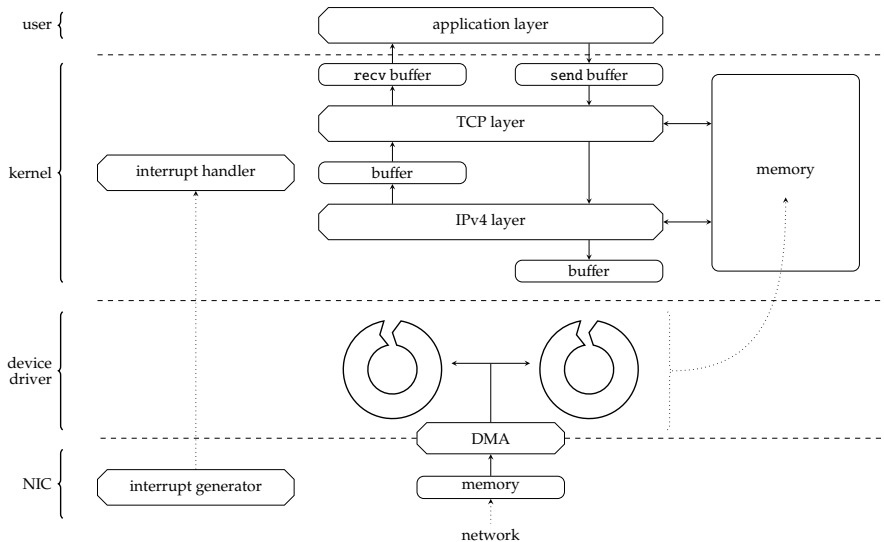
- make use of all possible hardware support,
- make use of effective data structures,
- minimise copying,
- optimise for common-case,
- ensure correctness for corner-cases,
- ...

POSIX sockets API (3) – An Implementation (in Linux)



http://www.linuxfoundation.org/collaborate/workgroups/networking/kernel_flow

POSIX sockets API (4) – An Implementation (in Linux)



POSIX sockets API (5) – An Implementation (in Linux)

- ▶ **Fact:** ports are basically buffers within network stack.
- ▶ **Implication #1:**
 - ▶ packets and segments might be received out-of-order, *but*
 - ▶ buffering enforces in-order delivery to the application.

POSIX sockets API (5) – An Implementation (in Linux)

- ▶ **Fact:** ports are basically buffers within network stack.
- ▶ **Implication #2:** send and transmission are decoupled ...
- ▶ ... transmission *could* occur
 1. when a complete segment is accumulated, or
 2. when transmission is forced, e.g., via
 - ▶ use of the PSH flag, or
 - ▶ some sort of time-out timer

so basically needs to realise a trade-off:

- ▶ less efficient wrt. latency (wait more time) but more efficient wrt. bandwidth (transmit complete segments more often), or
- ▶ more efficient wrt. latency (wait less time) but less efficient wrt. bandwidth (transmit complete segments less often).

POSIX sockets API (5) – An Implementation (in Linux)

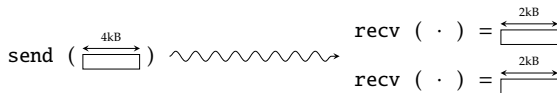
- ▶ **Fact:** ports are basically buffers within network stack.
- ▶ **Implication #3:** send and recv are decoupled ...
- ▶ ... any one of



is possible.

POSIX sockets API (5) – An Implementation (in Linux)

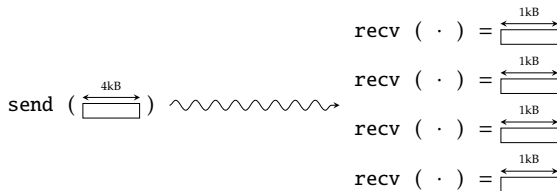
- ▶ **Fact:** ports are basically buffers within network stack.
- ▶ **Implication #3:** send and recv are decoupled ...
- ▶ ... any one of



is possible.

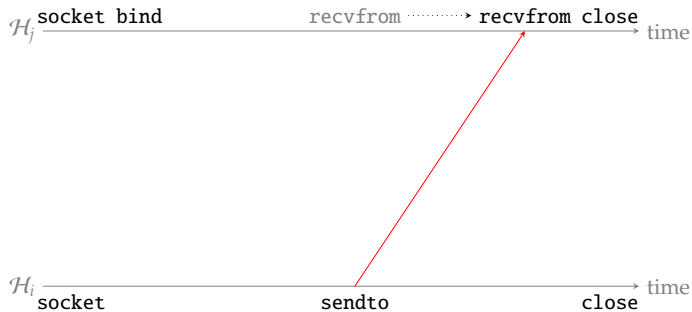
POSIX sockets API (5) – An Implementation (in Linux)

- **Fact:** ports are basically buffers within network stack.
- **Implication #3:** send and recv are decoupled ...
- ... any one of

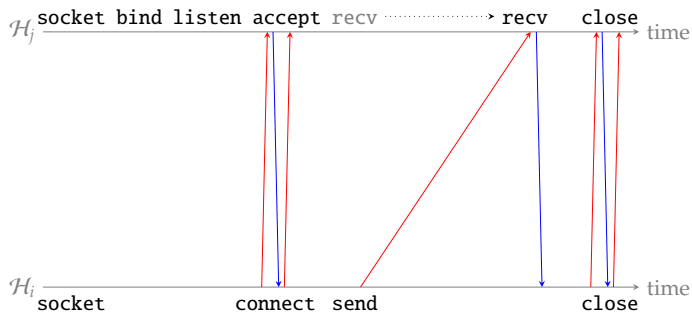


is possible.

Using POSIX sockets (1) – UDP



Using POSIX sockets (2) – TCP



Using POSIX sockets (3) – An “echo uppercase” TCP client/server

Listing

```
1 #include <sys/socket.h>
2 #include <arpa/inet.h>
3 #include <unistd.h>
4
5 void handle( int cs ) {
6     char t[ 1024 ];
7
8     while( true ) {
9         // terminal -> t
10        fgets( t, 1024, stdin );
11        // server  <- t
12        send( cs, t, strlen( t ), 0 );
13        // server  -> t'
14        t[ recv( cs, t, 1023, 0 ) ] = '\0';
15        // terminal <- t'
16        fputs( t, stdout );
17    }
18
19    // close connection
20    close( cs );
21 }
```

Listing

```
1 int main( int argc, char* argv[] ) {
2     struct sockaddr_in sa; socklen_t sl = sizeof( sa );
3     struct sockaddr_in ca; socklen_t cl = sizeof( ca );
4
5     memset( &sa, 0, sl );
6
7     sa.sin_family      = AF_INET;
8     sa.sin_addr.s_addr = inet_addr( argv[ 1 ] );
9     sa.sin_port        = htons( atoi( argv[ 2 ] ) );
10
11    // open  socket
12    int cs = socket( AF_INET, SOCK_STREAM, 0 );
13    // open  connection
14    connect( cs, ( struct sockaddr* )( &sa ), sl );
15    // handle connection
16    handle( cs );
17
18    return 0;
19 }
```

Using POSIX sockets (3) – An “echo uppercase” TCP client/server

Listing

```
1 #include <sys/socket.h>
2 #include <arpa/inet.h>
3 #include <unistd.h>
4
5 void handle( int cs ) {
6     char t[ 1024 ];
7
8     while( true ) {
9         // client -> t
10        t[ recv( cs, t, 1023, 0 ) ] = '\0';
11        // t' = toupper( t )
12        for( int i = 0; i < strlen( t ); i++ ) {
13            t[ i ] = toupper( t[ i ] );
14        }
15        // client <- t'
16        send( cs, t, strlen( t ), 0 );
17    }
18
19    // close connection
20    close( cs );
21
22 }
23
24 }
25 }
```

Listing

```
1 int main( int argc, char* argv[] ) {
2     struct sockaddr_in sa; socklen_t sl = sizeof( sa );
3     struct sockaddr_in ca; socklen_t cl = sizeof( ca );
4
5     memset( &sa, 0, sl );
6
7     sa.sin_family      = AF_INET;
8     sa.sin_addr.s_addr = inet_addr( argv[ 1 ] );
9     sa.sin_port        = htons( atoi( argv[ 2 ] ) );
10
11    // open socket
12    int ss = socket( AF_INET, SOCK_STREAM, IPPROTO_IP );
13    // bind socket
14    bind( ss, ( struct sockaddr* )( &sa ), sl );
15    // listen for connections
16    listen( ss, 10 );
17
18    while( true ) {
19
20        // open connection
21        int cs = accept( ss, &ca, &cl );
22        // handle connection
23        handle( cs );
24    }
25
26    // close socket
27    close( ss );
28
29    return 0;
30
31 }
```


Using POSIX sockets (3) – An “echo uppercase” TCP client/server

Listing

```
1 #include <sys/socket.h>
2 #include <arpa/inet.h>
3 #include <unistd.h>
4
5 void* handle( void* __cs ) {
6     char t[ 1024 ];
7
8     int cs = *( int* )( __cs );
9
10    while( true ) {
11        // client -> t
12        t[ recv( cs, t, 1023, 0 ) ] = '\0';
13        // t' = toupper( t )
14        for( int i = 0; i < strlen( t ); i++ ) {
15            t[ i ] = toupper( t[ i ] );
16        }
17        // client <- t'
18        send( cs, t, strlen( t ), 0 );
19    }
20
21    // close connection
22    close( cs );
23
24    return NULL;
25 }
```

Listing

```
1 int main( int argc, char* argv[] ) {
2     struct sockaddr_in sa; socklen_t sl = sizeof( sa );
3     struct sockaddr_in ca; socklen_t cl = sizeof( ca );
4
5     memset( &sa, 0, sl );
6
7     sa.sin_family      = AF_INET;
8     sa.sin_addr.s_addr = inet_addr( argv[ 1 ] );
9     sa.sin_port        = htons( atoi( argv[ 2 ] ) );
10
11    // open socket
12    int ss = socket( AF_INET, SOCK_STREAM, IPPROTO_IP );
13    // bind socket
14    bind( ss, ( struct sockaddr* )( &sa ), sl );
15    // listen for connections
16    listen( ss, 10 );
17
18    while( true ) {
19        pthread_t id;
20
21        // open connection
22        int cs = accept( ss, &ca, &cl );
23        // handle connection
24        pthread_create( &id, NULL, &handle, &cs );
25    }
26
27    // close socket
28    close( ss );
29
30    return 0;
31 }
```

► Take away points:

- Ultimately, the POSIX sockets API is an abstraction of the network ...
- ... even so, it's hard to argue you can totally avoid having to understand the underlying technology.
- As with any design, it has **good** and **bad** features: for example,
 - it offers a uniform interface to analogous concepts (cf. **domain sockets** [2], for IPC),
 - it allows special-case implementation choices such as use of **TCP offload**,
 - the abstraction offered is still low-level so can be hard to use (directly),
 - numerous requirements have changed over time (e.g., network vs. host order, new protocols, new use-cases), but the API hasn't,
 - ...

Additional Reading

- ▶ *Wikipedia: Raw socket*. URL: http://en.wikipedia.org/wiki/Raw_socket.
- ▶ *Wikipedia: Berkeley sockets*. URL: http://en.wikipedia.org/wiki/Berkeley_sockets.
- ▶ *Wikipedia: Winsock*. URL: <http://en.wikipedia.org/wiki/Winsock>.
- ▶ *Wikipedia: Domain socket*. URL: http://en.wikipedia.org/wiki/Unix_domain_socket.
- ▶ *Wikipedia: TCP offload engine*. URL: http://en.wikipedia.org/wiki/TCP_offload_engine.
- ▶ W.R. Stevens, B. Fenner, and A.M. Rudoff. *UNIX Network Programming Volume 1: The Sockets Networking API*. . 3rd ed. Addison Wesley, 2003.
- ▶ *Standard for Information Technology - Portable Operating System Interface (POSIX)*. . Institute of Electrical and Electronics Engineers (IEEE) 1003.1-2008. 2008. URL: <http://standards.ieee.org>.

References

- [1] *Wikipedia: Berkeley sockets*. URL: http://en.wikipedia.org/wiki/Berkeley_sockets (see p. 19).
- [2] *Wikipedia: Domain socket*. URL: http://en.wikipedia.org/wiki/Unix_domain_socket (see pp. 18, 19).
- [3] *Wikipedia: Raw socket*. URL: http://en.wikipedia.org/wiki/Raw_socket (see pp. 1, 2, 19).
- [4] *Wikipedia: TCP offload engine*. URL: http://en.wikipedia.org/wiki/TCP_offload_engine (see p. 19).
- [5] *Wikipedia: Winsock*. URL: <http://en.wikipedia.org/wiki/Winsock> (see p. 19).
- [6] W.R. Stevens, B. Fenner, and A.M. Rudoff. *UNIX Network Programming Volume 1: The Sockets Networking API*. 3rd ed. Addison Wesley, 2003 (see p. 19).
- [7] *Standard for Information Technology - Portable Operating System Interface (POSIX)*. Institute of Electrical and Electronics Engineers (IEEE) 1003.1-2008. 2008. URL: <http://standards.ieee.org> (see p. 19).