



FINAL DESIGN

SDMay25-01

James Minardi, Eli Ripperda, Lindsey Wessel, Mason Inman

SDMay25-01

James Minardi [CprE]

Eli Ripperda [CprE]

Lindsey Wessel [SE]

Mason Inman [SE]

JR Spidell [Client]

Dr. Zambreno



Dr. Jones



OUR CLIENT


JR SPIDELL

- ❖ Sr. Principal Systems Engineer
- ❖ This is not affiliated with his work - this is an independent project.
- ❖ Formerly volunteered to help with individuals with cerebral palsy and is motivated to help them further.
- ❖ **Wants** to develop **assistive wheelchair tech** with features including mobility assistance **and real-time seizure detection.**





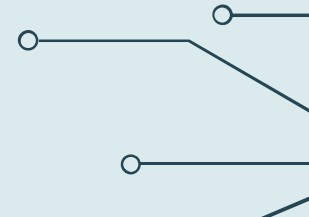
PROBLEM

- ❖ People with **mobility** and **cognitive impairments** face many challenges including maintaining **independence** and **safety**
 - ❖ Lack of advanced wheelchair technologies leads to **gaps in autonomy** , communication, etc
- 



OBJECTIVE

Develop a fast and accurate algorithm inference machine learning computer vision subsystem on an FPGA board to support our client's vision of advanced assistive technologies.



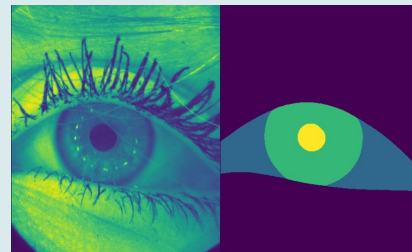
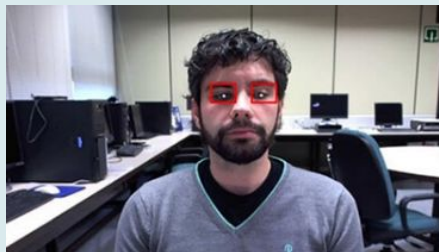
PROJECT OVERVIEW

SYSTEMS

- ❖ Camera
- ❖ Vision-Based ML algorithm
- ❖ Non-Vision-Based ML model
- ❖ Ultra96 v2 FPGA development board
- ❖ Display

REQUIREMENTS

- ❖ Real-time
- ❖ Accurate and performant to [NDA] fps
- ❖ Display model outputs and debugging information



**SDMay25-01 is under an NDA that prevents openly sharing specific metrics (eg. requirements).*

CLIENT PERSONA



INTERESTS & GOALS

- ❖ Helping people in wheelchairs
- ❖ Providing opportunities to SE & CprE Students



NEEDS

- ❖ Decrease latency in detecting health problems
- ❖ Iterative development of overall project



CHALLENGES

- ❖ Lacks abundant free time

CLIENT NEEDS



PRODUCT

- ❖ Embedded vision & non-vision ML algorithms
- ❖ Affordable System
- ❖ A fast and accurate system



PROGRESS

- ❖ Each team needs to improve upon their predecessors



COLLABORATION

- ❖ One of many teams
- ❖ Learn from previous teams
- ❖ Teach to future teams



REQUIREMENTS

REQUIREMENTS



TECHNICAL

- ❖ [NDA] FPS
- ❖ [NDA]% accuracy
- ❖ FPGA: Ultra96v2
- ❖ Tensil.ai to generate FPGA based ML accelerator



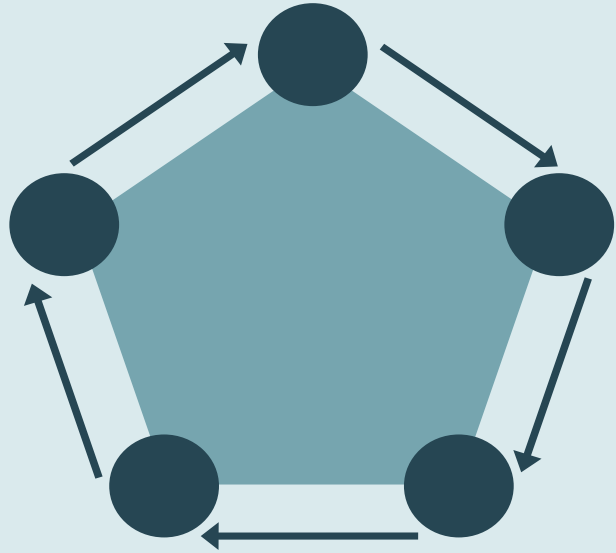
TRANSITION

- ❖ Aid in the handoff process to help future teams
- ❖ Well-documented work



CLIENT

- ❖ Low latency
- ❖ High accuracy
- ❖ Improve care of wheelchair bound individuals



PROJECT MANAGEMENT

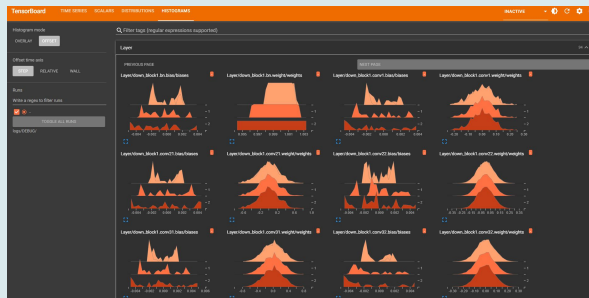
PROJECT MANAGEMENT STYLE

Agile

- ❖ Weekly client, team and advisor meetings.
- ❖ Requirements are flexible and changing per client request.
 - Concurrent work with other Sr. Design teams.



MILESTONES



Train the ML model and gather more metrics.

Optimize
Milestone 2

Real time video preprocessing and ML algorithm execution

Run Real-Time System
Milestone 4

Milestone 1

Obtain Baseline Metrics

Obtain baseline metrics for latency, accuracy, and FPS from client-provided model.

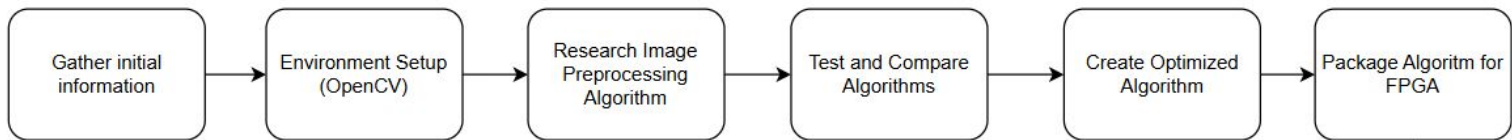
Milestone 3

Port to FPGA

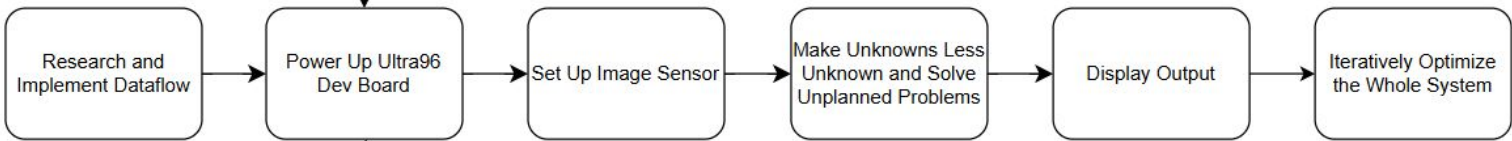
Use Tensil.ai framework to compile the model onto the FPGA.

TASK DECOMPOSITION

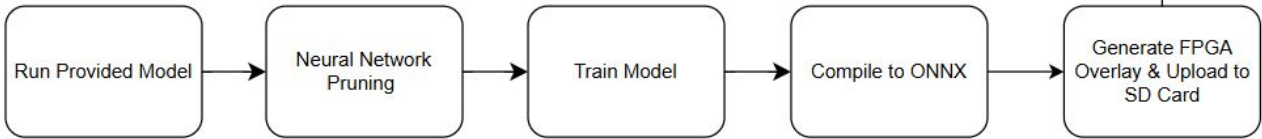
Image Preprocessing



Hardware Pipeline



ML Algorithm Optimization



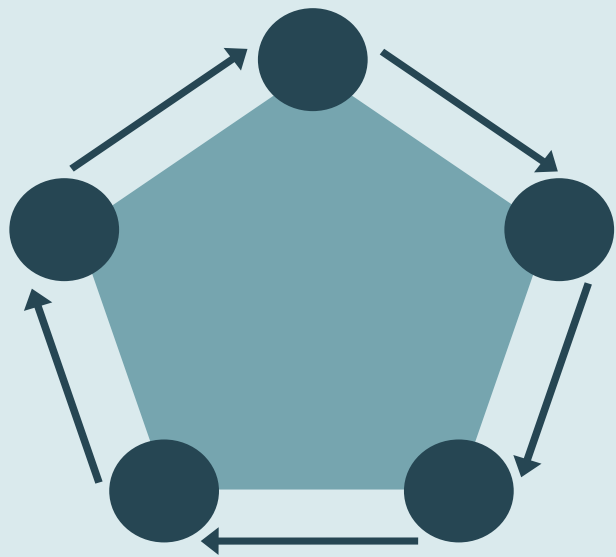
METRICS & EVALUATION CRITERIA

METRICS

- ❖ Test Latency Script (ms)
- ❖ Compare to Training Data
- ❖ Compare embedded ML prediction to ground truth
- ❖ Track FPS

EVALUATION CRITERIA

- ❖ Latency (ms)
- ❖ Accuracy (IoU values)
- ❖ Speed (FPS)



OUR DESIGN

MAJOR COMPONENTS

01 ML Algorithm Optimization

Team member: Mason

02 ML Image Processing

Team member: Lindsey

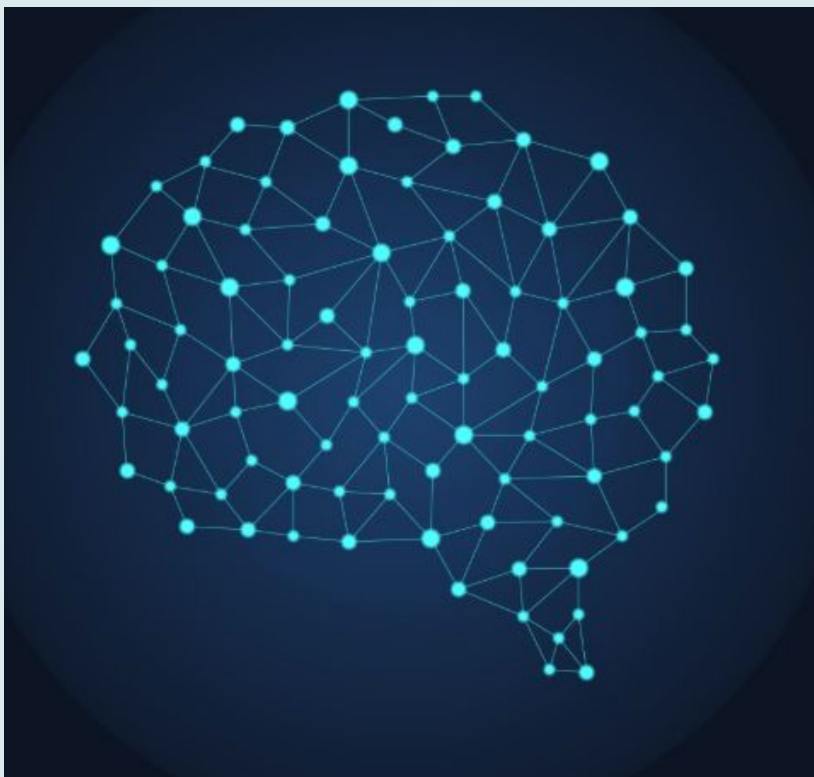
03 ML Acceleration: Tensil

Team member: Eli

04 Ultra96v2 Board

Team member: James





01

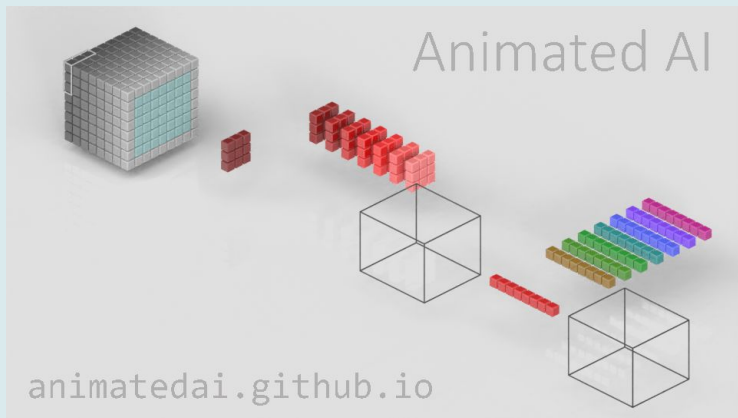
ML Algorithm Optimization



OPTIMIZATION STRATEGIES

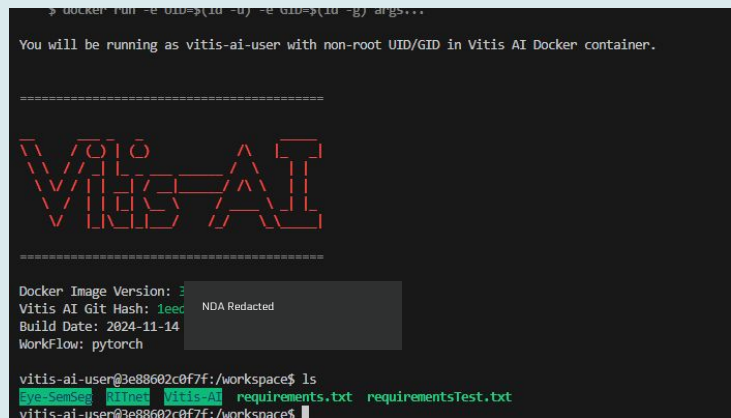
DEPTHWISE SEPARABLE CONVOLUTION

- ❖ Significant **parameter reduction** through smaller convolutions.



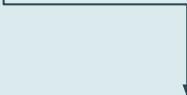
NEURAL NETWORK PRUNING

- ❖ Implement Vitis-AI tooling to use **quantization aware training**.





Test & Choose an Algorithm



* Algorithm 1 Result



* Algorithm 2 Result

02

Image PreProcessing ML Algorithm



CHOOSING AN ALGORITHM

Input

Algorithm 1

Algorithm 2

ACCURACY & TIME [s]



```
0.10858750343
0.11860680580
0.11986422538
0.10754251480
0.11414432525
0.10358309745
0.12024736404
0.11383533477
0.16153049468
0.11439013481
0.10884189605
0.10570573806
0.11348056793
0.16413927078
0.10365796089
0.10923194885
0.11400699615
0.10471320152
0.12034702301
```



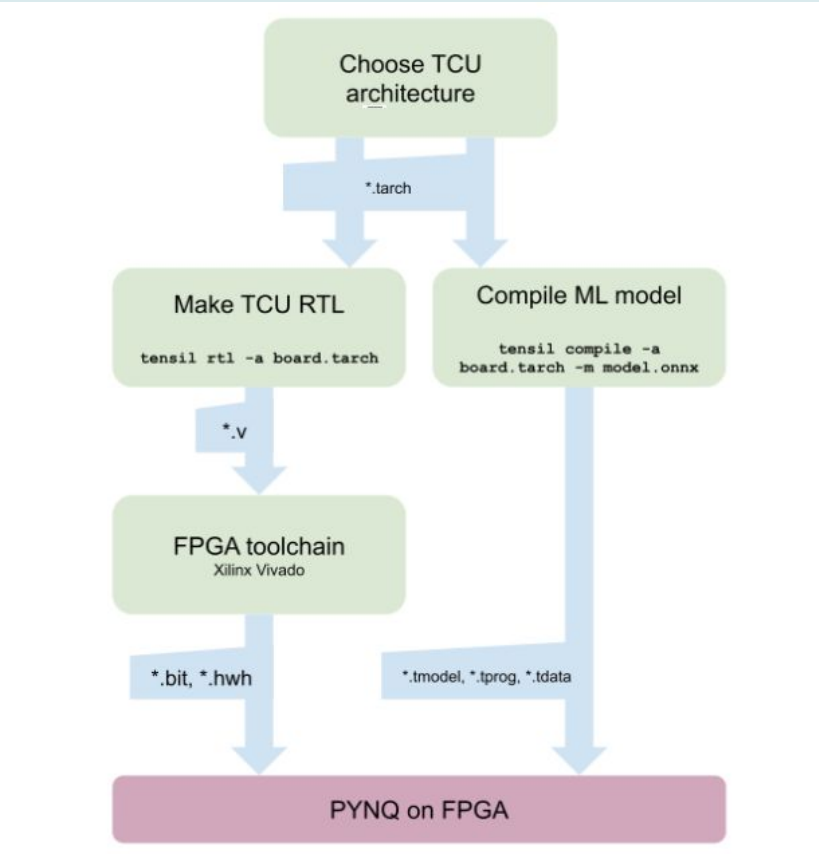
```
0.334112881
0.450552701
0.478029011
0.346270081
0.191206931
0.201151841
0.213173861
0.324286461
0.184012651
0.248700611
1.617017501
0.408772461
0.313836811
0.541065211
0.306583641
0.489648581
```



*input images to test accuracy

*Algorithm 1 is fast but lacks some accuracy

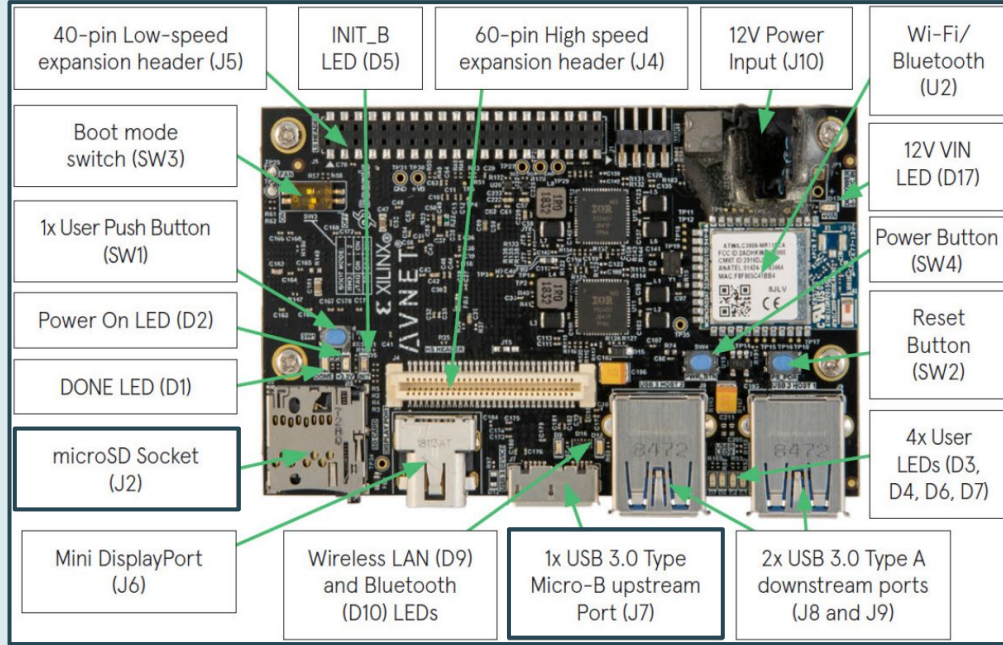
*Algorithm 2 is more accurate but slow



03

Tensil.AI



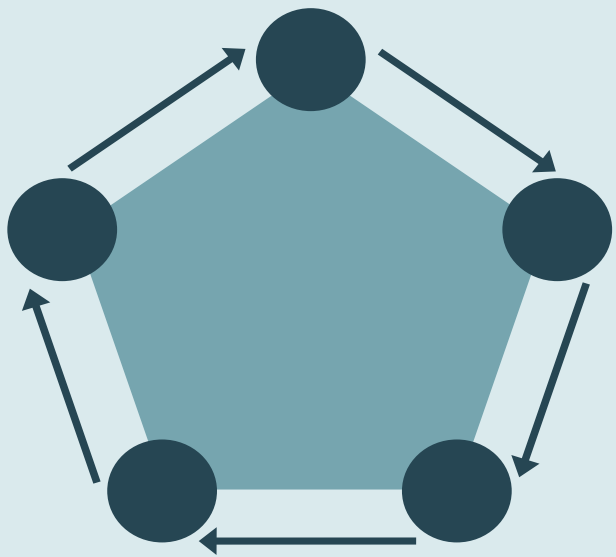


Topology

04

Ultra96v2





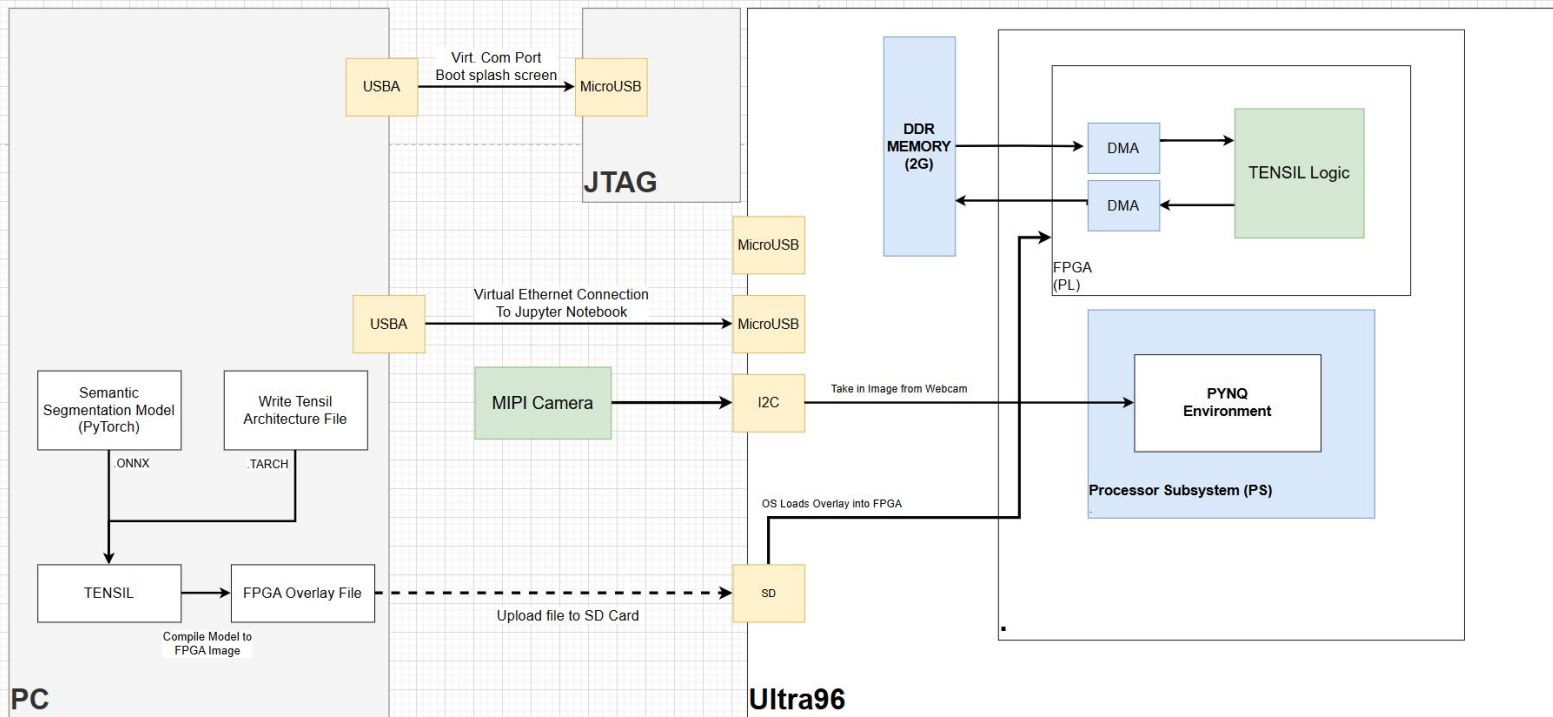
DESIGN & VISUALS

SYSTEM FLOW DIAGRAM

This represents the flow of data (an image/frame) will go through in our system.



SYSTEM BLOCK DIAGRAM



SYSTEM BLOCK DIAGRAM DESCRIPTION

ULTRA 96v2 Board

- ❖ AMD Zynq UltraScale+ processor
- ❖ 2GB DDR Memory
- ❖ PYNQ OS
 - Jupyter Notebook-based OS with Python APIs
 - Runs the Image preprocessing algorithm
- ❖ Uses AXI to send video data between the OS and model on the FPGA

THE MODEL

- ❖ Trained Model asynchronously compiled through Tensil-AI
- ❖ Model can now be used to generate and FPGA overlay for the board

I/O Devices

- ❖ MIPI Camera
- ❖ Ethernet connection via micro-USB for Jupyter notebook development
- ❖ Displayport output to monitor



Ethics

Ethical Concerns

Bias

- ❖ Will the Image Processing algorithm have unintentional gender or ethnicity bias?
- ❖ Training data must be diverse.

Liability

- ❖ Our system in the long-term may be used to determine life saving decisions.
- ❖ We must produce and be confident on the metrics we present.

Contextual Limitation

- ❖ How will the system's performance (accuracy) be impacted in non-ideal environments?
- ❖ Thoroughly test sub-systems in different environments.



CONCLUSION

As a result

Of our given problem and considerations of our project

We will

Increase the performance of an existing FPGA system

To achieve

Throughput high enough to make real-time decisions.

Linking to Our Client's Problem

This increase in data throughput will supplement our client's system, unlocking the ability to predict when end-users might have health-affecting events such as a seizure.



A decorative graphic in the top right corner consisting of several thin, dark grey lines that branch out and end in small open circles, resembling a circuit board or network diagram.

Thanks!

A decorative graphic in the bottom left corner consisting of several thin, dark grey lines that branch out and end in small open circles, resembling a circuit board or network diagram.

Any questions?