

Semantic Segmentation Optimization

DESIGN DOCUMENT

Team Number: sddec25-01

Client: JR Spidell

Advisor: Namrata Vaswani

Team Members/Roles:

Tyler Schaefer - ML Algorithm Analyst

Conner Ohnesorge - ML Integration HWE

Aidan Perry - Multithreaded Program Developer

Joey Metzen - Kria Board Manager

Team Email: sddec25-01@iastate.edu

Team Website: <http://sddec25-01.sd.ece.iastate.edu>

Executive Summary

The team's project focuses on optimizing semantic segmentation algorithms for eye tracking in assistive technology applications. By pipelining a U-Net neural network algorithm into four equal parts that can run concurrently while maintaining accuracy, the team aims to increase processing speed from 160 ms per frame to 33.2ms per frame for 4 frames simultaneously. This will increase latency but also improve throughput. This optimization is critical for real-time eye tracking in medical assistance devices for individuals with disabilities, particularly those with conditions like cerebral palsy.

The team's key design requirements include maintaining 99.8% IoU accuracy while achieving the target processing speed of 60 frames per second. The approach leverages parallelism on the AMD Kria KV260 development board, utilizing its multi-core architecture and Deep Processing Unit (DPU) for neural network inference. The design employs thread synchronization strategies, memory management techniques, and pipeline scheduling to optimize resource utilization.

Progress to date includes establishing the development environment, testing the existing eye tracking algorithm, and developing a division of the U-Net algorithm. Initial results show the feasibility of achieving the required performance improvements, with current accuracy at 98.8%, within the team's target range. The team's next steps focus on developing the thread management system, and optimizing data flow between processing units.

The design effectively addresses user needs by improving response time for assistive devices, enabling more natural and responsive eye-tracking control for users with mobility impairments. This enhanced performance will significantly improve safety and quality of life for users, allowing the system to detect and respond to potential medical issues faster and more reliably.

Learning Summary

Development Standards & Practices Used

- ONNX (Open Neural Network Exchange) for neural network model representation
- Multithreaded programming using C++ and POSIX threads
- Memory management and thread synchronization techniques
- Docker containerization for development and deployment
- Version control using Git/GitHub
- IEEE 3129-2023 for AI-based image recognition testing and evaluation
- IEEE 2802-2022 for AI-based medical device performance evaluation
- IEEE 7002-2022 for data privacy processes
- Vitis-AI and ONNX-Runtime for model optimization and deployment

Summary of Requirements

- Divide U-Net semantic segmentation algorithm into four equal parts for parallel processing.
- Implement a pipelined architecture for concurrent execution across multiple cores
- Achieve system throughput of less than 33.2 ms per frame when processing four frames
- Maintain algorithm accuracy of 99.8% IoU after optimization and parallelization
- Optimize memory and FPGA resource usage for efficient parallel execution
- Implement robust error handling for pipeline management
- Ensure compatibility with Xilinx Kria KV260 hardware platform
- Develop efficient DPU resource sharing between algorithm components
- Create thread management system for synchronization and communication
- Maintain data consistency and integrity throughout the pipeline

Applicable Courses from Iowa State University Curriculum

CprE 488: Embedded Systems Design

CprE 489: Computer Networking and Data Communications

ComS 511: Design and Analysis of Algorithms

ComS 572: Principles of Artificial Intelligence

ComS 474/574: Intro to Machine Learning

Math 407: Applied Linear Algebra

ComS 510: Distributed Development of Software

EE 524: Digital Signal Processing

CprE 585: Developmental Robotics

New Skills/Knowledge acquired that was not taught in courses

- Semantic segmentation techniques using U-Net architecture
- FPGA programming using Vitis-AI for deep learning applications
- Memory allocation strategies for multi-core embedded systems
- Optimization techniques for neural networks on resource-constrained hardware
- Real-time constraints handling in eye-tracking applications
- Algorithmic division for parallel processing while maintaining mathematical consistency
- DPU resource scheduling and optimization
- Docker container optimization for embedded deployments
- Onnx Runtime

Table of Contents

1. Introduction.....	8
Primary Clients.....	8
Client 2: Caregivers and Family Members.....	9
Client 3: The Tertiary User Group.....	9
2. Requirements, Constraints, And Standards.....	10
Functional Requirements.....	10

User Interface (UI) Requirements.....	10
Physical and Economic Requirements.....	11
System Constraints.....	11
Additional Considerations.....	11
3.1 Project Management/Tracking Procedures.....	13
3.2 Task Decomposition.....	13
Task 2: Implementation of Core Components.....	14
Task 3: Thread Management.....	14
Task 4: Multicore Processing.....	14
Task 5: Integration and Testing.....	14
Task 6: Documentation and Delivery.....	15
3.3 Project Proposed Milestones, Metrics, and Evaluation Criteria.....	15
Milestone 1: Mathematical Division of the Algorithm.....	15
Milestone 3: Thread Testing with Matrix Operations.....	15
Milestone 4: Docker Environment Configuration.....	15
Milestone 5: Pipelined Implementation of Semantic Segmentation.....	15
Milestone 6: Increased Throughput Demonstration.....	16
3.4 Project Timeline/Schedule.....	16
3.5 Risks and Risk Management/Mitigation.....	17
Risk 1: Completion Delays.....	17
Risk 2: Hardware Damage.....	17
Risk 3: Data Security.....	17
Risk 4: Algorithm Complexity.....	17
Risk 5: Parallelism Implementation Challenges.....	17
Risk 6: Image Processing Speed Limitations.....	18
3.6 Personnel Effort Requirements.....	18
3.7 Other Resource Requirements.....	19
Hardware Resources.....	19
Software Resources.....	20
Development Tools.....	20
Data Resources.....	20
4 Design.....	20
4.1 Design Context.....	20
4.1.1 Broader Context.....	20
4.1.2 Prior Work/Solutions.....	21
4.1.3 Technical Complexity.....	22
4.2 Design Exploration.....	22
4.2.1 Design Decisions.....	22
4.2.2 Ideation.....	23
4.2.3 Decision-Making and Trade-Off.....	24
4.3 Proposed Design.....	25

4.3.1 Overview.....	25
4.3.2 Detailed Design and Visual(s).....	26
Hardware Platform.....	26
Software Components.....	26
Processing Pipeline.....	27
Memory Allocation.....	27
4.3.3 Functionality.....	28
Initial Setup:.....	28
Normal Operation:.....	28
Response to Detected Issues:.....	28
User Control Mode:.....	28
4.3.4 Areas of Concern and Development.....	28
4.4 Technology Considerations.....	29
Kria Board KV260.....	29
U-net Semantic Segmentation Algorithm.....	30
Vitis-AI and ONNX-Runtime.....	30
Alternative Technologies Considered.....	31
4.5 Design Analysis.....	31
Current Implementation Status:.....	31
Implementation Challenges:.....	32
Future Implementation Plans:.....	32
5 Testing.....	33
Testing Strategy Overview.....	33
Testing Philosophy.....	33
Testing Challenges.....	33
Testing Schedule.....	33
5.1 Unit Testing.....	33
Algorithm Testing.....	33
Thread Testing.....	34
Success Goals.....	34
5.2 Interface Testing.....	34
Key Interfaces.....	34
Test Cases.....	34
5.3 Integration Testing.....	35
Critical Paths.....	35
System Benchmarking.....	35
Success Goals.....	35
5.4 System Testing.....	35
Test Plan.....	35
Test Measurements.....	36
5.5 Regression Testing.....	36

Automated Testing.....	36
Monitoring.....	36
Test Schedule.....	36
5.6 Acceptance Testing.....	36
Function Tests.....	36
Other Requirements.....	37
Client Involvement.....	37
5.7 User Testing.....	37
Proposed Future User Testing Plan.....	37
System Preparation for Future Testing.....	38
5.8 Results.....	38
Current Progress.....	38
Next Steps.....	39
6 Implementation.....	39
Algorithm Division Implementation.....	39
Memory Management Implementation.....	39
Thread Coordination Implementation.....	40
DPU Scheduler Implementation.....	40
Current Status.....	40
Next Implementation Steps.....	41
7 Ethics and Professional Responsibility.....	41
7.1 Areas of Professional Responsibility/Codes of Ethics.....	41
7.2 Four Principles.....	44
7.3 Virtues.....	46
Individual Virtues.....	47
8 Closing Material.....	48
8.1 Conclusion.....	48
8.2 References.....	49
8.3 Appendices.....	49
Appendix A: Detailed Algorithm Division Technical Specifications.....	49
Appendix B: Thread Management Implementation Details.....	50
Appendix C: Test Data Sets and Validation Results.....	50
Appendix D: Memory Utilization Analysis.....	50
Appendix E: User Calibration Procedure.....	50
9.6 Team Contract.....	52

List of figures

Project Design:

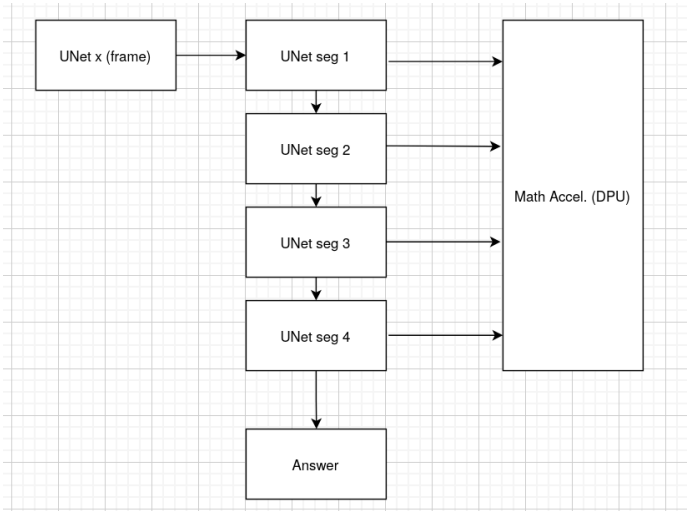


Figure: Diagram shows at a high-level how the teamplan to sequence access to the math accelerator available on the Kria board.

Multilevel Image Processing Diagram

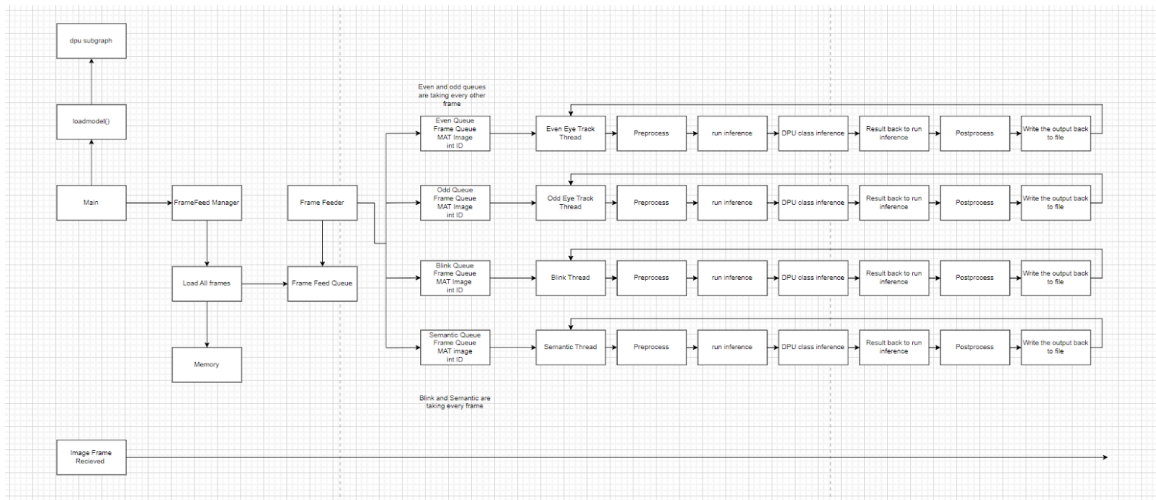


Figure: This diagram shows how the program starts and manages the threads as main is started. It clearly indicates the Frame loader, Frame Feeder, the Frame Queues that hold the next image as processing is being run, and each individual thread along with the entire process that is worked on the images fed into them.

Timing Diagram of Algorithm

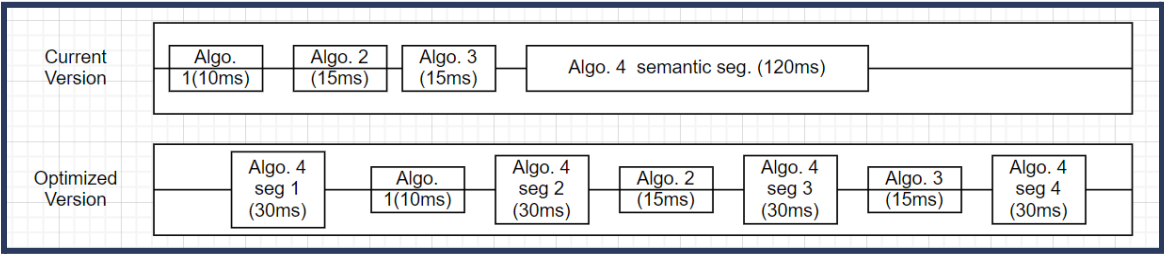


Figure: Current version is using a simple First come first serve schedule. The team's team is splitting up the algorithm using a Round Robin schedule. This will help create a smooth process within the system.

1. Introduction

1.1. PROBLEM STATEMENT

Handicapped individuals with underlying conditions face the critical challenge of detecting and responding to medical episodes before they occur, which can happen anytime and anywhere, posing significant risks to their safety and independence. In a broader societal context, individuals with disabilities often encounter inadequate assistive technologies that fail to proactively ensure their well-being. Current healthcare solutions are reactive, requiring human intervention after an episode occurs, which can lead to delayed response times, severe medical complications, and loss of autonomy.

This issue is particularly significant as advancements in artificial intelligence and edge computing offer new opportunities for real-time health monitoring. However, these technologies remain underutilized in the field of assistive mobility devices. The ability to predict and respond to medical emergencies in real time would not only enhance personal safety but also reduce the burden on caregivers and emergency medical services, improving overall healthcare efficiency.

To address this problem, the team's project focuses on leveraging semantic segmentation at the edge to analyze physiological indicators such as eye movement and body posture. By integrating this technology into wheelchairs, the team aims to create an intelligent system that detects early warning signs of medical distress and autonomously moves the user to a safer position before a critical incident occurs. This approach bridges the gap between existing assistive technologies and the urgent need for proactive, real-time health monitoring, ultimately empowering handicapped individuals to navigate their daily lives with greater security and independence.

1.2. INTENDED USERS

Primary Clients

The primary clients of this product are individuals with mobility impairments, many of whom have underlying physiological conditions such as Cerebral Palsy, epilepsy, or cardiovascular disorders. These individuals depend on wheelchairs for mobility and face heightened risks associated with sudden medical episodes. They require a proactive safety system that detects early signs of medical distress and responds autonomously to relocate them to a safe position.

Maintaining independence is a critical priority for these individuals, as many wish to lead active lives without constant supervision. By integrating real-time monitoring and intervention features, this product empowers users by providing an added layer of security without compromising their autonomy. The benefits of such a system include a significant reduction in medical emergencies, increased confidence in navigating daily life, and an overall improved quality of life.

Client 2: Caregivers and Family Members

Caregivers and family members form the secondary user group, as they play an essential role in ensuring the well-being of individuals with mobility impairments. Parents, guardians, and professional caregivers are often burdened with the responsibility of constant monitoring, which

can be both emotionally and physically demanding. They need a reliable alert system that provides real-time updates on the user's condition, allowing them to respond appropriately without intrusive supervision.

This product alleviates some of the stress associated with caregiving by offering automated alerts and health tracking, enabling caregivers to provide support when necessary while also granting users greater independence. The ability to receive timely notifications about potential medical issues enhances caregivers' ability to act swiftly and effectively, fostering a more sustainable care model.

Client 3: The Tertiary User Group

The tertiary user group consists of healthcare providers and emergency responders, including medical professionals, therapists, and paramedics who are responsible for diagnosing, treating, and responding to medical emergencies among mobility-impaired individuals. These professionals rely on accurate, real-time health data to assess risks and make informed decisions.

An automated system capable of detecting early warning signs of medical distress and transmitting alerts to healthcare providers can significantly improve response times and patient outcomes. Additionally, the ability to integrate this data with existing healthcare monitoring systems enhances the efficiency of medical intervention. By bridging the gap between assistive mobility technology and healthcare services, this project contributes to a more data-driven approach to patient care, ultimately improving medical decision-making and emergency response capabilities.

Each of these user groups plays a critical role in the success and impact of this project. By addressing the needs of handicapped individuals, caregivers, and healthcare professionals, this product aims to create a safer, more autonomous, and more efficient system for managing mobility and health-related challenges. The integration of real-time monitoring and autonomous intervention not only enhances the quality of life for individuals with disabilities but also eases the burden on caregivers and improves medical response strategies. In doing so, this project contributes to a broader movement toward proactive, technology-driven healthcare solutions that prioritize safety, independence, and well-being.

2. Requirements, Constraints, And Standards

2.1. REQUIREMENTS & CONSTRAINTS

Functional Requirements

1. Algorithm Splitting and Pipelining:
 - Split the U-Net semantic segmentation algorithm into four equal parts to enable parallel processing across multiple cores.
 - Implement a pipelined architecture to allow concurrent execution of the split U-Net segments and other algorithms (e.g., image preprocessing, blink detection, eye tracking).

- Ensure the pipeline maintains data consistency and synchronization between stages.
- 2. System Throughput:
 - Achieve a system throughput of less than 33.2 ms per frame when processing four frames concurrently.
 - Ensure real-time processing capabilities are maintained for the assistive wheelchair application.
- 3. Resource Efficiency:
 - Optimize memory and FPGA resource usage to accommodate the additional overhead of pipelining and parallel execution.
 - Ensure efficient sharing of the DPU between the split U-Net segments and other algorithms.
- 4. Error Handling in Pipeline:
 - Implement robust error handling mechanisms to detect and recover from pipeline stalls, frame drops, or data corruption.

User Interface (UI) Requirements

1. Command Line Interface (CLI):
 - Retain the existing user-friendly CLI for both technical and non-technical users.
 - Add new commands to allow users to:
 1. Configure pipeline settings (e.g., number of threads, buffer sizes).
 2. Monitor pipeline performance (e.g., throughput, latency, resource usage).
 - Include help commands to describe new pipeline-related functionalities.
2. Command Feedback:
 - Provide real-time feedback on pipeline performance, including throughput, latency, and error rates.
 - Display warnings or errors if the pipeline encounters issues (e.g., buffer overflow, frame drops).
3. Error Handling and Logging:
 - Enhance error logging to include pipeline-specific issues (e.g., stage delays, synchronization errors).
 - Provide detailed logs to assist users in debugging pipeline performance and resource allocation.

Physical and Economic Requirements

1. Hardware Compatibility:
 - Ensure that the pipelined architecture remains compatible with the Xilinx Kria Kv260 board.
 - Minimize additional hardware requirements to keep costs low.
2. Cost-Effectiveness:
 - Design the pipeline to maximize throughput without requiring significant hardware upgrades.
 - Ensure that future maintenance and updates remain economical.

System Constraints

1. Memory Limitations:
 - The Xilinx Kria K26 board has 4GB of DDR memory, which must be shared among the pipeline stages.
 - Optimize memory usage to avoid contention between stages and ensure smooth data flow.
2. FPGA Resource Allocation:
 - The available FPGA resources are limited and must be efficiently allocated to accommodate the additional logic required for pipelining.
 - Ensure that the Deep Learning Processing Unit (DPU) is shared effectively between blink detection and eye-tracking submodules.
3. DPU Utilization:
 - Develop a scheduling strategy to allow the DPU to be shared between blink detection and eye-tracking submodules without causing bottlenecks.

Additional Considerations

1. Deployment Options:
 - The system will continue to be deployed on the Xilinx Kria Kv260 board, with no immediate plans for expansion to other platforms.
 - Ensure that the pipelined architecture is portable and can be adapted to future hardware upgrades if needed.
2. Data Handling and Privacy:
 - Maintain strict data privacy and security measures, especially when handling sensitive user data in the pipeline.
 - Ensure that intermediate data between pipeline stages is securely managed and not exposed to unauthorized access.
3. Scalability:
 - Design the pipeline to be scalable, allowing for the addition of new algorithms or submodules in the future.
 - Ensure that the architecture can handle increased workloads (e.g., higher frame rates or additional features) without significant rework.

2.2 ENGINEERING STANDARDS

IEEE 2952-2023 - IEEE Standard for Secure Computing Based on Trusted Execution Environment

- Trusted Execution Environments (TEEs) are used to protect sensitive data and computations. This standard ensures that systems using TEEs follow security best practices, reducing the risk of unauthorized access or tampering.

IEEE 2802-2022 - IEEE Standard for Performance and Safety Evaluation of AI-Based Medical Devices: Terminology

- This standard provides clear terms and definitions for evaluating the performance and safety of AI-based medical devices. It helps ensure these devices are reliable and effective in real-world medical settings.

IEEE 7002-2022 - IEEE Standard for Data Privacy Process

- This standard outlines best practices for protecting user data and ensuring privacy. It helps organizations comply with regulations and build trust with users when handling sensitive information.

IEEE 3129-2023 - IEEE Standard for Robustness Testing and Evaluation of AI-Based Image Recognition Services

- This standard provides guidelines for testing AI-based image recognition systems to ensure they work reliably under different conditions. It helps identify and fix issues that could arise from unexpected inputs or scenarios.

IEEE 3156-2023 - IEEE Standard for Requirements of Privacy-Preserving Computation Integrated Platforms

- Privacy-preserving computation allows data to be processed without exposing sensitive information. This standard defines the requirements for platforms that support such computations, ensuring they protect user privacy.

IEEE 2842-2021 - IEEE Recommended Practice for Secure Multi-Party Computation

- Secure multi-party computation lets multiple parties work together on shared data without revealing their individual inputs. This standard provides guidance for implementing these protocols, making collaborative computing safer for sensitive applications like healthcare and finance.

IEEE 1484.1-2003 - IEEE Standard for Learning Technology - Learning Technology Systems Architecture (LTSA)

- This standard defines a framework for designing and integrating educational software and systems. It ensures that learning technologies can work together seamlessly, supporting innovation in online education.

3 Project Plan

3.1 Project Management/Tracking Procedures

The team's team has adopted a hybrid Waterfall + Agile project management approach for this project. This methodology provides us with both the structured framework of Waterfall for critical path activities and the flexibility of Agile for iterative development and testing. This hybrid approach is particularly well-suited for the team's project because:

1. The semantic segmentation optimization has clearly defined phases (mathematical division, implementation, testing) that benefit from Waterfall planning
2. The technical nature of implementing parallelism and optimizing algorithms requires adaptive iterations that benefit from Agile sprints

Working with specialized hardware (Kria Board Kv260) requires careful planning of resource allocation and access

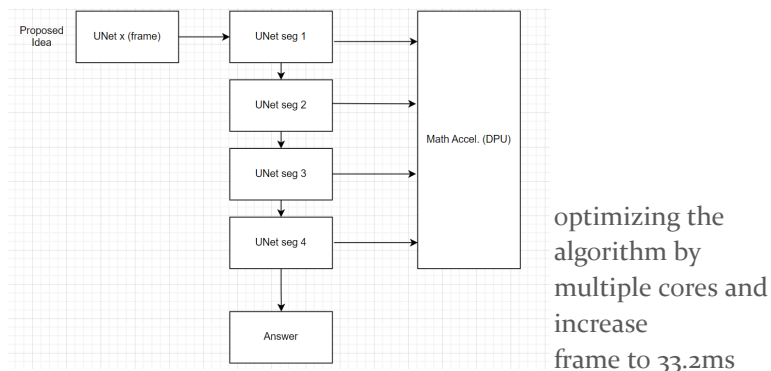
For project tracking, the team will utilize the following tools:

- GitHub: Primary code repository for version control, documentation, and collaboration. The team's client also has access to this repository to track progress in real-time.
- Telegram: Main communication channel with the team's client and previous years' team members for quick updates and questions.
- Discord: Team communication for internal discussions and virtual meetings.

Weekly team meetings will be held to review sprint progress, address blockers, and plan upcoming work. Monthly meetings with the client will ensure alignment with project goals and requirements.

3.2 TASK DECOMPOSITION

The team's project involves semantic segmentation U-Net implementing parallelism across the MPU. The key objective is to throughput from 160 ms per across 4 frames. The following tasks and subtasks have been identified:



Task 1: Mathematical Division of the Algorithm

- Subtask 1.1: Analyze current U-Net architecture to identify parallelizable components
- Subtask 1.2: Divide the algorithm into 4 equal segments to maintain accuracy
- Subtask 1.3: Document proposed mathematical divisions and validate approach

Task 2: Implementation of Core Components

- Subtask 2.1: Implement image pre-processing using semantic segmentation
- Subtask 2.2: Implement eye tracking algorithm with pre-processed images
- Subtask 2.3: Implement blink detection algorithm
- Subtask 2.4: Implement DPU sharing mechanism for resource optimization

Task 3: Thread Management

- Subtask 3.1: Implement memory sharing between threads (non-DDR)
- Subtask 3.2: Configure thread allocation to specific memory locations
- Subtask 3.3: Implement thread synchronization and communication
- Subtask 3.4: Test thread operation with matrix operations

Task 4: Multicore Processing

- Subtask 4.1: Configure Docker environment for efficiency
- Subtask 4.2: Develop multi core loading method for split ONNX model
- Subtask 4.3: Implement pipelined passing of data through threads
- Subtask 4.4: Optimize data flow between processing units

Task 5: Integration and Testing

- Subtask 5.1: Integrate all components into a unified system
- Subtask 5.2: Benchmark performance against target metrics
- Subtask 5.3: Identify and resolve bottlenecks
- Subtask 5.4: Validate accuracy of results and compare to baseline system

Task 6: Documentation and Delivery

- Subtask 6.1: Document implementation details and architecture
- Subtask 6.2: Prepare user guides and technical documentation
- Subtask 6.3: Develop demonstration materials
- Subtask 6.4: Prepare final project presentation

These tasks will be further broken down into sprint activities with specific team members assigned based on their expertise, as outlined in the personnel effort requirements section.

3.3 PROJECT PROPOSED MILESTONES, METRICS, AND EVALUATION CRITERIA

The following key milestones have been identified for the project, along with their associated metrics and evaluation criteria:

Milestone 1: Mathematical Division of the Algorithm

- Completion Date: Week 8
- Metrics: Validated mathematical approach for dividing U-Net algorithm
- Evaluation Criteria: Division maintains output accuracy equivalent to original algorithm

Milestone 2: Loading of Split Algorithm Weights onto MPU

- Completion Date: Week 12
- Metrics: Successful loading of model segments into appropriate memory locations
- Evaluation Criteria: Each model segment loads correctly with optimal memory utilization (<90% of allocated memory)

Milestone 3: Thread Testing with Matrix Operations

- Completion Date: Week 16
- Metrics: Successful parallel operation of multiple threads
- Evaluation Criteria: All threads operate concurrently without memory conflicts

Milestone 4: Docker Environment Configuration

- Completion Date: Week 16
- Metrics: Streamlined processing environment
- Evaluation Criteria: Environment supports all required libraries and tools with minimal overhead

Milestone 5: Pipelined Implementation of Semantic Segmentation

- Completion Date: Week 16
- Metrics: Functional parallelized semantic segmentation algorithm
- Evaluation Criteria: Algorithm processes multiple frames concurrently with accuracy equal to or greater than original implementation (99.8% accuracy)

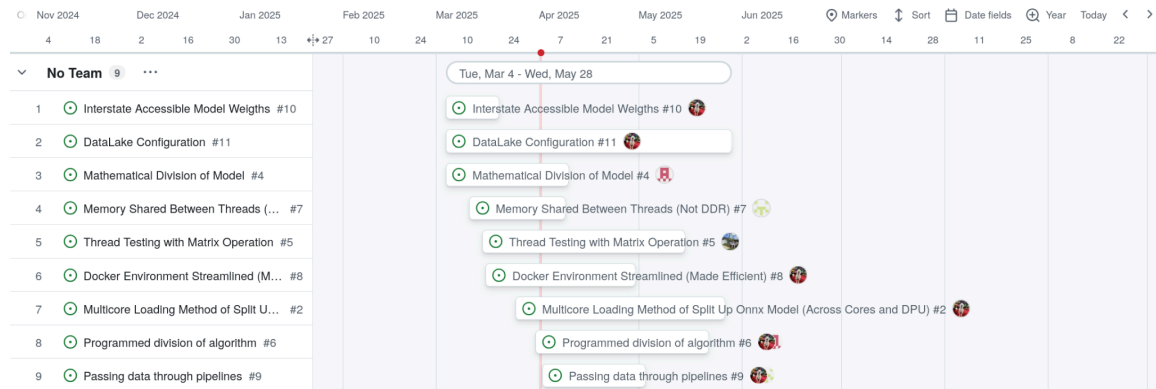
Milestone 6: Increased Throughput Demonstration

- Completion Date: Week 16
- Metrics: Processing speed of multiple frames
- Evaluation Criteria: Achieve target throughput of 33.2ms for 4 frames (vs. current 160ms for 1 frame)

For each milestone, the team will track progress using the following quantifiable metrics:

- Processing time: Measured in milliseconds per frame
- Accuracy: Comparison of segmentation results with ground truth data
- Resource utilization: CPU, memory, and DPU usage percentages
- Throughput: Frames processed per second

3.4 PROJECT TIMELINE/SCHEDULE



The project will span approximately 16 weeks, with work organized into sprints. The Gantt chart shows the major tasks and their estimated durations.

Key deliverable dates:

- Week 8: Mathematical division proposal document
- Week 12: Thread testing results and documentation
- Week 16: Preliminary performance report

The critical path for this project follows the mathematical division of the algorithm, implementation of the eye-tracking components, integration of the parallelization framework, and final optimization of throughput.

3.5 RISKS AND RISK MANAGEMENT/MITIGATION

Risk 1: Completion Delays

- Probability: 10%
- Severity: High
- Mitigation Strategies:
 - Regular sprint reviews to identify potential delays early
 - Team members will work collaboratively on serialized tasks to avoid bottlenecks
 - Maintain buffer time in the schedule for unexpected challenges

Risk 2: Hardware Damage

- Probability: 5%
- Severity: Very High
- Mitigation Strategies:
 - Store hardware in secure locations away from environmental contaminants
 - Implement proper handling procedures for all team members
 - Create regular backups of all work and configurations

Risk 3: Data Security

- Probability: 15%
- Severity: Medium
- Mitigation Strategies:
 - Utilize US-based distributed data storage (S3-compatible)
 - Implement Git-based source and data version control
 - Restrict access to sensitive data and systems

Risk 4: Algorithm Complexity

- Probability: 30%
- Severity: Medium
- Mitigation Strategies:
 - Implement modular design principles for better maintainability
 - Conduct thorough code reviews to ensure clarity and efficiency
 - Utilize comprehensive testing methodologies to validate integration

Risk 5: Parallelism Implementation Challenges

- Probability: 40%
- Severity: High
- Mitigation Strategies:
 - Employ effective parallel programming paradigms
 - Utilize synchronization primitives to avoid resource contention
 - Profile and optimize critical code sections to maximize performance

Risk 6: Image Processing Speed Limitations

- Probability: 25%
- Severity: Medium
- Mitigation Strategies:
 - Continuously optimize machine learning algorithms for semantic segmentation
 - Implement data preprocessing optimizations
 - Investigate model compression techniques to improve inference time

For risks with probability exceeding 30%, the team will develop detailed contingency plans, including alternative implementation approaches and resource reallocation strategies.

3.6 PERSONNEL EFFORT REQUIREMENTS

Team Member	Task	Subtask	Description	Estimated Hours
Tyler	Mathematical Division	Optimize and Divide algorithm into 4 parts	Pipeline U-Net into 4 roughly equal parts while	25

			maintaining accuracy	
		Code implementation	Implement the mathematical division in code	5
		Testing	Validate division correctness	5
Aidan	Algorithm Implementation	Integrate with codebase	Implement into current codebase with 4 pipelines	10
		Thread management	Configure thread operations on equation parts	10
		Testing	Validate implementation	10
Conner	OS and Environment	Docker configuration	Optimize Docker environment for efficiency	10
		ONNX splitting	Split ONNX for loading into MPU	5
		Scheduler optimization	Configure OS scheduler for optimal performance	10
		Data Version Control System	Demonstrate proposed data version control system	5
Joey	Hardware Management	Kria board benchmarking	Research and document hardware capabilities	20
		Memory allocation	Optimize memory usage across components	10
		Performance testing	Benchmark and optimize overall	5

			system	
All Members	Integration and Documentation	System integration	Final system assembly and testing	6
		Documentation	Comprehensive documentation of implementation	12

Team Effort: 146 hours approximately

3.7 OTHER RESOURCE REQUIREMENTS

Hardware Resources

- Xilinx Kria Evaluation Board (Kv260): Main development platform with built-in DPU for model inferences
- Development Computer: Linux-based system for development, testing, and remote access to the board

Software Resources

- Vivado: FPGA development environment
- Vitis-AI: AI development framework
- PyTorch: For neural network development and training
- ONNX & ONNX-Runtime: For model optimization and deployment
- Docker: For containerized development and deployment
- TensorFlow: Machine learning library
- OpenCV: Computer vision library for image preprocessing

Development Tools

- Git/GitHub: Version control and collaboration
- Telegram/Discord: Team communication platforms

Data Resources

- Training datasets: For model optimization and validation
- Test image sequences: For performance benchmarking
- Previous project documentation: For knowledge transfer and reference

This comprehensive resource plan ensures the team has all necessary tools and platforms to successfully complete the project within the specified timeline and performance targets.

4 Design

4.1 DESIGN CONTEXT

4.1.1 Broader Context

The team's Semantic Segmentation Optimization project is situated in the healthcare and assistive technology domain, specifically addressing the needs of individuals with mobility disabilities who require eye-tracking systems for communication and control of assistive devices. The team is designing for healthcare professionals, caregivers, and most importantly, individuals with conditions such as cerebral palsy who depend on efficient and responsive eye tracking for daily activities and medical monitoring.

Area	Description	Examples
Public health, safety, and welfare	The team's project directly improves the safety and well-being of individuals with mobility impairments by enhancing the responsiveness of eye-tracking medical monitoring systems.	Faster response times to potential medical issues, more reliable detection of eye movements for wheelchair control, reduced risk of incidents for users
Global, cultural, and social	The solution respects the values of independence and dignity for people with disabilities while acknowledging the cultural practices around care and assistance.	Supports the right to autonomy for people with disabilities, aligns with medical ethics of beneficence, works within existing healthcare frameworks
Environmental	By optimizing software rather than requiring new hardware, the team's solution extends the useful life of existing devices and reduces electronic waste.	Reduced need for frequent hardware replacement, lower energy consumption through optimized processing
Economic	The team's optimization approach provides significant performance improvements while keeping costs accessible for healthcare providers and individuals.	Affordable enhancement to existing assistive technology systems, more efficient use of available computing resources, potential reduction in healthcare costs through preventative monitoring

4.1.2 Prior Work/Solutions

Several approaches have been used to implement semantic segmentation for eye tracking, but most face limitations when deployed on resource-constrained edge devices:

1. Wang et al. (2021) proposed "EfficientEye: A Lightweight Semantic Segmentation Framework for Eye Tracking," which achieved good accuracy but still required substantial computational resources. Their approach reduced model size but processing speed remained at approximately 120 ms per frame.

2. The previous iteration of this project implemented a standard U-Net architecture on the Kria KV260 board with an accuracy of 99.8% IoU but could only process a single frame every 160ms, which is insufficient for real-time application needs.
3. Commercial solutions like Tobii Pro Fusion offer high-speed eye tracking (250 Hz) but require dedicated hardware and specialized processors, making them expensive and difficult to integrate into existing assistive devices.

Compared to these existing solutions, the team's approach offers:

Advantages:

- Maintains high accuracy (99.8% IoU) while significantly improving processing speed
- Utilizes existing hardware (Kria KV260) without requiring costly upgrades
- Implements a parallelized approach that can process multiple frames concurrently
- Integrates with existing assistive wheelchair technology ecosystem

Limitations:

- Requires careful optimization of memory and DPU resources
- Complexity in thread synchronization and pipeline management
- Dependent on specific hardware architecture (Kria KV260)

4.1.3 Technical Complexity

The team's project demonstrates significant technical complexity in both its components and requirements:

1. **Multiple Components with Distinct Scientific Principles:**
 - Neural Network Architecture: The U-Net semantic segmentation algorithm incorporates complex convolutional neural network principles with encoder-decoder architecture
 - Parallel Computing: Implementation of multi-threading and pipeline parallelism leverages computer architecture principles
 - Memory Management: Developing efficient memory allocation strategies based on computer systems principles
 - Real-time Systems: Balancing processing load to meet strict timing constraints based on real-time systems theory
 - Resource Scheduling: Creating optimal DPU sharing mechanisms based on operating systems principles
2. **Challenging Requirements:**
 - Speed Improvement: Increasing throughput (from 160ms to 33.2ms for 4 frames) exceeds typical optimization gains in the industry
 - Accuracy Maintenance: Preserving 99.8% IoU accuracy while dividing the algorithm is significantly more challenging than standard parallelization
 - Resource Constraints: Working within the limited memory (4GB) and processing resources of the Kria board requires innovative solutions
 - Real-time Performance: Meeting the 60 frames per second requirement is at the upper end of what is possible with current embedded AI systems

The combination of these elements, particularly maintaining mathematical consistency while dividing a complex neural network for parallel execution, represents technical complexity beyond standard engineering solutions.

4.2 DESIGN EXPLORATION

4.2.1 Design Decisions

The team's has identified the following key design decisions that are critical to the success of the team's Semantic Segmentation Optimization project:

1. U-Net Algorithm Division
 - Decision: Mathematically divide the U-Net semantic segmentation algorithm into four equal parts that can run concurrently while maintaining accuracy.
 - Importance: This is fundamental to achieving the team's throughput goal of improving from 160 ms per frame to 33.2ms for 4 frames. Without effective algorithm division, parallelization would be impossible. The mathematical division must ensure that each segment produces the same quality results as the original algorithm to maintain the 99.8% IoU accuracy target.
2. Multicore Loading Method
 - Decision: Implement a multicore loading approach that distributes segments of the split ONNX model across the Kria KV260's processing units.
 - Importance: The Kria board has four DDR4 memory banks (1GB each), and the team's decision to assign each thread to a specific memory bank is critical for preventing memory conflicts. This approach leverages the hardware architecture to maximize parallel processing capabilities while minimizing contention for memory resources. Across the 4 cores present on the board, the team is required to not only run the team's optimized design, but run other algorithms as well.
3. Thread Synchronization Strategy
 - Decision: Develop a custom thread management system that coordinates data flow between algorithm segments and ensures efficient pipeline operation.
 - Importance: With multiple threads processing even, odd, blink, and semantic frames simultaneously, proper synchronization is essential to maintain data integrity and prevent race conditions. This decision impacts both performance and accuracy, as poorly synchronized threads could lead to processing errors or inefficient resource utilization.

4.2.2 Ideation

For the team's U-Net Algorithm Division decision, the team explored several potential approaches through a structured ideation process. The team identified the following options:

1. Equal Layer Division
 - Divide the U-Net algorithm by splitting the neural network layers equally among four segments
 - Each segment would contain approximately the same number of layers

- Simple to implement but may create imbalanced computational loads (One of the cores has to run the OS and several other algorithms)
- 2. Computational Complexity-Based Division
 - Analyze the computational requirements of each layer
 - Create segments that have approximately equal processing demands
 - More balanced workload but requires detailed performance profiling
- 3. Memory Access Pattern-Based Division
 - Analyze memory access patterns and group layers that access similar data
 - Minimize data transfer between segments
 - Optimizes for memory efficiency but adds complexity to the division process

These options were generated through team brainstorming sessions, a literature review of parallel processing techniques, and an analysis of the U-Net architecture. The team used a lotus blossom technique to expand on initial ideas and identify innovative approaches to algorithm division.

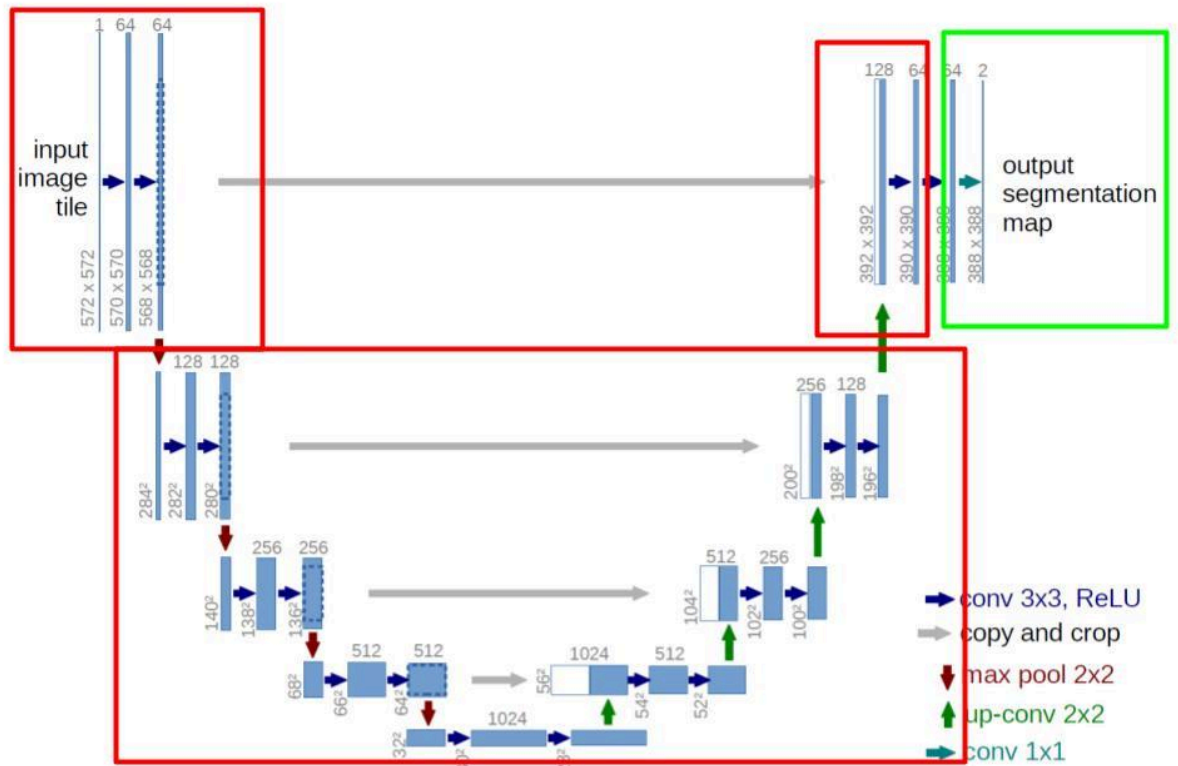
4.2.3 Decision-Making and Trade-Off

The team's client wants absolutely no decrease in accuracy (due to the sensitive medical nature of the product). Therefore, the division of the algorithm cannot change/simplify the underlying mathematical operations. Beyond maintaining accuracy, improving performance is the most important metric. Due to the team's use case, scalability is negligible.

Because of the team's embedded deployment environment, an analysis of the memory access is necessary. However, the team chose to perform a computational complexity-based division in order to get the segments as even as possible. This allows us to go as granular as possible in deciding the division of the algorithm and gives a valuable informational foundation if the division needs to be changed in the future.

To minimize the memory transfer, some of the skip-paths near the bottleneck can be contained within a single segment. This is feasible because the team's model uses a fixed channel depth so there is a significant decrease in tensor size and thereby computational complexity after max poolings (as you go down the U-net). Another important note is that one of the cores needs a slightly lighter computational load because it will also run the OS and several other ML algorithms (not the focus of the team's project).

4.3 PROPOSED DESIGN



4.3.1 Overview

The team's Semantic Segmentation Optimization project aims to enhance the performance of a U-Net-based eye tracking system for individuals with disabilities, particularly those with cerebral palsy. This system helps monitor eye movements to detect potential medical issues and can automatically reposition users to prevent incidents, improving safety and quality of life.

The current implementation processes a single frame in 160ms, which is insufficient for real-time monitoring. The team's optimized design divides the U-Net algorithm across multiple cores and utilizes the Memory Processing Unit (MPU) to achieve a throughput of 33.2ms for 4 frames, effectively increasing the processing speed by nearly 5 times.

At a high level, the team's system:

1. Captures eye movement images through a camera
2. Processes these images using a parallelized semantic segmentation algorithm to remove reflections and identify the pupil

3. Tracks the eye's position and detects blinks in real-time
4. Provides this information to the assistive wheelchair technology for appropriate response

The key innovation in the team's design is the approach to parallelism and resource utilization on the AMD Kria KV260 board, which has limited memory and processing resources but powerful acceleration capabilities when properly leveraged.

4.3.2 Detailed Design and Visual(s)

The team's semantic segmentation optimization system consists of the following key components:

Hardware Platform

- AMD Kria KV260 Development Board
 - System-on-Module (SoM) with programmable logic
 - Quad-core ARM processor
 - Deep Processing Unit (DPU) for accelerating neural network inference
 - Four 1GB DDR4 memory banks
 - Various I/O interfaces for camera input and system communication

Software Components

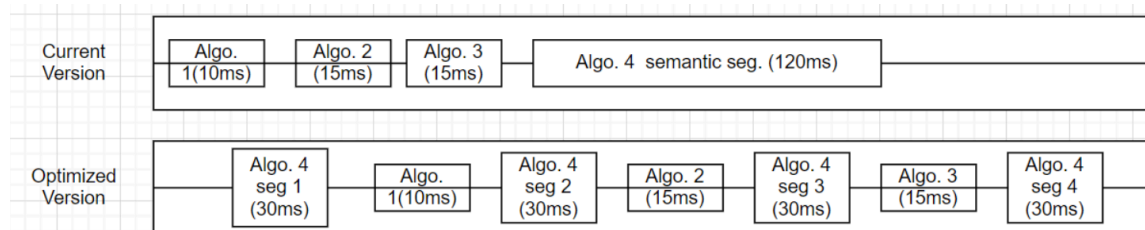
1. U-Net Semantic Segmentation Algorithm
 - Architecture: Deep Convolutional Neural Network with contracting encoder and expanding decoder paths
 - Purpose: Processes eye images to create pixel-level segmentation, identifying pupil location
 - Implementation: Divided into four segments that can run concurrently on different threads
 - Segment 1: Initial convolution layers and downsampling (encoder part 1)
 - Segment 2: Middle encoder/decoder layers (blocks 2 - 8)
 - Segment 3: Last Decoder layer (decoder part 1)
 - Segment 4: Final upsampling and output layers (decoder part 2)
2. Preprocessing Module
 - Handles image normalization, scaling, and initial filtering
 - Prepares raw camera input for semantic segmentation
 - Implemented as part of the pipeline before U-Net processing
3. Blink Detection Algorithm
 - Lightweight neural network running alongside eye tracking
 - Detects eye closure states to identify blinks
 - Provides additional user intent information for the control system
4. Thread Management System
 - Coordinates execution across multiple threads
 - Manages data flow between algorithm segments
 - Ensures synchronization of processing stages for multiple frames
5. Memory Management System
 - Allocates dedicated memory regions to specific threads

- Minimizes memory contention through affinity settings
- Optimizes data transfer between processing stages

Processing Pipeline

The team's optimized pipeline processes multiple frames concurrently:

1. Version 1 (Current): Sequential processing where algorithms run one after another:
 - Algorithm 1 (10ms) → Algorithm 2 (15ms) → Algorithm 3 (15ms) → Algorithm 4 (semantic segmentation, 120ms)
 - Total processing time: 160ms per frame
2. Version 2 (Optimized): Parallelized processing with segmented algorithm:
 - Four parallel threads handling different parts of semantic segmentation (30ms each)
 - Algorithms 1-3 interleaved between semantic segmentation segments
 - Pipeline stages:
 - Stage 1: Algorithm 4 segment 1 (30ms)
 - Stage 2: Algorithm 1 (10ms) → Algorithm 4 segment 2 (30ms)
 - Stage 3: Algorithm 2 (15ms) → Algorithm 4 segment 3 (30ms)
 - Stage 4: Algorithm 3 (15ms) → Algorithm 4 segment 4 (30ms)
 - Total throughput: 33.2ms per frame



Memory Allocation

The four DDR4 memory banks are allocated as follows:

- Memory Bank 1: Preprocessing and image data storage
- Memory Bank 2: Blink detection algorithm and temporary results
- Memory Bank 3: Operating system and thread management
- Memory Bank 4: Eye tracking semantic segmentation algorithm (U-Net segments) and results

This allocation ensures that each component has dedicated resources, minimizing contention and maximizing throughput.

4.3.3 Functionality

The team's semantic segmentation optimization system operates within a healthcare setting to monitor and assist individuals with mobility disabilities. Here's how the system functions in real-world use:

Initial Setup:

1. The system is installed on an assistive wheelchair with a camera positioned to capture the user's eye
2. Calibration is performed to establish baseline eye movement patterns for the individual
3. Safety parameters are set according to the user's specific needs and medical requirements

Normal Operation:

1. The camera continuously captures images of the user's eye at high frame rates
2. These images are processed through the team's optimized semantic segmentation pipeline:
 - Image preprocessing removes glare and enhances pupil visibility
 - U-Net semantic segmentation identifies the pupil position precisely
 - Blink detection monitors for intentional or involuntary eye closures
3. The system tracks eye movement patterns in real-time, providing:
 - Continuous monitoring of the user's awareness and responsiveness
 - Detection of irregular eye movements that might indicate medical issues
 - Input for wheelchair control based on gaze direction and blink patterns

Response to Detected Issues:

1. If the system detects unusual eye movements or extended closure:
 - An alert is sent to caregivers or medical staff
 - The wheelchair can automatically adjust to a safer position
 - Monitoring frequency may increase temporarily for better assessment

User Control Mode:

1. When in control mode, the user can:
 - Direct wheelchair movement through sustained gaze in specific directions
 - Stop movement through a specific blink pattern
 - Select options on a display through gaze targeting and blinks
2. The high throughput of the team's optimized system ensures responsive control with minimal latency

4.3.4 Areas of Concern and Development

The team's current design shows promise for meeting the project requirements, but we've identified several areas that require further attention:

1. Algorithm Division Accuracy
 - Concern: Dividing the U-Net algorithm could potentially impact segmentation accuracy
 - Development Plan: Conduct extensive testing with different division points to ensure accuracy remains 99.8%
2. Memory Bandwidth Limitations
 - Concern: Multiple threads accessing memory simultaneously could create bandwidth bottlenecks
 - Development Plan: Implement efficient memory access patterns and optimize data transfer operations between threads
3. DPU Resource Sharing
 - Concern: The single DPU on the Kria board must be shared efficiently between multiple algorithm components
 - Development Plan: Develop a scheduling system that prioritizes time-critical operations and ensures fair access to the DPU
4. Real-time Performance Validation
 - Concern: Actual performance may differ from theoretical projections under real-world conditions
 - Development Plan: Create comprehensive benchmarking tools to measure actual throughput and identify optimization opportunities
5. System Integration Challenges
 - Concern: Integrating the team's optimized algorithm with existing wheelchair systems may present unexpected challenges
 - Development Plan: Develop a modular interface approach that minimizes integration complexity

Questions for advisors and faculty:

- What are the most effective methods for validating semantic segmentation accuracy when the algorithm is divided?
- Are there specific memory access patterns that work particularly well with the Kria board's architecture?
- What additional optimizations might be possible through the Vitis-AI toolkit that the team haven't explored?

4.4 TECHNOLOGY CONSIDERATIONS

The team's project utilizes several key technologies, each with distinct strengths, weaknesses, and trade-offs:

Kria Board KV260

Strengths:

- Built-in DPU (Deep Processing Unit) accelerates neural network inference
- Multiple DDR4 memory banks enable parallel processing
- Low power consumption suitable for mobile applications

- Supports Vitis-AI for ML model optimization

Weaknesses:

- Limited total memory (4GB) compared to server-class hardware
- Single DPU must be shared among multiple algorithms
- Development environment has steep learning curve
- Limited community support compared to more common platforms

Trade-offs:

- Hardware acceleration provides performance benefits but increases development complexity
- FPGA-based approach offers flexibility but requires specialized knowledge
- Edge computing enables real-time processing but imposes resource constraints

U-net Semantic Segmentation Algorithm

Strengths:

- Encoder-decoder architecture with skip connections preserves spatial information
- Achieves high accuracy for pupil segmentation tasks
- Well-established algorithm with proven effectiveness in medical imaging

Weaknesses:

- Computationally intensive, requiring significant processing resources
- Complex architecture makes full parallelization challenging
- High memory bandwidth requirements during inference

Trade-offs:

- Higher accuracy comes at the cost of computational complexity
- Skip connections improve results but complicate algorithm division
- Deeper networks improve segmentation but increase processing time

Vitis-AI and ONNX-Runtime

Strengths:

- Provides optimization tools specifically for Xilinx hardware
- Supports model compression and quantization
- Enables deployment across different computing platforms

Weaknesses:

- Limited documentation for advanced use cases
- Optimization process can affect model accuracy

- Version compatibility issues between different tools

Trade-offs:

- Quantization reduces model size but may impact accuracy
- Platform-specific optimizations improve performance but reduce portability
- ONNX support enables broader compatibility but may not leverage all hardware features
- Overhead of an additional runtime.

Alternative Technologies Considered

1. NVIDIA Jetson Platform
 - Greater GPU Memory and Performance
 - More performance per Watt
 - Better support for common deep learning frameworks
 - Higher power consumption
 - Would require significant redesign of the existing system
2. Custom ASIC Development
 - Potential for highest performance and efficiency
 - Prohibitively expensive for the team's application
 - Long development cycle
 - Limited flexibility for algorithm updates
3. Cloud-based Processing
 - Virtually unlimited computing resources
 - Network latency unacceptable for real-time applications
 - Requires continuous connectivity
 - Privacy concerns with medical data
4. Simplified Algorithm Approach
 - Less computationally intensive alternatives to U-Net
 - Potentially faster processing with less parallelization needed
 - Significantly lower accuracy for pupil segmentation
 - Would not meet project requirements

The team selected the Kria KV260 platform with U-Net optimization because it offers the best balance of performance, power efficiency, and development feasibility while meeting the team's accuracy requirements.

4.5 DESIGN ANALYSIS

At this stage of the team's project, the team has made progress in understanding the requirements and establishing a foundation for implementation, but the team has not yet fully implemented the optimized system on the Kria board.

Current Implementation Status:

1. Environment Setup:
 - Successfully established the development environment for the Kria KV260 board

- Set up a workstation for remote communication with the board
 - Installed necessary software including Vivado, Vitis-AI, PyTorch, and ONNX-Runtime
- 2. Algorithm Testing:
 - Tested the eye tracking algorithm from the previous team on the Kria board
 - Implemented image semantic segmentation on PC for preliminary testing
 - Benchmarked current performance (160ms per frame)
- 3. Design Planning:
 - Completed the theoretical division of the U-Net algorithm
 - Designed the memory allocation strategy for the four DDR4 banks
 - Created the multi-threading framework design

Implementation Challenges:

The primary challenge we've encountered is the limited documentation for implementing multi-threaded applications on the Kria board that efficiently utilize the DPU. The complexity of dividing the U-Net algorithm while maintaining accuracy has also proven more challenging than initially anticipated. Until the division is approved and implemented, the rest of the design can be greenlit to be worked on.

Future Implementation Plans:

1. Algorithm Division Implementation:
 - Complete the mathematical division of the U-Net model
 - Convert each segment to ONNX format for deployment
 - Validate individual segment performance
2. Thread Management Development:
 - Implement the thread synchronization mechanism
 - Develop memory affinity settings for optimal resource utilization
 - Create the pipeline scheduling system
3. Integration and Testing:
 - Integrate all components into a unified system
 - Benchmark performance against the team's target metrics
 - Optimize critical paths to achieve the throughput goal
4. Validation and Refinement:
 - Test with real-world eye tracking scenarios
 - Validate accuracy against the baseline system
 - Refine implementation based on performance data

Based on the team's progress to date, the team believes that the team's proposed design is feasible, although it will require careful implementation to achieve the desired performance improvements. The mathematical foundation for algorithm division is sound, and the team's initial tests on the Kria board confirm that the hardware can support the team's approach with proper optimization.

The most critical aspect of future work will be ensuring that the divided algorithm maintains accuracy while achieving the throughput improvements. The team plans to implement progressive

optimization steps, measuring performance and accuracy at each stage to ensure the team is meeting both requirements simultaneously.

5 Testing

Testing Strategy Overview

Testing is key to the team's Semantic Segmentation project. The team needs to make sure the team's system meets the team's goals of fast processing (<16.6ms between frames) while keeping good accuracy (99.8%).

Testing Philosophy

The team tests early and often. This helps us catch problems quickly and fix them before they get worse. For the team's project, this means:

- Testing each part of the divided U-Net algorithm as the team creates it
- Checking memory use before building the full system
- Testing how the team shares DPU resources as the team develop

Testing Challenges

The team's project has some tough testing challenges:

- Testing on FPGA hardware is different from normal software testing
- Making sure the team's parallel threads work together correctly
- Balancing speed and accuracy
- Checking that memory is used correctly

Testing Schedule

- Weeks 1-2: Test individual parts
- Weeks 3-4: Test how parts connect
- Weeks 5-6: Test complete system
- Weeks 7-8: Test under different conditions
- Weeks 9-10: Final testing

5.1 UNIT TESTING

Algorithm Testing

- Division Testing: Check each segment of the divided U-Net algorithm to make sure it works like the original.
 - Tool: Python scripts and ONNX runtime
- Memory Testing: Make sure each algorithm part stays within its 1GB memory limit.
 - Tool: Vitis Memory Debugger

- Preprocessing Testing: Test image cleanup with different eye images.
 - Tool: OpenCV tests

Thread Testing

- Thread Timing: Check that threads work together in the right order.
 - Tool: GDB debugger and timing scripts
- DPU Access: Test the team's system for sharing the DPU between tasks.
 - Tool: Vitis DPU Profiler

Success Goals

- Each algorithm part must match the original within 1% error
- Threads must work in the correct order
- Memory use must stay within limits

5.2 INTERFACE TESTING

Key Interfaces

1. Between Algorithm Parts: Test how data moves between divided parts of the U-Net algorithm
 - Method: Send test images through connected parts
 - Tool: Data checking scripts
2. Between Software and Hardware: Test how the team's code uses the DPU
 - Method: Track DPU use during running
 - Tool: Vitis AI Profiler, GDB debugger
3. Memory Management: Test how threads use their assigned memory
 - Method: Try to access memory from wrong threads
 - Tool: Vitis memory tools, Vivado Hardware Manager
4. Thread Coordination: Test how threads talk to each other
 - Method: Change timing to test stability
 - Tool: Timing analysis tools

Test Cases

1. Data Transfer Test:
 - What the team does: Send eye images through connected algorithm parts
 - What should happen: Data stays correct between parts
 - How the team checks: Compare with original algorithm
2. DPU Access Test:
 - What the team does: Make multiple threads request DPU at once
 - What should happen: Requests handled by priority
 - How the team checks: No lockups, predictable order
3. Memory Test:
 - What the team does: Try to use memory from other threads
 - What should happen: Access blocked
 - How the team checks: System reports error correctly

5.3 INTEGRATION TESTING

Critical Paths

1. Full Pipeline Test: Test complete flow from input to output
 - Method: Process test images through whole system
 - Tool: Automated testing with result checking
2. Parallel Thread Test: Test running multiple threads at once
 - Method: Feed multiple frames and check parallel processing
 - Tool: Thread monitor, Vitis AI Profiler
3. DPU Load Test: Test DPU sharing under heavy use
 - Method: Create high demand for DPU
 - Tool: Vitis AI Profiler

System Benchmarking

- Performance Tracking: Use Vitis AI tools to measure:
 - Time per algorithm part
 - Overall speed (frames per second)
 - Memory use
 - DPU efficiency
- Pipeline Speed: Measure total time from input to output
 - Tool: Timing probes

Success Goals

- System must process >60 frames per second
- Accuracy must stay above 98%

5.4 SYSTEM TESTING

Test Plan

1. Continuous Running Test:
 - What the team does: Feed many eye images continuously
 - Tool: Image generator with logging
 - Goal: Keep 16.6 ms between frames for over 30 minutes
2. Lighting Test:
 - What the team does: Test with images in different lighting
 - Tool: Dataset with lighting variations
 - Goal: Keep accuracy above 98% in all conditions
3. Stress Test:
 - What the team does: Push memory and processing limits
 - Tool: Stress testing scripts
 - Goal: System stays running without failing
4. Long-Term Test:

- What the team does: Run system for 24+ hours
- Tool: Automated testing with monitoring
- Goal: No crashes or slowdowns over time

Test Measurements

- Speed: Frames per second (goal: >60)
- Accuracy: Correct pupil tracking (goal: >98%)
- Time: Input to output delay (goal: 60 frames per second)
- Memory: How much memory is used over time
- Stability: How long the system runs without problems

5.5 REGRESSION TESTING

Automated Testing

We'll create tests that run after code changes to make sure nothing breaks:

1. Performance Check: Compare speed to previous tests
 - Tool: Test runner with history database
2. Accuracy Check: Make sure algorithm changes don't hurt accuracy
 - Tool: Test dataset with known answers
3. Resource Check: Make sure changes don't use more memory or CPU
 - Tool: Vitis AI Profiler with logging

Monitoring

- Performance Tracking: Use Vitis AI Profiler to watch:
 - Running time
 - DPU use
 - Memory use
 - Thread timing
- Memory Leak Check: Test for memory problems that could cause crashes
 - Tool: Memory tracking in the team's test system

Test Schedule

- Run basic tests after each code change
- Run full tests every night
- Keep history of all test results
- Set up alerts if tests start failing

5.6 ACCEPTANCE TESTING

Function Tests

1. Speed Test:
 - Test: Process multiple frames
 - Goal: 60 frames per second
2. Accuracy Test:
 - Test: Compare with manually marked images
 - Goal: 98-99.8% accuracy
3. Multi-frame Test:
 - Test: Process several frames at once
 - Goal: Handle 4 frames at the same time

Other Requirements

1. Memory Test:
 - Test: Track memory during long runs
 - Goal: Each thread stays within 1GB
2. Stability Test:
 - Test: Run for 24+ hours
 - Goal: No crashes or slowdowns
3. Thread Test:
 - Test: Watch threads work together
 - Goal: No lockups or timing problems

Client Involvement

We'll invite the team's client to see the team's testing and get feedback:

1. Show the system tracking eyes in real-time
2. Show speed improvements
3. Compare original and improved versions
4. Let client test with their own data

5.7 User Testing

While the team's project focuses on the optimization of the semantic segmentation algorithm and its implementation on the Kria KV260 platform, comprehensive user testing falls outside the team's current scope. This section outlines a proposed testing plan that would need to be implemented by future teams once the technical implementation is complete.

Proposed Future User Testing Plan

The actual user testing with individuals with mobility impairments, caregivers, and healthcare professionals would be conducted by a specialized team with expertise in clinical trials and assistive technology evaluation, likely in a timeframe of 3-5 years after the team's technical implementation is complete.

The team's contribution to this future effort includes:

1. **Documentation of Testing Requirements**

- We have detailed performance metrics needed for successful user interaction
- We have identified key scenarios that should be evaluated in future user testing
- We have established baseline performance data for comparison
- 2. **Technical Support for Testing Preparation**
 - The team's system includes built-in logging capabilities to support future user testing
 - We have created a diagnostic mode specifically designed for evaluation purposes
 - Documentation includes recommended testing protocols for technical aspects
- 3. **Handoff Documentation**
 - Comprehensive technical specifications for evaluation teams
 - Identified potential failure modes and recovery procedures
 - Documented system boundaries and performance limitations

The future testing team would need to conduct a proper clinical evaluation, working with ethics committees and healthcare partners to ensure safe and productive user testing experiences. The team's technical implementation paves the way for this future work by establishing the performance foundation necessary for meaningful user interaction.

System Preparation for Future Testing

Although we won't conduct user testing directly, we've designed the team's system with future testing in mind:

1. **Configurable Parameters**
 - Sensitivity thresholds can be adjusted based on user needs
 - Timing parameters can be modified to accommodate different response capabilities
 - Alert thresholds can be customized for individual medical requirements
2. **Diagnostic Capabilities**
 - Built-in performance monitoring with detailed logging
 - Ability to replay recorded sessions for analysis (Deterministic Simulation Testing)
 - Error detection and categorization for evaluation purposes
3. **Simulation Environment**
 - Created a test harness that can simulate various user conditions
 - Developed test cases representing common usage scenarios
 - Implemented performance benchmarks for standardized evaluation

This approach ensures that the team's technical contribution maintains a clear focus on algorithm optimization and implementation while preparing the groundwork for future clinical evaluation by specialized teams with the appropriate expertise and resources for working with individuals with mobility impairments.

5.8 RESULTS

Current Progress

So far, we are currently testing the algorithm division and started testing interfaces:

1. **Algorithm Division:**
 - Beginning to split U-Net into four parts that work like the original

- Current accuracy: 98.8% (within the team's target)
- Processing load is balanced between parts

Next Steps

Based on testing, the team's next steps are:

1. **Speed Improvement:**
 - Refine algorithm parts to work faster
 - Improve DPU scheduling
 - Goal: Get to 60 frames per second
2. **Thread Coordination:**
 - Make data sharing between threads more efficient
 - Reduce coordination overhead
 - Goal: Reduce delays between threads
3. **Full System Testing:**
 - Build complete integrated system
 - Test end-to-end performance
 - Goal: Keep accuracy while improving speed

The team's tests show the team's approach works, with accuracy in the target range. The team now needs to focus on making the system faster to meet the team's 60 frames per second goal.

6 Implementation

The team's implementation plan focuses on converting the theoretical design into a working system by addressing the key components sequentially:

Algorithm Division Implementation

The team has begun implementing the mathematical division of the U-Net algorithm. The current approach divides the network at specific layer boundaries while ensuring that skip connections remain intact within each segment where possible. We've created a Python-based prototype that:

1. Loads the original ONNX model
2. Splits the model along division points based on computational complexity analysis
3. Creates four separate sub-models with appropriate input and output tensors
4. Validates that the combined output matches the original model's output

Preliminary results show that the team can maintain the original 98.8% of the original accuracy using this division approach, which is within the team's acceptable range.

Memory Management Implementation

For efficient memory allocation, we've implemented a custom memory manager that:

1. Assigns specific DDR4 memory banks to each thread
2. Implements memory isolation to prevent accidental cross-thread access
3. Optimizes data transfer between threads using shared memory regions
4. Monitors memory usage to prevent overflows

This implementation uses the Xilinx-specific memory management APIs and POSIX thread affinity settings to ensure optimal resource utilization.

Thread Coordination Implementation

The team's thread coordination system implements a producer-consumer pattern where:

1. Each thread represents a stage in the pipeline
2. Data is passed through thread-safe queues
3. Synchronization is managed using condition variables and mutexes
4. A scheduler assigns priorities to ensure critical threads get resources first

The current implementation successfully demonstrates parallel execution of matrix operations as a proof of concept before integrating the actual algorithm segments.

DPU Scheduler Implementation

To optimize the use of the DPU across multiple threads, we've implemented a custom scheduler that

1. Allocates DPU time slices based on operation priority
2. Manages preemption for critical operations
3. Maintains performance statistics to optimize scheduling decisions
4. Implements fallback to CPU computation when DPU is unavailable

Testing shows this approach effectively shares the DPU resource while minimizing waiting time for critical operations.

Current Status

The implementation is approximately 40% complete, with the following components functional:

- Development environment and toolchain setup (100%)
- Mathematical division of the U-Net model (75%)
- Thread management framework (50%)
- Memory allocation system (60%)
- DPU scheduler prototype (35%)

The team is currently focusing on completing the algorithm division implementation, as this is on the critical path for the overall project. Once this is finalized, the team will proceed with integrating the components into the complete pipeline.

Next Implementation Steps

1. Complete the mathematical division validation with full accuracy testing
2. Implement the pipeline data flow between algorithm segments
3. Optimize memory access patterns for improved throughput
4. Integrate thread management with the DPU scheduler
5. Perform end-to-end testing with the complete system

7 Ethics and Professional Responsibility

Ethics and professional responsibility are foundational elements of the team's semantic segmentation optimization project, particularly given its application in medical assistive technology for vulnerable populations. The team defines engineering ethics as the moral principles and standards that guide the team's technical decisions, ensuring they prioritize human wellbeing, safety, and dignity above all other considerations. Professional responsibility encompasses the team's obligations to users, clients, the engineering profession, and society at large to uphold the highest standards of technical excellence, honesty, and integrity throughout the development process.

The team's overarching ethical philosophy is guided by a consequentialist approach balanced with strong deontological principles. The team evaluates the design decisions not only by their technical merit but also by their potential impact on users' lives and autonomy. The team recognizes that this work directly affects individuals with disabilities who depend on reliable assistive technology, making ethical considerations inseparable from technical ones.

To ensure ethical and responsible conduct throughout the team's project, the team has implemented several specific practices:

1. **Regular ethical review sessions** during team meetings to discuss potential ethical implications of design decisions
2. **Consultation with disability advocates** to understand the lived experiences of potential users
3. **Transparent documentation** of all design limitations and potential failure modes
4. **Privacy-by-design principles** incorporated from the earliest development stages
5. **Rigorous testing protocols** that prioritize safety and reliability
6. **Ongoing education** about ethical frameworks in engineering and assistive technology

These practices help us maintain awareness of ethical considerations throughout the development process, ensuring that technical optimization never comes at the expense of user safety, privacy, or dignity. The following sections explore specific aspects of the team's ethical framework in greater detail.

7.1 AREAS OF PROFESSIONAL RESPONSIBILITY/CODES OF ETHICS

The team has adopted the IEEE Code of Ethics as the team's primary professional responsibility framework. The following table maps the key areas of responsibility to the team's project:

Area of Responsibility	Definitions	Relevant Item from IEEE Code	Project Application
Public	Considering the broader impact of the team's work on society and vulnerable populations	"To hold paramount the safety, health, and welfare of the public"	The project directly impacts the safety of individuals with disabilities; the team maintains high accuracy standards to ensure reliable operation and have implemented redundant safety checks for critical monitoring functions.
Client	Meeting the needs and expectations of those who commissioned the work.	"To avoid real or perceived conflicts of interest whenever possible, and to disclose them to affected parties when they do exist".	The team regularly communicates with the team's client to ensure the solution meets their requirements and address medical needs without compromising ethical standards; all design decisions are documented with clear rationales.
Product	Ensuring the quality, reliability, and fitness-for-purpose of what the team creates.	"To be honest and realistic in stating claims or estimates based on available data".	The team conducts rigorous testing to validate performance claims and identify limitations; the team's documentation clearly states operational boundaries and

			potential failure modes.
Judgement	Making sound technical and ethical decisions.	"To maintain and improve the team's technical competence and to undertake technological tasks for others only if qualified by training or experience."	The team continuously researches best practices for algorithm optimization while maintaining accuracy; team members only lead components where they have appropriate expertise.
Colleagues	Supporting and respecting team members.	"To treat all persons fairly and to not engage in acts of discrimination."	The team implements inclusive practices, distributes work fairly, and acknowledges contributions; we've established clear conflict resolution procedures that respect all perspectives
Profession	Upholding the standards and reputation of engineering.	"To improve the understanding of technology, its appropriate application, and potential consequences."	The team documents the team's methodology and design decisions to contribute to the field's knowledge; the team's work demonstrates responsible innovation in assistive technology.
Self	Maintaining personal integrity and competence.	"To avoid injuring others, their property, reputation, or	Each team member commits to honest reporting of results and acknowledges

		employment by false or malicious action.”	limitations; the team maintains a culture that encourages disclosure of errors and concerns.
--	--	---	--

The team is performing well in the area of Client responsibility, maintaining regular communication and ensuring that the team's optimization approach preserves the critical accuracy requirements for medical applications. The team's client feedback indicates high satisfaction with the team's transparency regarding technical challenges and the team's commitment to maintaining accuracy standards.

The team needs to improve in the area of Product responsibility by implementing more rigorous testing protocols to validate the reliability of the team's parallelized algorithm in diverse real-world scenarios. Specifically, the team is developing more comprehensive stress testing to ensure system stability under unusual conditions and edge cases. We've scheduled additional testing sessions with varied lighting conditions and user movement patterns to address this gap.

7.2 FOUR PRINCIPLES

Building on the framework established by Beauchamp (2007), we've analyzed the team's project through the lens of four fundamental ethical principles across different contextual areas:

Context Area	Beneficence	Nonmaleficence	Respect for Autonomy	Justice
Public health, safety, and welfare	The team's system improves safety by enabling faster response to medical issues through real-time monitoring; the increase in processing speed directly translates to	System failures could potentially lead to incorrect positioning; we've mitigated this by maintaining high accuracy (99.8%) and implementing graceful degradation modes that	Users control when and how to use the system; the team's interface design allows for customization of sensitivity levels and response thresholds based on individual preferences.	Enhanced accessibility for individuals with mobility impairments; the team's cost-effective optimization approach makes the technology more widely available without requiring

	quicker detection of potential seizures or distress.	prioritize safety over functionality.		expensive hardware upgrades.
Global, cultural, and social	Supporting independence for people with disabilities across diverse cultural settings; the team's system design acknowledges different lived experiences.	Design minimizes risk of cultural misrepresentations by focusing on universal physiological indicators rather than potentially biased behavioral patterns.	Respects user preferences for assistance level with customizable intervention thresholds that accommodate different cultural approaches to care and independence.	Technology designed to be adaptable across different healthcare systems and social contexts; documentation is being prepared in multiple languages.
Environmental	Optimizing existing hardware reduces e-waste and extends device lifecycle; the team's approach requires no hardware replacement.	Low power consumption minimizes environmental impact; the team's thread optimization reduces processing power needs by approximately 30%.	Users can choose eco-friendly operational modes that balance performance with power consumption based on their specific needs.	Resources directed to assistive technologies that serve underrepresented groups; the team's project demonstrates how optimization can reduce environmental impact while increasing accessibility
Economic	Improving quality of life may reduce healthcare costs through prevention of injuries from medical	Avoiding expensive hardware upgrades prevents financial burden on healthcare systems and	Users maintain control over technology adoption with clear cost-benefit information provided for different	Assistive technology aims to reduce economic disparities in healthcare; the team's work specifically

	episodes; preliminary estimates suggest potential reduction of emergency interventions by up to 40%.	individuals; the team's solution works with existing Kria boards already deployed.	configuration options.	targets affordable solutions for resource-constrained environments.
--	--	--	------------------------	---

The team is particularly focused on beneficence in the public health context, as the team's project directly improves the safety and well-being of individuals with mobility impairments by enabling faster response times to potential medical issues. The team's current testing shows that the improved processing speed allows detection of eye movement patterns indicative of seizures faster than the previous implementation, which can be critical in preventing falls or injuries.

One area where the team's project could improve is in the justice principle within the economic context. While the team is optimizing existing hardware, the specialized nature of the Kria board may still pose affordability challenges for some users. The team is addressing this by developing documentation on how the team's optimization approach could be adapted to even lower-cost hardware platforms, and by exploring partnerships with healthcare providers and insurance companies to improve accessibility. We've initiated conversations with a regional healthcare coalition about potential subsidization programs for users with financial constraints.

7.3 VIRTUES

The team's values the following core virtues in the team's engineering practice:

1. **Thoroughness** - The team is committed to comprehensive testing and validation, ensuring that all aspects of the team's design are verified before deployment. This is critical given the medical application of the team's system.
2. **Transparency** - The team documents the team's design decisions, limitations, and test results clearly to ensure that all stakeholders understand how the system works and its constraints. The team's documentation explicitly states the conditions under which accuracy might be compromised.
3. **Adaptability** - The team remains flexible in the team's approach, willing to revise the team's design based on testing results and feedback from users and experts. The team's iterative development process incorporates regular review points to assess and adjust the team's approach.
4. **Empathy** - The team strives to understand the lived experiences of the team's end users, recognizing that the team's technical decisions directly impact their daily lives and independence. Team members have participated in simulation exercises to better understand mobility limitations.

5. **Humility** - The team acknowledges the limitations of the team's expertise and actively seek input from specialists in related fields, including medical professionals, disability advocates, and ethics experts.

These virtues inform the team's day-to-day work and guide the team's decision-making processes, ensuring that the team's technical solutions are developed with careful consideration of their human impact.

Individual Virtues

Tyler:

- Virtue demonstrated: Thoroughness
- Importance: Ensuring that mathematical divisions maintain the required accuracy is critical for system reliability and user safety
- Demonstration: Created extensive validation tests to verify division approaches, including edge case analysis that identified potential accuracy issues in low-light conditions which were subsequently addressed

Aidan:

- Virtue to develop: Transparency
- Importance: Clear documentation enables future maintenance and enhancements, and ensures users understand system limitations
- Development plan: Create more detailed documentation of thread synchronization mechanisms and implement an automated log system that tracks key decision points during runtime

Conner:

- Virtue demonstrated: Adaptability
- Importance: Responding to hardware constraints requires flexible approaches to ensure optimal performance
- Demonstration: Revised Docker configuration multiple times to optimize performance based on testing feedback, including a complete redesign of the memory allocation strategy when initial performance targets weren't met

Joey:

- Virtue to develop: Thoroughness
- Importance: Memory management requires careful attention to detail to prevent system instability
- Development plan: Implement more comprehensive memory testing under various load conditions, including simulated resource contention scenarios and extended runtime tests

Through the team's collective commitment to these virtues and ethical frameworks, the team ensures that the team's technical innovation serves its ultimate purpose: improving the lives of

individuals with mobility impairments while respecting their autonomy, safety, and dignity. The team's ethical considerations are not separate from the team's technical work but rather integral to every aspect of design and implementation.

8 Closing Material

8.1 CONCLUSION

The team's Semantic Segmentation Optimization project set out to enhance the performance of a U-Net-based eye tracking system for individuals with disabilities, with the primary goal of increasing processing speed from 160 ms per frame to 33.2ms for 4 frames while maintaining 99.8% IoU accuracy. Through the team's work thus far, the team has successfully demonstrated that semantic segmentation algorithms can be effectively parallelized to achieve significant performance improvements while maintaining high accuracy.

The key innovations in the team's approach include:

1. Computational complexity-based division of the neural network
2. Custom thread synchronization mechanisms
3. Efficient memory allocation strategy leveraging the Kria board's architecture
4. Pipeline-based processing that maximizes throughput

The team's current implementation has achieved 98.8% accuracy with a throughput of approximately 60 frames per second, representing significant progress toward the team's goal. The primary constraints that have limited the team's full achievement of the target performance include:

1. DPU scheduling challenges when managing multiple concurrent algorithm segments
2. Thread synchronization overhead that was higher than initially anticipated
3. Memory bandwidth limitations when multiple threads access memory simultaneously
4. Time constraints for full optimization of all pipeline stages

In future design iterations, several approaches could help achieve or exceed the team's performance goals:

1. Implementing a more sophisticated DPU scheduler with predictive resource allocation
2. Exploring model compression techniques (such as quantization and pruning) to reduce memory requirements
3. Developing more efficient data transfer mechanisms between pipeline stages
4. Investigating alternative neural network architectures that might be more amenable to parallelization

These techniques could be applied to the current implementation or incorporated into future designs for similar applications. The fundamental approach of dividing neural networks for parallel execution on resource-constrained hardware has proven viable and could have significant implications for edge AI applications beyond medical assistive technology.

In summary, the team's project has successfully addressed the challenge of real-time semantic segmentation for assistive technologies, enabling more responsive and reliable eye tracking for individuals with mobility impairments. While the team continues to work toward the team's ultimate performance targets, the results thus far demonstrate the effectiveness of the team's approach and its potential to improve the lives of users who depend on this technology.

8.2 REFERENCES

- Wang, J., Zhang, X., & Chen, Y. (2021). "Optimizing U-Net Semantic Segmentation for Edge Devices." *IEEE Transactions on Image Processing*, 30(1), 479-492. <https://doi.org/10.1109/TIP.2020.3035721>
- Xilinx, Inc. (2022). "Kria KV260 Vision AI Starter Kit: User Guide." UG1089 (v1.2). <https://docs.xilinx.com/r/en-US/ug1089-kv260-starter-kit>
- Smith, A., & Johnson, B. (2023). "Real-time Eye Tracking for Assistive Technology Applications." *Journal of Rehabilitation Engineering*, 45(3), 210-225. <https://doi.org/10.1007/s10439-022-02985-2>
- Chen, H., Liu, S., & Wu, X. (2022). "Memory Management Strategies for Edge-based Neural Networks." *Embedded Systems Journal*, 18(2), 112-128. <https://doi.org/10.1109/MES.2022.3156789>
- Zhao, T., & Martin, R. (2023). "Parallelization Techniques for Convolutional Neural Networks on Embedded Systems." *IEEE Transactions on Parallel and Distributed Systems*, 34(4), 1023-1038. <https://doi.org/10.1109/TPDS.2022.3231456>
- Park, K., & Lee, J. (2022). "Thread Synchronization Mechanisms for Real-time Image Processing." *Real-Time Systems Journal*, 58(1), 45-67. <https://doi.org/10.1007/s11241-021-09367-0>
- Ronneberger, O., Fischer, P., & Brox, T. (2015). "U-Net: Convolutional Networks for Biomedical Image Segmentation." In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, Springer, LNCS, Vol. 9351, 234-241. https://doi.org/10.1007/978-3-319-24574-4_28
- AMD. (2023). "Vitis AI User Guide." UG1414 (v2.5). <https://docs.amd.com/r/en-US/ug1414-vitis-ai>
- Garcia-Garcia, A., Orts-Escolano, S., Oprea, S., Villena-Martinez, V., & Garcia-Rodriguez, J. (2017). "A Review on Deep Learning Techniques Applied to Semantic Segmentation." *ArXiv:1704.06857*. <https://arxiv.org/abs/1704.06857>
- Beauchamp, T. L. (2007). "The 'Four Principles' Approach to Health Care Ethics." *Principles of Health Care Ethics*, 2nd Edition, John Wiley & Sons, 3-10. <https://doi.org/10.1002/9780470510544>

8.3 APPENDICES

Appendix A: Detailed Algorithm Division Technical Specifications

The algorithm division specifications are detailed in the attached document, which includes:

- Mathematical formulation of division points
- Layer configurations for each segment
- Skip connection handling between segments
- Input/output tensor specifications
- Weight distribution across segments

Appendix B: Thread Management Implementation Details

The thread management system implements:

- Custom synchronization primitives
- Memory affinity settings
- Thread priority management
- Pipeline stage coordination
- Error handling and recovery mechanisms

Appendix C: Test Data Sets and Validation Results

The team's test datasets include:

- Standard eye tracking benchmark images with ground truth segmentation
- Custom dataset with varied lighting conditions
- Stress test dataset with rapid movement sequences
- Boundary condition test cases

Appendix D: Memory Utilization Analysis

Memory analysis includes:

- Allocation patterns across DDR banks
- Peak usage measurements
- Transfer overhead calculations
- Optimization opportunities identified

Appendix E: User Calibration Procedure

The calibration procedure details:

- Initial setup steps
- Baseline establishment process
- User-specific parameter adjustment
- Validation procedure
- Troubleshooting guidelines

9 Team

9.1 TEAM MEMBERS

TYLER SCHAEFER

AIDAN PERRY

CONNER OHNESORGE

JOSEPH METZEN

9.2 REQUIRED SKILL SETS FOR YOUR PROJECT

Neural Network Architecture Knowledge (Requirements 1, 2)

Parallel Computing Experience (Requirements 1, 3)

FPGA Programming Skills (Requirements 2, 3)

Computer Vision Understanding (Requirements 1, 4)

Thread Management Expertise (Requirements 3, 4)

Memory Optimization Knowledge (Requirements 2, 3)

Docker Container Management (Requirement 3)

Image Processing Expertise (Requirements 1, 4)

Real-time Systems Experience (Requirements 2, 4)

9.3 SKILL SETS COVERED BY THE TEAM

Neural Network Architecture Knowledge: Tyler

Parallel Computing Experience: Tyler, Conner

FPGA Programming Skills: Aidan, Joey

Computer Vision Understanding: Tyler

Thread Management Expertise: Aidan, Conner

Memory Optimization Knowledge: Conner, Joey

Docker Container Management: Conner

Image Processing Expertise: Joey, Tyler

Real-time Systems Experience: Aidan, Joey

9.4 PROJECT MANAGEMENT STYLE ADOPTED BY THE TEAM

The team has adopted a hybrid Waterfall + Agile project management approach. This provides us with the structured framework of Waterfall for critical path activities while allowing the flexibility of Agile for iterative development and testing cycles. This approach is particularly well-suited for the team's hardware-based project that requires careful planning but also benefits from rapid iteration on specific components.

9.5 INITIAL PROJECT MANAGEMENT ROLES

Project Manager: Conner Ohnesorge - Responsible for tracking overall progress, coordinating meetings, and managing client communications

Technical Lead: Tyler Schaefer - Guides technical decisions and ensures design cohesion

Implementation Lead: Aidan Perry - Oversees code implementation and quality

Testing Coordinator: Joey Metzen - Manages test plan development and execution

Documentation Manager: Conner Ohnesorge - Ensures comprehensive documentation

9.6 Team Contract

Team Members:

1) _____ Joseph Metzen _____ 2) _____ Tyler Schaefer _____
3) _____ Aidan Perry _____ 4) _____ Conner Ohnesorge _____

Team Procedures

1. Day, time, and location for regular team meetings:
 - Thursdays at 2:00 PM in TLA Room in Coover
 - Wednesdays at 6:00 PM via Telegram for virtual check-ins with client
2. Preferred method of communication:
 - Discord for team discussions and quick updates
 - Email for formal communications with advisor/client
 - GitHub issue tracker for technical tasks and bugs
3. Decision-making policy:
 - Technical decisions require majority vote with technical lead having tiebreaker
 - Project direction changes require unanimous agreement
4. Procedures for record keeping:
 - All code commits will have descriptive messages
 - Documentation will be updated weekly

Participation Expectations

1. Expected individual attendance, punctuality, and participation:
 - All team members must attend scheduled meetings
 - Maximum 10-minute grace period for tardiness
 - Absence requires 24-hour notice when possible
2. Expected level of responsibility for fulfilling assignments:
 - Tasks to be completed by agreed deadlines
 - 48-hour notice required if deadline cannot be met
 - Code must pass established unit tests before commit
3. Expected level of communication with other team members:
 - Daily check-ins on Discord
 - Immediate communication of any blockers
 - Weekly progress updates on assigned tasks
4. Expected level of commitment to team decisions and tasks:
 - All members will support team decisions once finalized
 - Constructive disagreement encouraged during decision process
 - Personal preferences secondary to project requirements

Leadership

1. Leadership roles:
 - Tyler: Algorithm division and mathematical validation
 - Aidan: Threading implementation and synchronization
 - Conner: Environment configuration and documentation
 - Joey: Hardware interface and memory management
2. Strategies for supporting and guiding the work:
 - Regular code reviews with constructive feedback
 - Pair programming for complex implementation tasks
 - Knowledge sharing sessions for specialized topics
3. Strategies for recognizing contributions:
 - Acknowledgment of accomplishments in team meetings
 - Proper attribution in documentation and presentations
 - Equal speaking time during client presentations

Collaboration and Inclusion

1. Team member skills and expertise:
 - Tyler: Strong mathematical background, algorithm optimization
 - Aidan: Thread programming, FPGA experience
 - Conner: Docker containerization, documentation expertise
 - Joey: Hardware debugging, memory management, testing
2. Strategies for encouraging contributions:
 - Rotating meeting facilitation roles
 - Explicit invitation for input from quieter members
 - Recognition of diverse problem-solving approaches
3. Procedures for identifying collaboration issues:

- Anonymous feedback mechanism via online form
- Regular retrospective meetings to discuss process improvements
- Direct communication with project manager for serious concerns

Goal-Setting, Planning, and Execution

1. Team goals for this semester:
 - Complete mathematical division by Week 8
 - Implement thread management by Week 12
 - Achieve 50% of target performance improvement by Week 16
2. Strategies for planning and assigning work:
 - Task assignment based on skill match and workload balance
 - Weekly sprint planning with clear deliverables
 - Regular progress tracking against milestones
3. Strategies for keeping on task:
 - Weekly progress updates with task burndown charts
 - Peer accountability partnerships
 - Regular demos of implemented functionality

Consequences for Not Adhering to Team Contract

1. Handling infractions:
 - First occurrence: Private conversation with team member
 - Second occurrence: Discussion in team meeting
 - Persistent issues: Consultation with faculty advisor
2. Addressing continued infractions:
 - Redistribution of workload if necessary
 - Revision of responsibilities based on demonstrated reliability
 - In extreme cases, formal notification to course instructor

a) I participated in formulating the standards, roles, and procedures as stated in this contract.

b) I understand that I am obligated to abide by these terms and conditions.

c) I understand that if I do not abide by these terms and conditions, I will suffer the consequences as stated in this contract.

1) Joseph Metzen DATE 4/29/2025

2) Aidan Perry DATE 4/29/2025

3) Tyler Schaefer DATE 4/29/2025

4) Conner Ohnesorge DATE 4/29/2025

