# Back to the Stone Age

A basic networking stack

Team 2: Eric Moss, Sarah Strickman, Will Merges
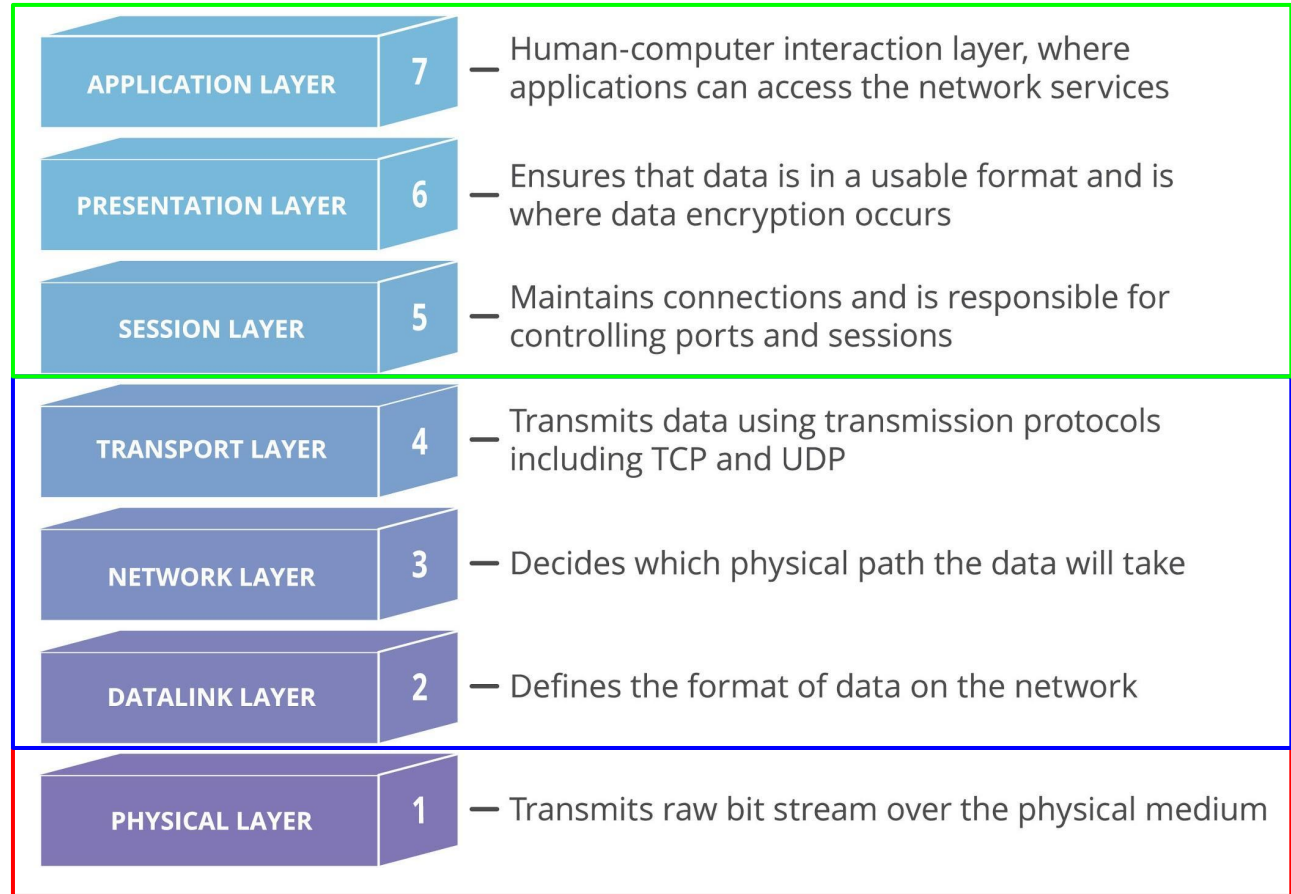
# Table of Contents

- Overview
    - Concept
    - Overall Design
- Projects
    - Messages and Syscalls
    - Packet headers
    - Ethernet Driver
    - Additional Features and Testing
- Conclusion
    - Overall difficulties
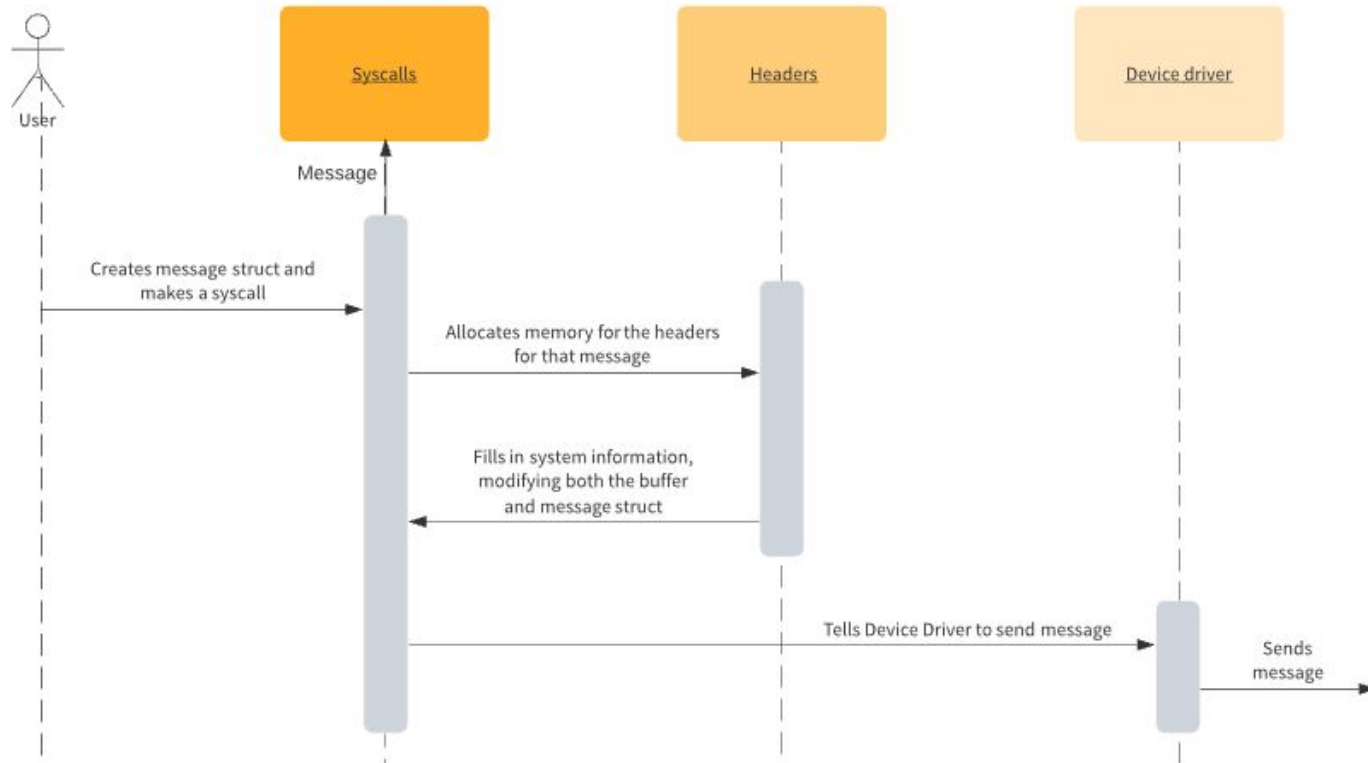    - Closing notes, what we learned

# Overview

- Project focus was Networking
  - Ethernet driver for the Intel 82557
  - Ethernet, IPv4, and UDP headers for packet construction
  - Syscalls and message abstraction to be used by user programs

- Projects were closely related
  - Led to close communication, as the projects relied on each other
  - Unique development timeline
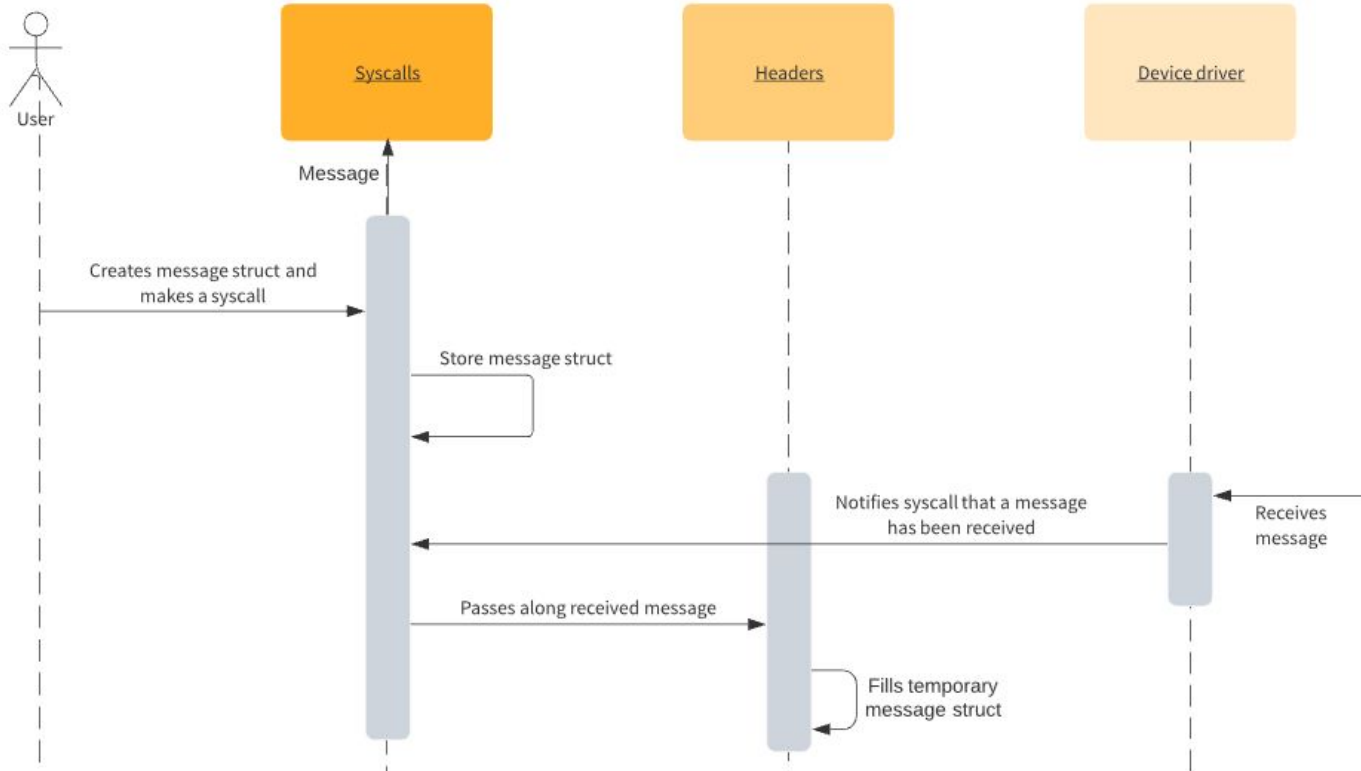    - One part had to work before the others could

# Design - OSI

- OSI model
- Device Driver / Hardware
- Header Functions
  - Ethernet II / IPv4 / UDP
- User processes / Syscalls

| | | |
|---|---|---|
| APPLICATION LAYER | 7 | Human-computer interaction layer, where applications can access the network services |
| PRESENTATION LAYER | 6 | Ensures that data is in a usable format and is where data encryption occurs |
| SESSION LAYER | 5 | Maintains connections and is responsible for controlling ports and sessions |
| TRANSPORT LAYER | 4 | Transmits data using transmission protocols including TCP and UDP |
| NETWORK LAYER | 3 | Decides which physical path the data will take |
| DATALINK LAYER | 2 | Defines the format of data on the network |
| PHYSICAL LAYER | 1 | Transmits raw bit stream over the physical medium |

# Design – Sending a message

# Design - Receiving a message

# Project 1: Messages and Syscalls

4 syscalls

- netsend
- netrecv
- setip
- setMAC

Send/Receive messages in a standardized message format

# Project 1: Message Format

**msg_t -** *net.h*

| src_port | dst_port | src_addr | |
|----------|----------|----------|---|
| dst_addr | | dst_MAC low | |
| dst_MAC high | | src_MAC low | |
| src_MAC high | | len | data* high |
| data* low | | | |

total size: 34 bytes

# Creating a message

Syscalls require some info from users, and fill in the rest themselves

- For sending:
  - Source port
  - Destination port
  - Destination address
  - Destination MAC
  - Data / Data length
- For receiving:
  - Receiving port
  - Data / Data length

# Setting MAC / IPv4 addresses

User programs can change MAC and IP addresses at runtime

- Setting MAC:
  - Take array of bytes to represent MAC
  - Extract relevant bytes
  - Set system MAC
- Setting IP:
  - User must convert address to 4 byte IP in network order (available with htons() call)
  - IP will be set as origin for sends and destination for receives
  - Can only set one system wide IP

OS will fill in this data during sends and receives

# Project 2: Packet Headers

- Two parts to this project
    - Adding headers when sending a packet
    - Handling headers when a packet is received

- Headers (both parts of the project deal with these)
    - Ethernet II Frame Header    - Link Layer
    - IPv4 Header                - Network Layer
    - UDP Header                 - Transport Layer

# Project 2: Packet Headers - Structs

**LINKhdr_t -** *link.h*

| dst_mac | | src_mac lo |
|---|---|---|
| src_mac hi | ethertype | |

total size: 14 bytes

# Project 2: Packet Headers - Structs

**NETipv4hdr_t -** *ip.h*

| ver_ihl | dscp_ecn | tot_len[2] | id | flags_offset |
|---------|----------|------------|-----------|--------------|
| ttl | protocol | checksum[2] | src_addr | |
| dst_addr | | | | |

Total size: 20 bytes

# Project 2: Packet Headers - Structs

**UDPhdr_t -** *transport.h*

| src_port | dst_port | len[2] | checksum |
|----------|----------|--------|----------|
|          |          |        |          |

total size: 8 bytes

# Sending a message

- Each layer is responsible for adding its own information
  - All functions modify the same buffer

- Only the uppermost layer (udp) copies data into buffer
  - Lower time complexity

```
__link_add_header(uint8_t* buff, uint16_t len, msg_t* msg)

    fills in information relating to the link layer
    size =

    __ipv4_add_header(uint8_t* buff, uint16_t len, msg_t* msg)

        fills in information relating to the link layer
        size =

        __udp_add_header(uint8_t* buff, uint16_t len, msg_t* msg)

            fills in information relating to the link layer
            memcpy data so it comes after the headers
            return size of udp header + size of data (or 0 if error)

        compute and fill in remaining fields (length and checksum)
        return size of ipv4 header + size (or 0 if error)

    pad payload if it's too small
    return size of link header + size (or 0 if error)
```

# Receiving a message - ipv4

```
__link_parse_frame(msg_t* msg, uint16_t len, const uint8_t* data)
    get mac address info from data (src, dst) and fill in msg with info
    return
        __ipv4_parse_frame(msg_t* msg, uint16_t len, const uint8_t* data)
            get ip address info (src, dst) from data and fill in msg with info
            return
                __udp_parse_frame(msg_t* msg, uint16_t len, const uint8_t* data)
                    get port information from data and fill in msg with info
                    point msg->data to the data in data
                    return 1 if packet if big enough (should be passed to user), 0 on error

        or 0 on error (i.e. protocol is not UDP)

    or 0 on error (i.e. ethertype is not ipv4 or arp)
```
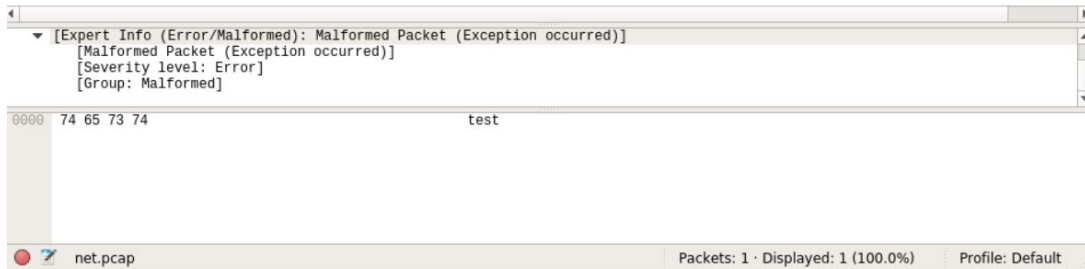
Note: We also have the ability to handle/respond to arp packets - this is additional functionality and will be talked about later on in the presentation

# Wireshark

[Expert Info (Error/Malformed): Malformed Packet (Exception occurred)]
   [Malformed Packet (Exception occurred)]
   [Severity level: Error]
   [Group: Malformed]

0000   74 65 73 74                              test

net.pcap           Packets: 1 · Displayed: 1 (100.0%)    Profile: Default

<-- Before

| No. | Time | Source | Destination | Protocol | Length | Info |
|---|---|---|---|---|---|---|
| 1 | 0.000000 | 0.0.0.0 | 1.0.0.127 | UDP | 128 | 0 → 0 Len=86 |

After -->

▸ Frame 1: 128 bytes on wire (1024 bits), 128 bytes captured (1024 bits)
▸ Ethernet II, Src: 00:00:00_00:00:00 (00:00:00:00:00:00), Dst: 00:00:00_00:00:00 (00:00:00:00:00:00)
▸ Internet Protocol Version 4, Src: 0.0.0.0, Dst: 1.0.0.127
▸ User Datagram Protocol, Src Port: 0, Dst Port: 0
▸ Data (86 bytes)

# Project 3: Ethernet Driver

- First thing we did
- The DSL lab computers have Intel 82557's, but our driver technically supports anything in the 8255x family
- Lots of technical reading

1 | of 175

*yikes*

# PCI

- The NIC is a PCI device
- Accessed PCI configuration space to obtain:
  - Device ID
  - Vendor ID
  - Interrupt line
  - Base address register
- Brute force scanned the entire PCI address space
  - Probably a better way to do this
- Had to make device a PCI master
- *pci.c / pci.h*

| Byte Offset (hexadecimal) | Byte 3 | Byte 2 | Byte 1 | Byte 0 |
|---|---|---|---|---|
| 0 | Device ID | | Vendor ID | |
| 4 | Status Register | | Command Register | |
| 8 | Class Code (200000h) | | | Revision ID |
| C | BIST | Header Type | Latency Timer | Cache Line Size |
| 10 | CSR Memory Mapped Base Address Register | | | |
| 14 | CSR I/O Mapped Base Address Register | | | |
| 18 | Flash Memory Mapped Base Address Register | | | |
| 1C | Reserved | | | |
| 20 | | | | |
| 24 | | | | |
| 28 | | | | |
| 2C | Subsystem ID | | Subsystem Vendor ID | |
| 30 | Expansion ROM Base Address Register | | | |
| 34 | Reserved | | | Cap_Ptr |
| 38 | Reserved | | | |
| 3C | Max_Latency (FFh) | Min_Grant (FFh) | Interrupt Pin (01h) | Interrupt Line |
| DC | Power Management Capabilities | | Next Item Pointer | Capability ID |
| E0 | Reserved | Data | Power Management CSR | |

# The NIC

- Broken up into Command Unit (CU) and Receive Unit (RU)
- Commands are executed using the System Control Block (SCB)
  - Load CU base - 0x00
  - Load RU base - 0x00
  - CU start
  - RU start
- CSR is accessed through the I/O address space
  - Offset is obtained from PCI BAR
- PORT commands control resets
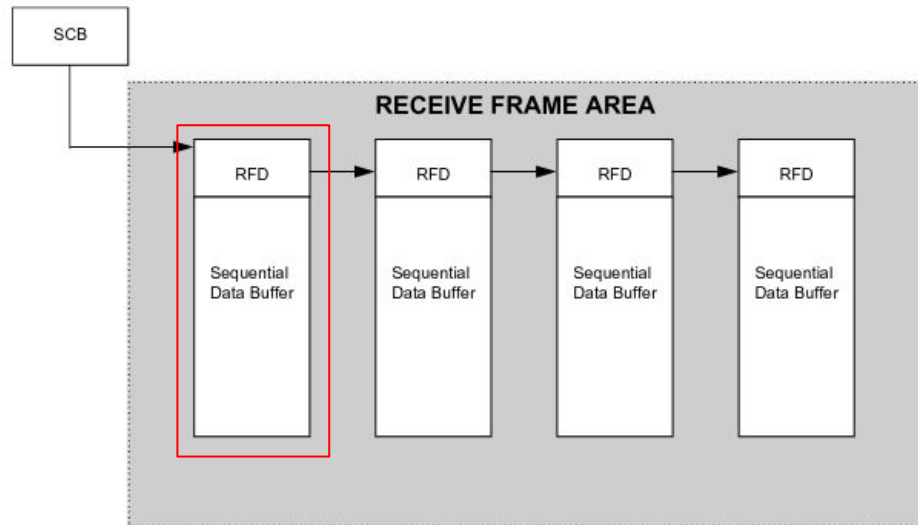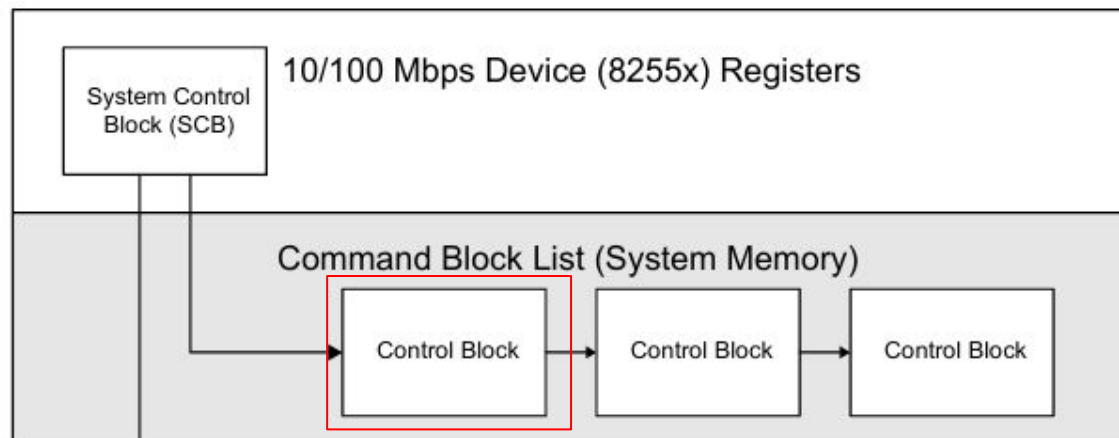- *eth.c / eth.h*

**Control / Status Register**

| Upper Word | | Lower Word | | Offset |
|---|---|---|---|---|
| 31 | | 16  15 | 0 | |
| SCB Command Word | | SCB Status Word | | 0h |
| SCB General Pointer | | | | 4h |
| PORT | | | | 8h |
| EEPROM Control Register | | Reserved | | Ch |
| MDI Control Register | | | | 10h |
| RX DMA Byte Count | | | | 14h |
| PMDR | Flow Control Register | | Reserved | 18h |
| Reserved | | General Status | General Control | 1Ch |
| Reserved | | | | 20h-2Ch |
| Function Event Register | | | | 30h |
| Function Event Mask Register | | | | 34h |
| Function Present State Register | | | | 38h |
| Force Event Register | | | | 3Ch |

# DMA – Command Block List

- A linked list of commands to execute on the CU
  - Load internal (MAC) address - makes sure packets addressed to us are not discarded
  - Transmit - send a frame
- Only ever have one command on the CBL
  - Kept our own queue of commands to be executed
  - When a command completes, change the CBL pointer and restart the CU
- Allocate command space from static memory
  - Frame data is right after the transmit command in "simple mode", variable sized buffers
  - 2 byte align everything
- Limitations
  - Not the fastest
  - Support 50 commands concurrently, up to 8192 bytes of command data

# DMA - Receive Frame Area

- RFA is made up of Receive Frame Descriptors
  - We only ever have one RFD in the RFA, makes the ISR simpler
  - 2 byte aligned
- In simple mode, receive data placed directly after RFD header
  - Flexible mode could support copying frame directly to user's memory space
- Restart the RU every "frame ready" interrupt

# Interrupts

- Device (should) interrupt when:
  - CU command completed (load address, transmit)
  - Received a frame
- Interrupt type determined by SCB status word
- Utilizes callback functions to get data from the driver
  - Commands are executed with an id
  - Pass a function pointer to ethernet driver to be called when commands are complete
  - Separate function pointer for received frame data to be passed for the OS to deal with
- Can also get "receive no resources" interrupt
  - Limitation of how the RFA is setup

# Additional Implementation

- Added ability to respond to ARP requests
  - "50% ARP"
- Handled by the receive callback function
- *arp.h / arp.c*

```
2 19.291386    52:55:0a:00:02:02    Broadcast          ARP    42 Who has 10.0.2.15? Tell 10.0.2.2
3 19.295766    0f:0f:0f:0f:0f:0f    52:55:0a:00:02:02  ARP    46 10.0.2.15 is at 0f:0f:0f:0f:0f:0f
4 19.296140    10.0.2.2            10.0.2.15          UDP    49 51740 → 8081 Len=7
```

# Difficulties

- Time and Energy
  - Hard to test
  - Some functionality of the hardware seemed to be assumed in the manual
- Never implemented TCP or ICMP
  - Always a stretch goal, UDP is much easier
- Sarah was remote
  - could not go to DSL labs - had to work on qemu
- Lots of protocols to read
  - PCI, Ethernet, IPv4, UDP
- Memory management
  - Headers and data would overwrite each other
  - Aligning data structures on allocation caused many bugs

# Closing Notes

- Happy with what we accomplished
  - Design feels clean and easy to add onto if we wished
- "Front loading" work to the beginning of the semester paid off
  - Gave us a good foundation to build off of
- Valuable experience
  - Wish we could have all been there in person, but what can you do 🤷‍♀️

# Questions?