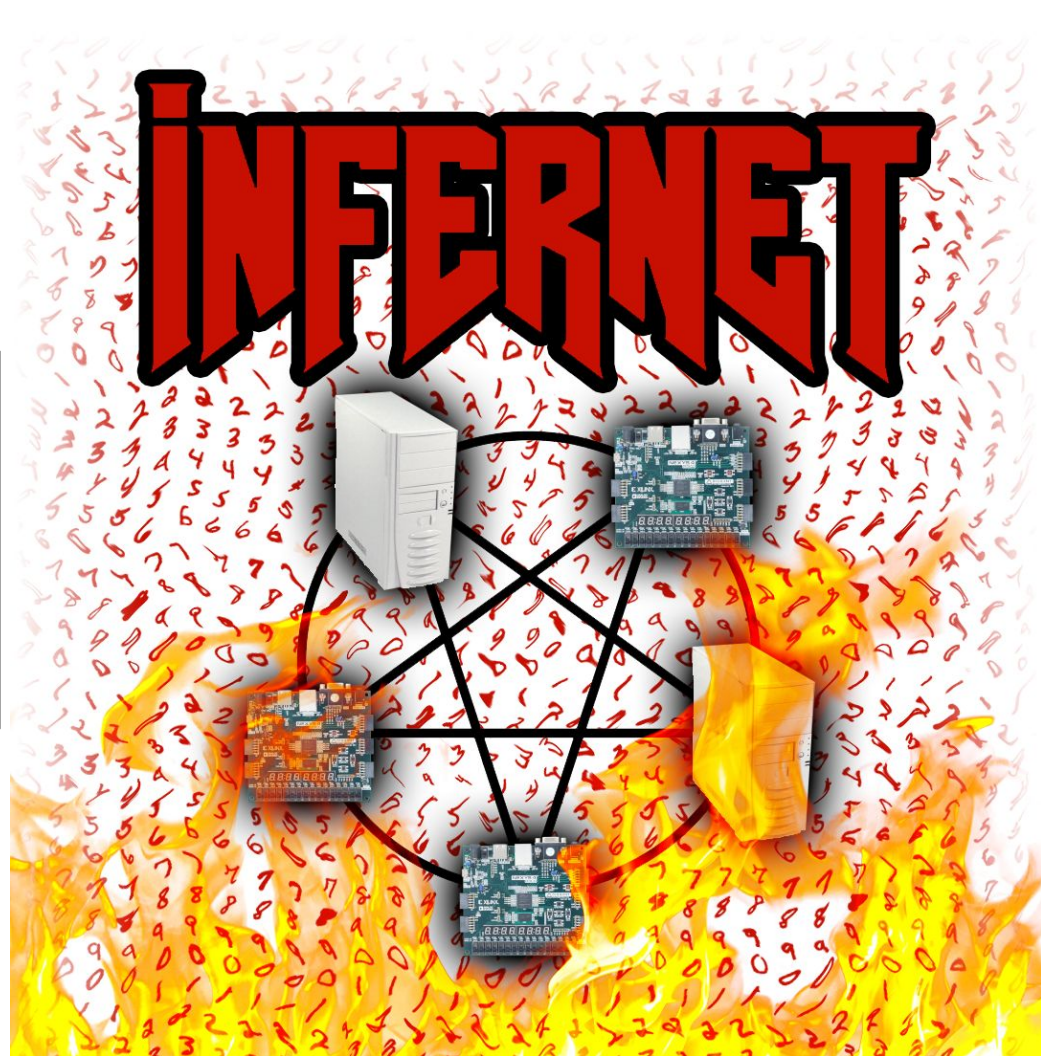


# Team 6

## ECE532 Project



# Meet the team



Andrew



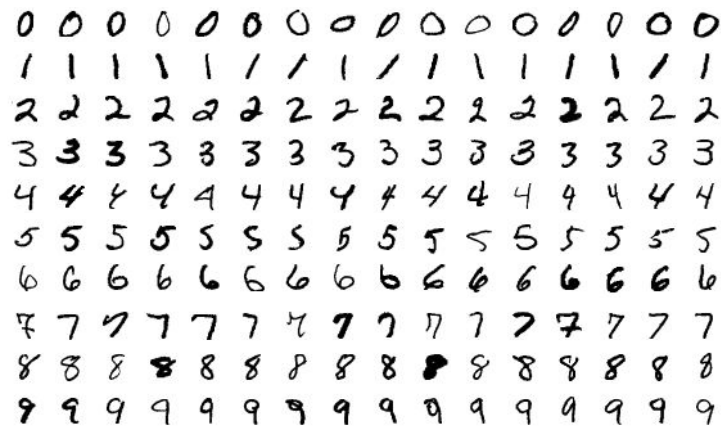
James



Quinn

# What is Infernet?

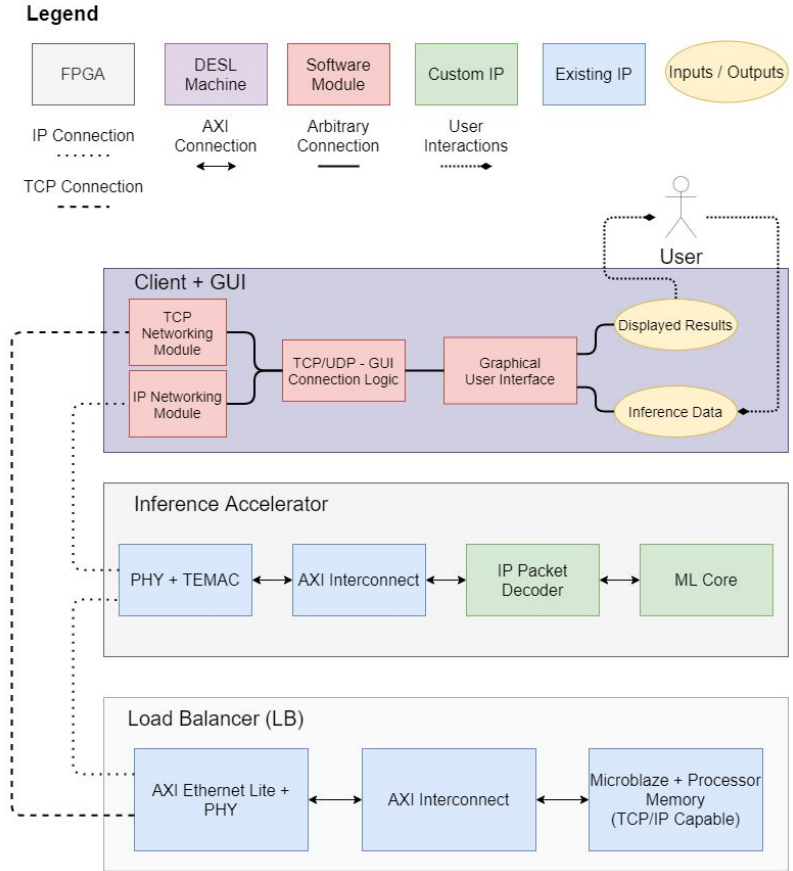
- ❑ Distributed network of machine learning inference modules targeting the MNIST dataset
- ❑ Multiple networked FPGAs with hardware inferencers available to desktop PC clients
  - ❑ For this project: FPGAs and PCs on FPGANet



*Sample images from the MNIST dataset*

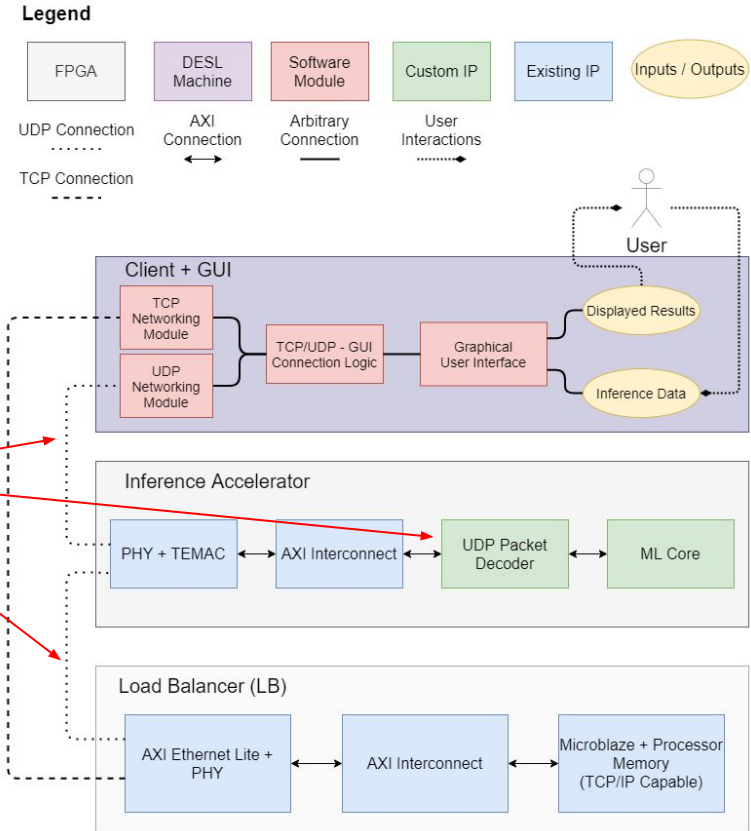
# System Overview - Old

- ❏ One or more inference accelerators (IA)
- ❏ A load balancer (LB)
- ❏ A desktop client (client)
- ❏ A desktop GUI, which will be implemented within the client to display results



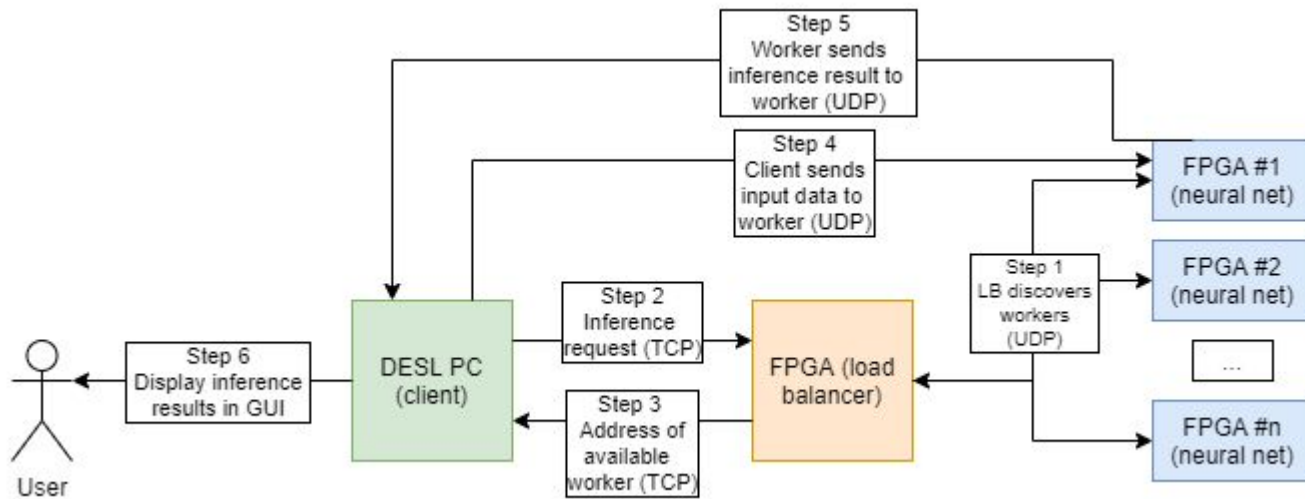
# System Overview - New

- ❏ All the same system components
- ❏ Switching from IP to UDP to circumvent RTOS issues
- ❏ Same amount of communication, slightly different protocol





# Updated Protocol



# Challenges So Far

## Client Software:

- Wrangling network interface names in windows
- No way to open a raw socket from python socket library because not admin

## Load Balancer:

- Due to Vivado bug, FreeRTOS and lwip socket API do not work together.
- No raw IP in lwip

## IP Packet Module

- Scary Integrations
- Network Order / Endianness
- Physical Interfaces

## Hardware Neural Net

- Lots of stuff to build
- Understanding DSP slices/columns and carry chains

# Future Challenges

## Client Software:

- Creating a good-looking GUI in python

## Load Balancer:

- Get UDP working with no-OS raw lwip.
- Catch and fix corner cases in state machine

## IP Packet Module

- Timing closure with big dataframes (workarounds exist)

## Hardware Neural Net

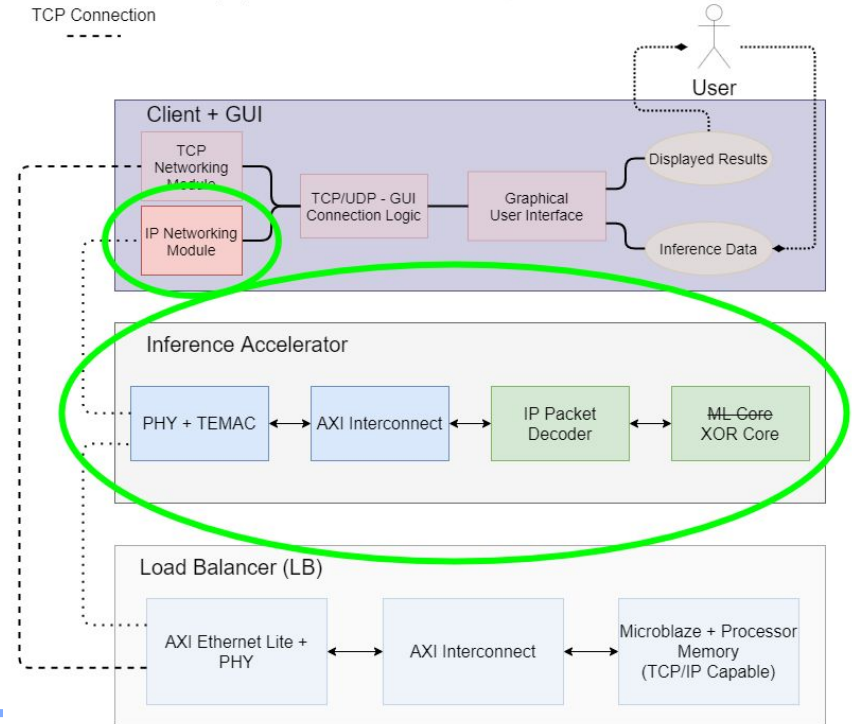
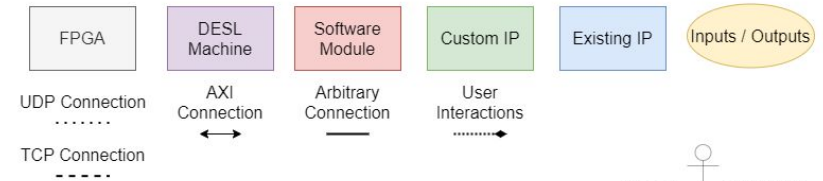
- Lots of stuff to build! :)
- Complicated control logic
- NN accuracy issues might be hard to debug
  - But also not fatal to the project



# Our Demo

- ❑ Focused on getting FPGA <=> DESL PC networking working in hardware
- ❑ Getting IP packet core from sim to H/W
- ❑ Figuring out Python IP networking on DESL PCs
- ❑ Replaced full-scale neural net with XOR neural net
- ❑ Uses same underlying dataflow arch (data width, mults)

## Legend





## What's Left

---

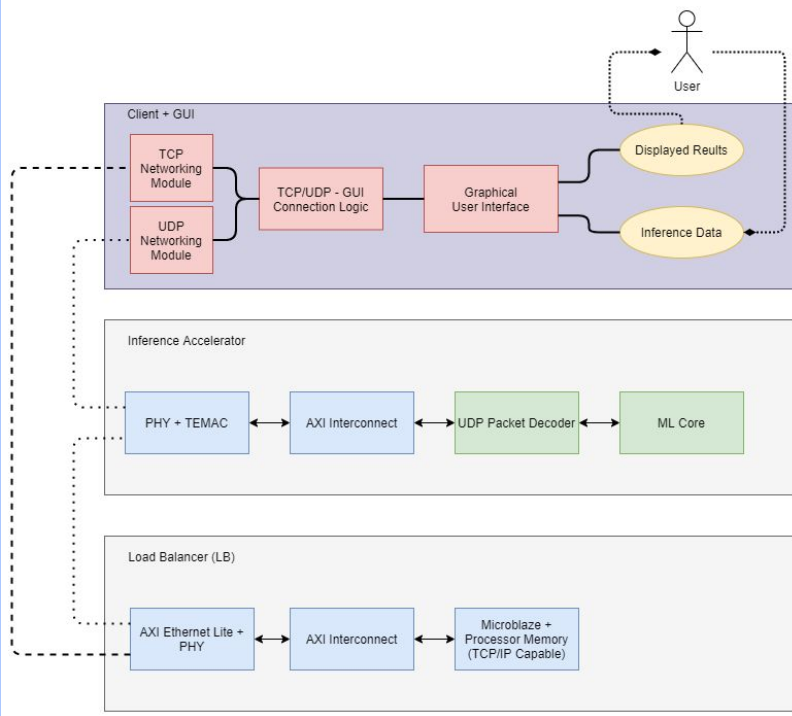
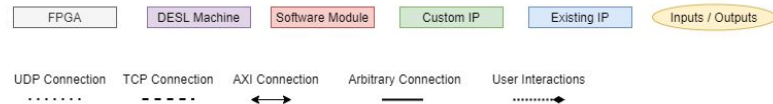
- ❑ Rebuild the LB using UDP, with the echo server as a template
- ❑ Convert IP core to UDP
- ❑ Networking validation
  - ❑ Validate LB<->IA protocol components
  - ❑ Validate LB<->Client protocol components
- ❑ Finish building the neural net in hardware
- ❑ Close timing on the neural net
- ❑ Write a pretty GUI

## Final Demo Plans

---

- Multiple DESL PCs each contacting the load balancer, acquiring a board, and streaming large data batches to the accelerators
  - We'll see how far it can scale
    - Probably will end up limited by DESL infrastructure (e.g. how many PCs can we control at once?)
  - Hopefully demonstrating good classification accuracy for the ML workload
    - And hopefully no crashes :)

## Legend



# Questions?