

Semantic Segmentation Optimization

By: SDDec25-01



SDDec25-01

Team Members:

Joey Metzen - CprE

Aidan Perry - CprE

Conner Ohnesorge - EE

Tyler Schaefer - SE

Advisor:

Dr. Namrata Vaswani

Client:

JR Spidell


Problem Statement



Problem

- People with disabilities face risks from undetected medical issues. Traditional methods lack real-time monitoring.
- Using eye movement tracking with semantic segmentation can detect warning signs and automatically reposition users to prevent incidents, improving safety needs.

Client

- Volunteered to help individuals with cerebral palsy.
- Create an assistive wheelchair technology. 

Team

- Update the system to increase throughput.



Resources/Data Management



Hardware/Software:

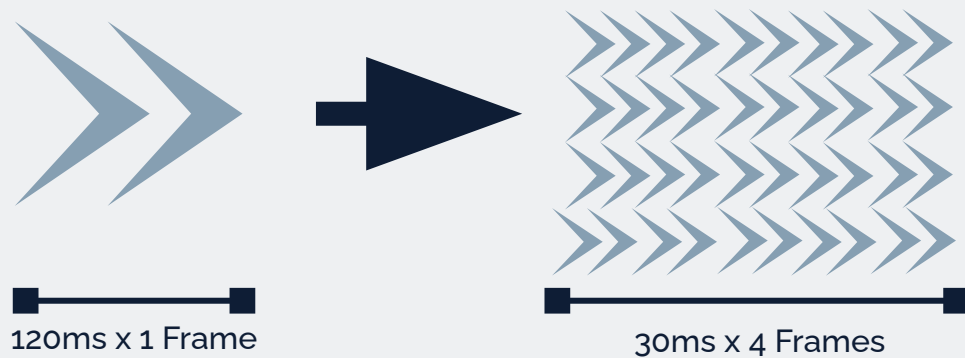
- Kria Board Kv260
- Xilinx
- Vitis-AI
- Pytorch
- ONNX & ONNX-Runtime
- Petalinux

Mem Components:

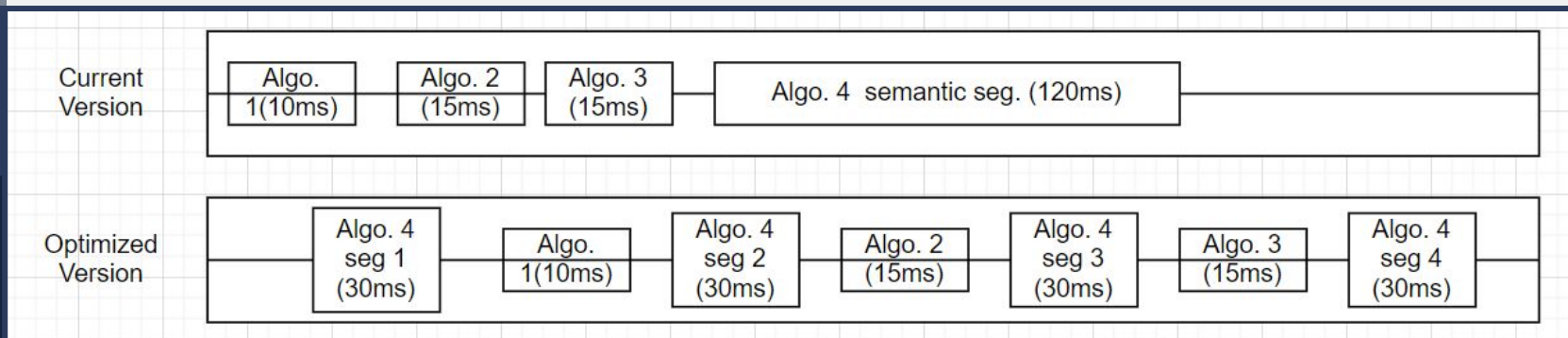
- DDR4 RAM (Main Memory) [4GB]
- Quad-SPI Flash (Boot ROM) [64MB]
- Cortex-A53 L1 Data Cache
- Cortex-A53 L2 Cache [1MB]



Project Overview



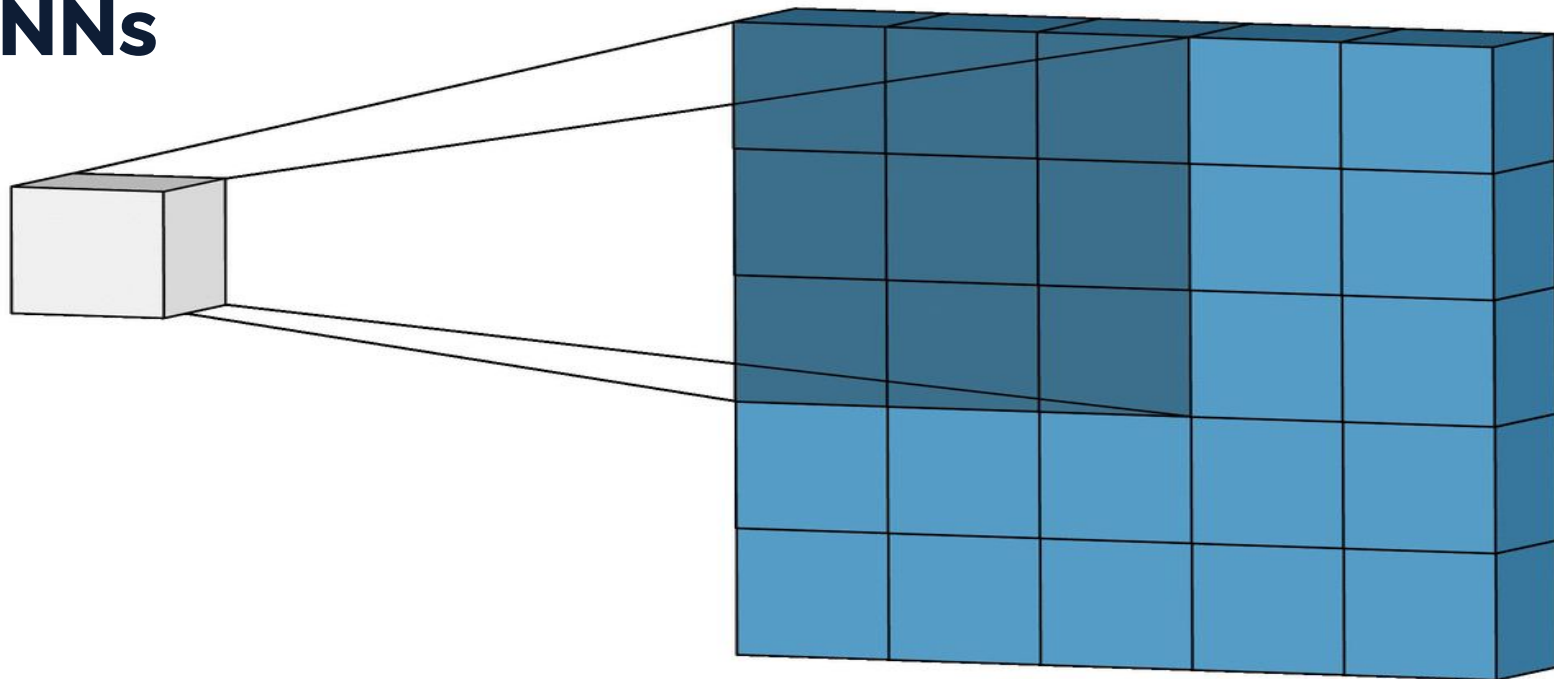
Increase throughput by splitting the U-net algorithm over 4 cores and across the DPU.



What is Semantic Segmentation



CNNs



- Convolution = sliding window
- NN basically Turns $\text{func}(3 \times 3)$ into $\text{func}(\dots \text{func}(\text{func}(3 \times 3)))$

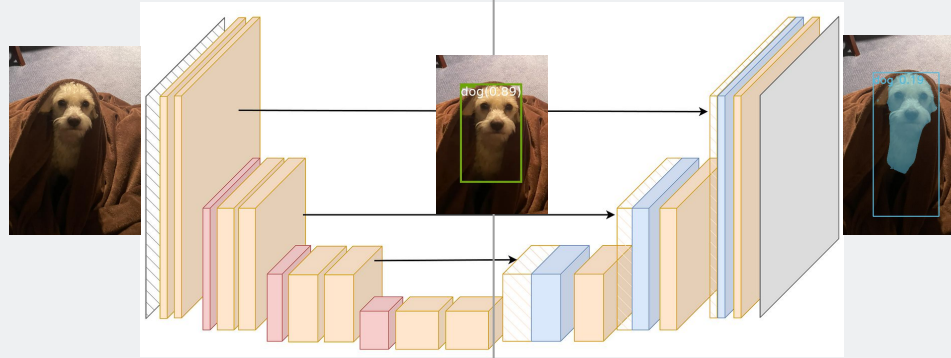
U-net Semantic Segmentation cont.

Contracting Encoder

- Downsampling (i.e. 2x2 Max Pool) compensated by the doubles # channels
- Transmit to across to decoder
Includes spatial info

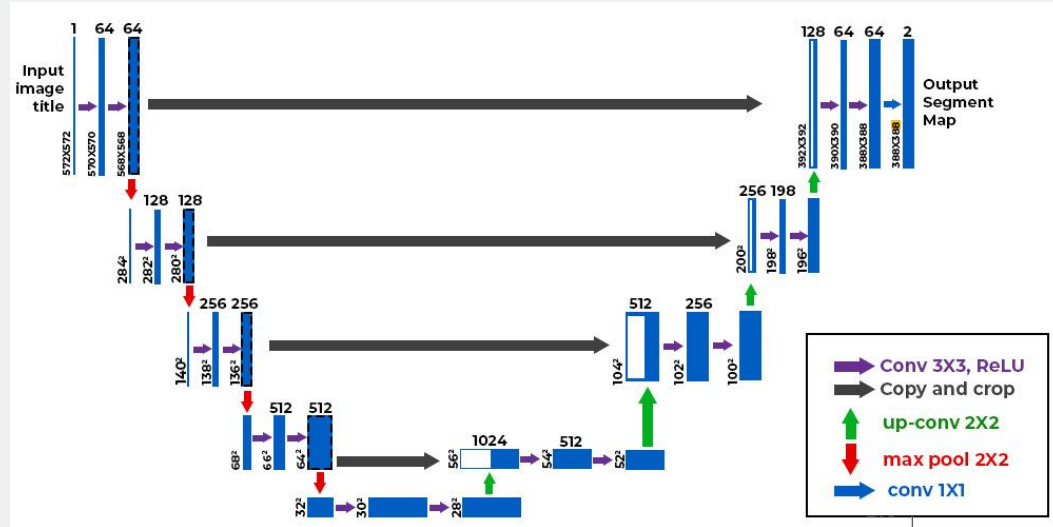
Expanding Decoder

- Receive from encoder
- Upsampling (i.e. 2x2 Convolution) decreases the amount of channels
Includes semantic Info



Each “forward step” applies a relu function to the output of a repeated convolutional layer application over input channels.

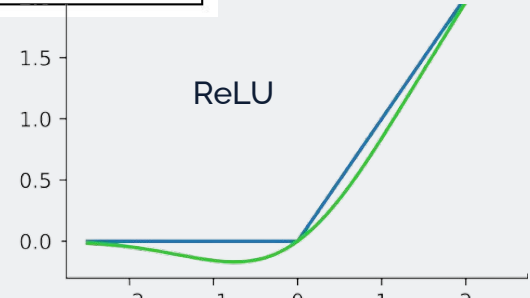
U-net Semantic Segmentation



2	2	7	3
9	4	6	1
8	5	2	4
3	1	2	6

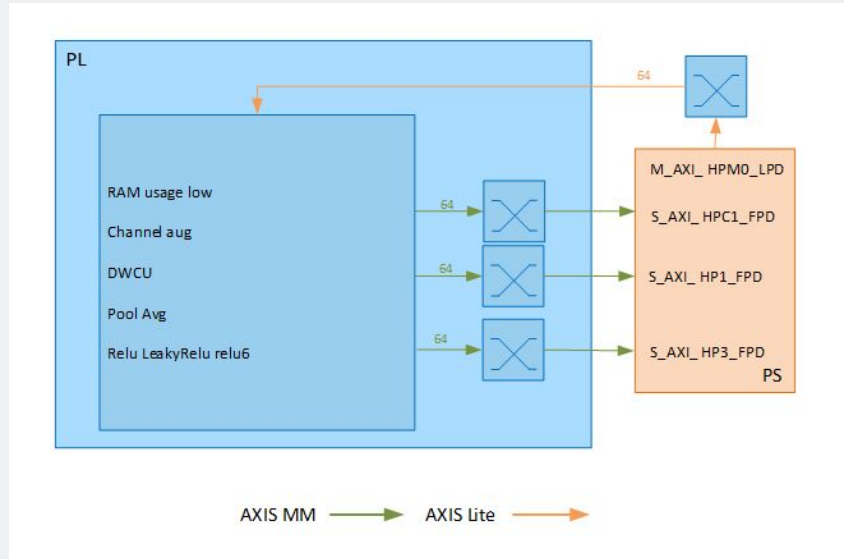
Max Pool
Filter - (2 x 2)
Stride - (2, 2)

9	7
8	6



Quantization

32 fp -> 8b int
Less precise
More processable
DPU needs integer inputs
Post Training Quantization



Data Version Control

- S3 compatible store
Versioned on top of Git
- Distributed (Supports cross country cooperation)
- Low Learning Curve



Questions before we continue?

Why Unet?

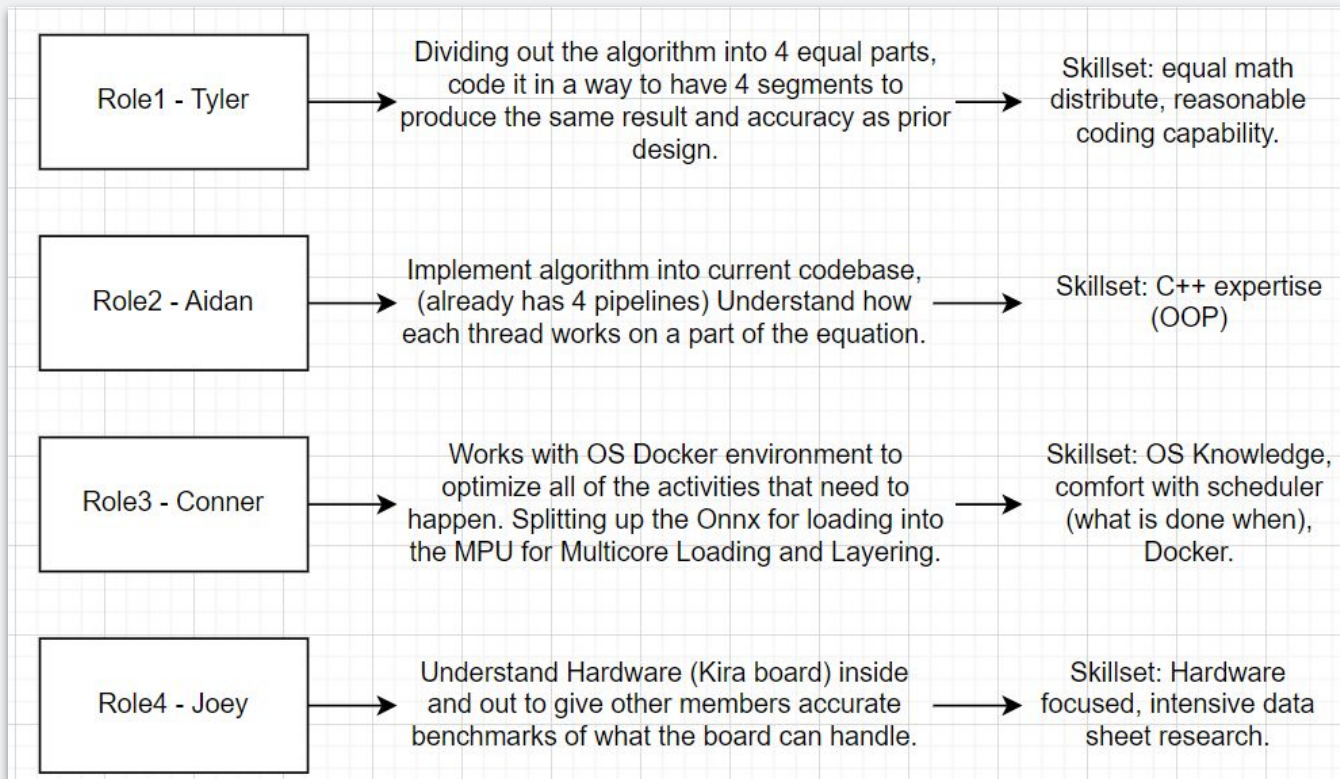
Why not Google drive?

What is S3?

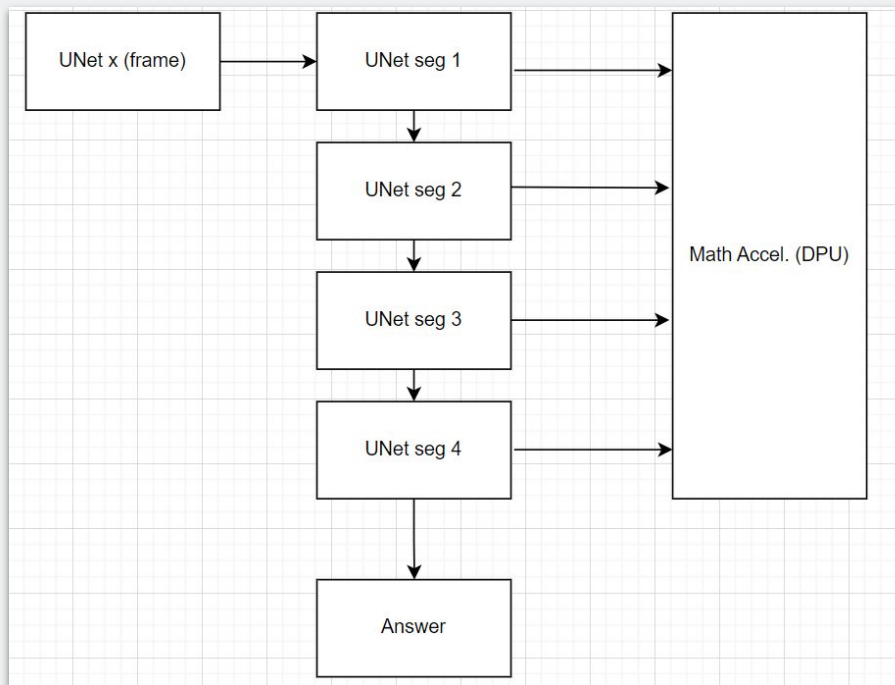


Task Decomposition and Decision Making

Member Roles



Proposed Pipelining

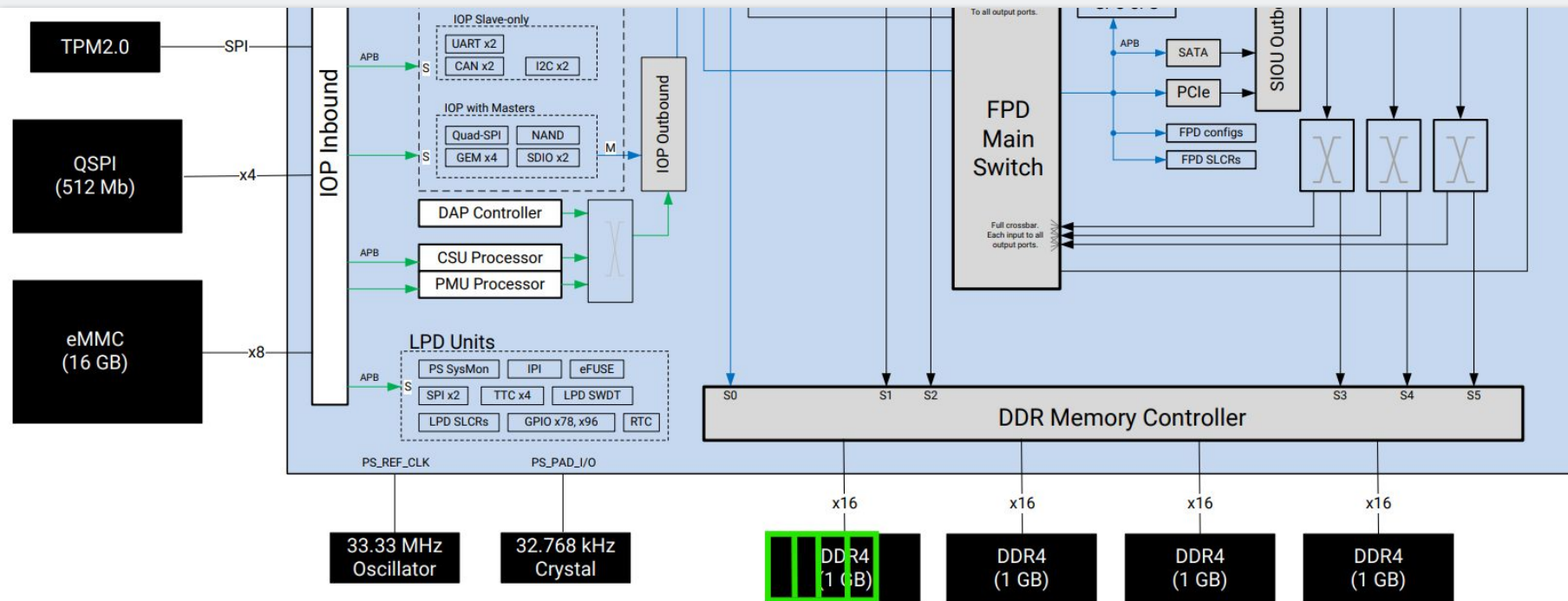


Tracked Metrics

- Throughput
- Accuracy
- Resource Utilization



Emphasis on Resource Utilization



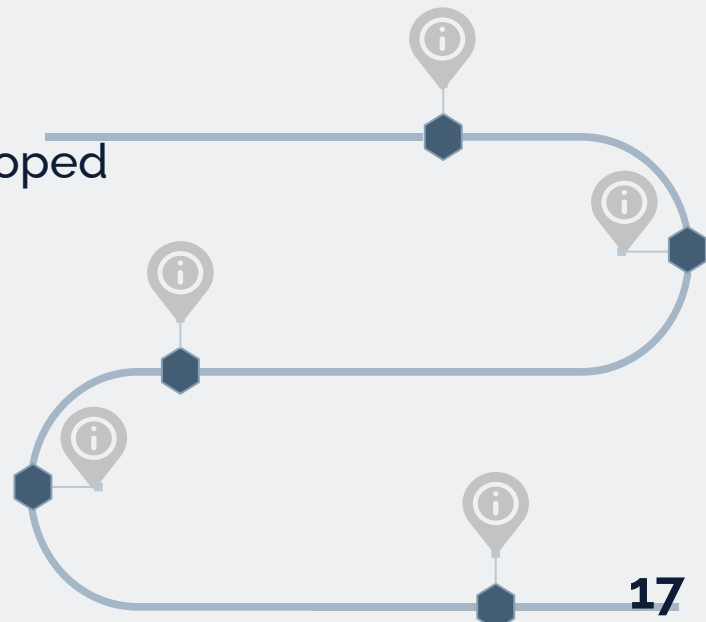
X24999-012122

Milestones

- Mathematical division of the Algorithm
- Loading of Split Algorithm weights onto DPU
- Pipelined Implementation of the Semantic

Segmentation algorithm across the 4 developed threads.

- Increased Throughput over multiple frames

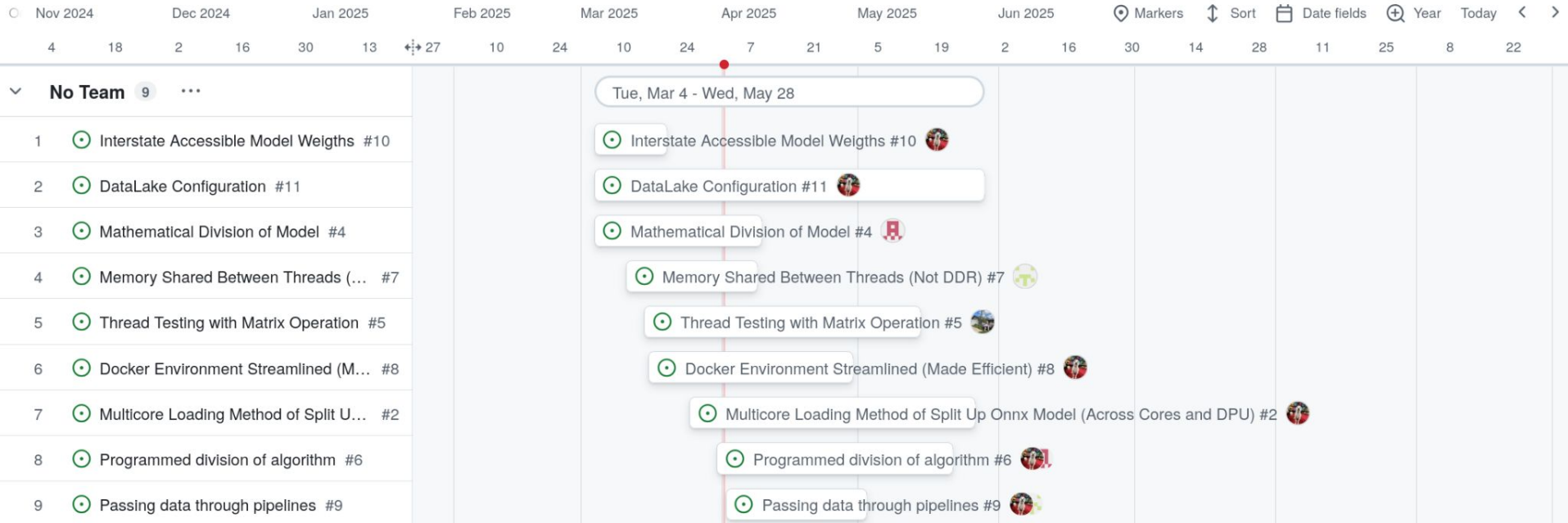


Project Management Style

Waterfall & Agile

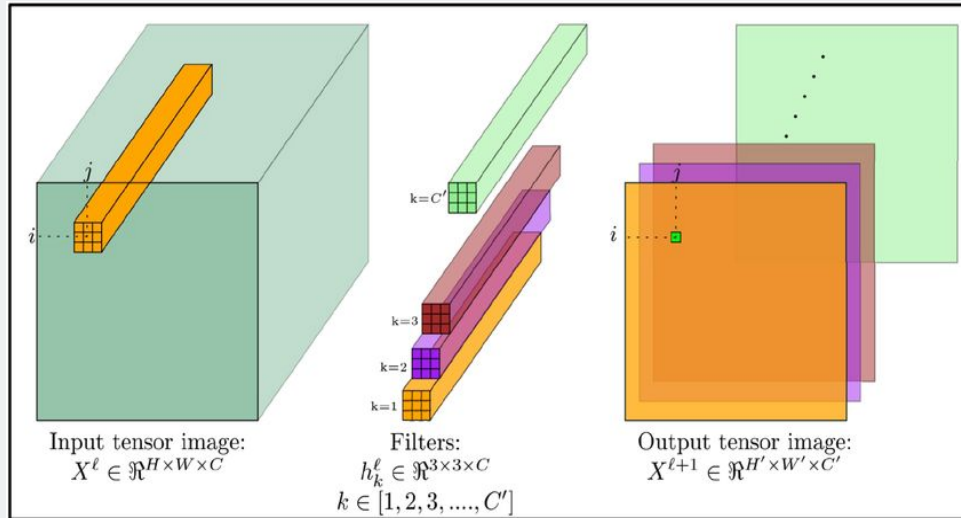


Gantt Chart



Segmenting the Model

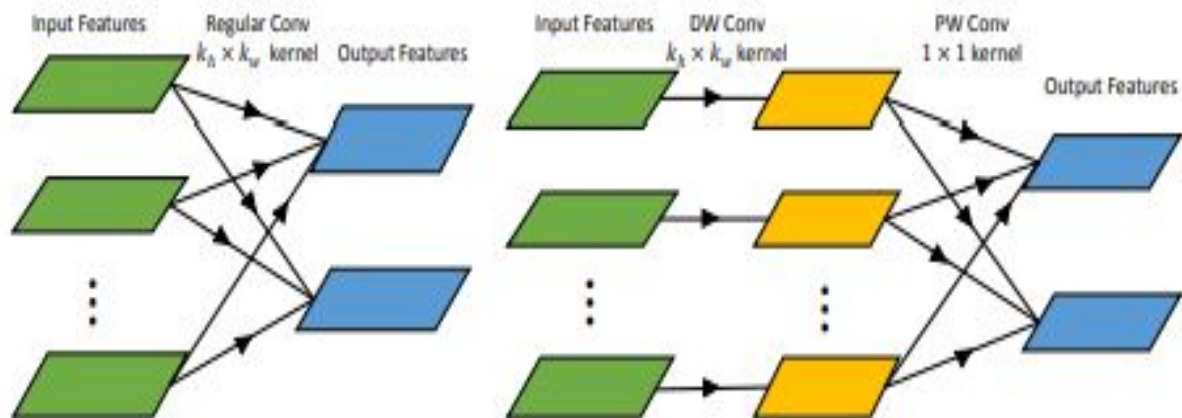
Computational Complexity of Model



Traditional Convolutions

$$O(H \times W \times C_{in} \times K \times C_{out})$$

Computational Complexity Analysis



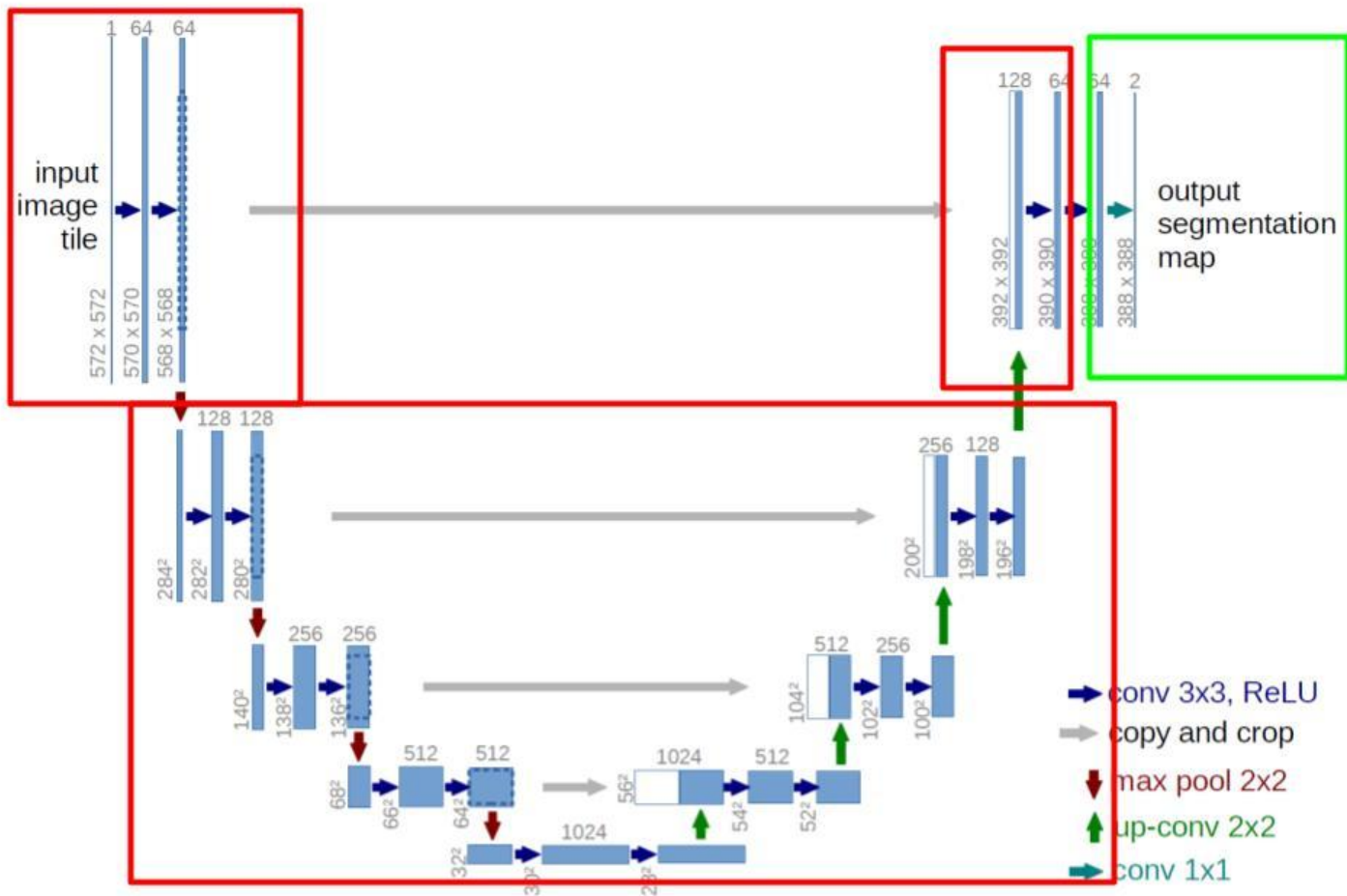
Depthwise Separable Convolutions

Depthwise

$$O(H \times W \times C_{in} \times K)$$

Pointwise

$$O(H \times W \times C_{in} \times C_{out})$$



Conclusion

Problem Solved: Real-time monitoring for individuals with disabilities using eye tracking.

Outcome: Improve safety and throughput with pipelined U-Net on DPU.

Next Steps: Optimize Performance and Thorough Safety Testing.



Thank You

:)

Questions

?