



**VCU**

College of Engineering

# Network Feature Extraction From High Fidelity Cyber Simulation Environment

#025-345

## Project Proposal

Prepared for

DoD Aspire/Brian Cavanagh

DoD

By

Christopher Castro, Ryan Collette, William Lagos, Sam Soltanian

Under the supervision of

Irfan Ahmed

10/11/2024

## Executive Summary

This document outlines an OS device classification system that utilizes machine learning. Rather than focusing on real-time breach detection, this system passively collects and analyzes web traffic, classifying operating systems within the network. By employing AI-driven classification through network feature extraction, we aim to provide enhanced visibility into device behavior, addressing gaps in both research and practical implementation of consumer-level hardware for OS identification. With wireless enabled devices becoming increasingly prevalent across consumer and critical environments, understanding network composition and behavior is crucial for anticipating and managing potential vulnerabilities. This project focuses on creating a scalable, easily replicable software solution capable of classifying devices on a network with minimal setup, operating as a passive extension to the router to streamline device monitoring.

**Note:** The Executive Summary should be updated between major reports as more knowledge is acquired and understanding of the project expands. For example, when submitting Preliminary Design Report in December 2024, make sure you update this page to reflect the progress on the project since the submission of Project Proposal in early October 2024.

## **Table of Contents**

Section A. Problem Statement	5
Section B. Engineering Design Requirements	7
B.1 Project Goals (i.e. Client Needs)	7
B.2 Design Objectives	7
B.3 Design Specifications and Constraints	8
B.4 Codes and Standards	9
Section C. Scope of Work	11
C.1 Deliverables	11
C.2 Milestones	12
C.3 Resources	12
Section D. Concept Generation	13
Section E. Concept Evaluation and Selection	14
Section F. Design Methodology	16
F.1 Computational Methods (e.g. FEA or CFD Modeling, example sub-section)	16
F.2 Experimental Methods (example subsection)	16
F.5 Validation Procedure	16
Section G. Results and Design Details	18
G.1 Modeling Results (example subsection)	18
G.2 Experimental Results (example subsection)	18
G.3 Prototyping and Testing Results (example subsection)	18
G.4. Final Design Details/Specifications (example subsection)	18
Section H. Societal Impacts of Design	20
H.1 Public Health, Safety, and Welfare	20
H.2 Societal Impacts	20
H.3 Political/Regulatory Impacts	20
H.4. Economic Impacts	20
H.5 Environmental Impacts	21
H.6 Global Impacts	21
H.7. Ethical Considerations	21
Section I. Cost Analysis	22

Section J. Conclusions and Recommendations	23
Appendix 1: Project Timeline	24
Appendix 2: Team Contract (i.e. Team Organization)	25
Appendix 3: [Insert Appendix Title]	26
References	27

## **Section A. Problem Statement**

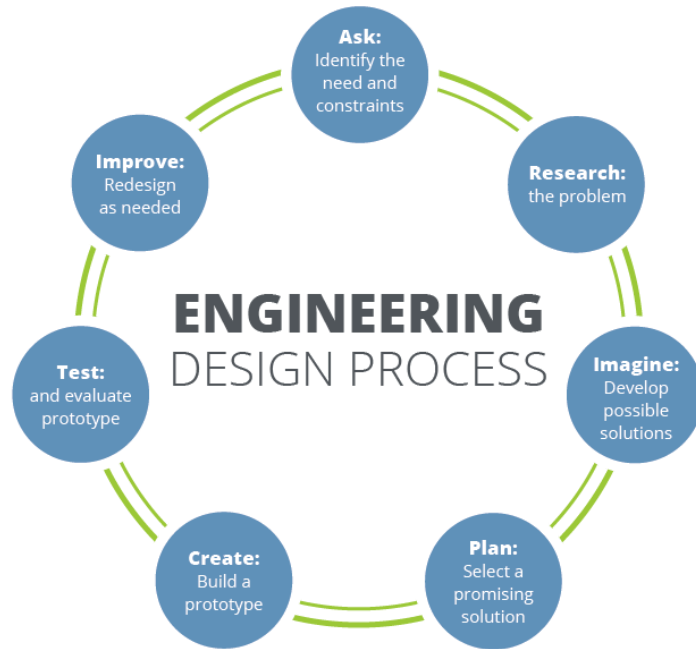
OsirisML is a machine learning model designed to identify the operating system origin of a network request. As network complexity increases and organizations around the world fail to update their systems, we are left with a wide range of possible candidates for where a request originated from. The ability to identify these operating systems using passive network monitoring can prove to be a useful ability in modern warfare. Acquisition of these details provides us with a massive advantage in the enumeration phase while remaining undetected.

By leveraging passive data collection and applying ML classification models, we can extract significant features from TCP IP packets that enable a comprehensive classification system for identifying devices based on network behavior. This non-intrusive model allows us to monitor without alerting potential malicious actors, maintaining a low profile while gathering essential data. This implementation additionally allows for a massive range of expansion in terms of deployment, usage of data acquired from this enumeration phase, and efficiency improvements to current implementations.

Current research in operating system classification implementations remains underexplored, especially in high-fidelity environments where such devices are prevalent. Implementing a low-cost, deployable model that runs on accessible consumer grade hardware allows for easy reproducibility. This model could be seamlessly introduced into a military environment with minimal setup, allowing for a broad scope of traffic capture across various sectors.

Practical challenges include ensuring compatibility with a range of routers and maintaining system resilience against potential countermeasures that may attempt to interfere with passive monitoring. To support the model's complexity, the device may require augmentation with external processing resources or edge computing, adding an external processing layer that introduces new security considerations.

This approach paves the way for advanced network device mapping and monitoring, enriching our ability to understand and manage the dynamic composition of multi-platform environments across conflict and civilian zones alike.



**Figure 1. The iterative nature of the engineering design process [2].**

## **Section B. Engineering Design Requirements**

### **B.1 Project Goals (i.e. Client Needs)**

- **Enhance Existing Model Functionality:**
  - Improve the current OS fingerprinting model to accurately identify not just the operating system type but also specific builds and versions from passive network traffic.
- **Increase Model Accuracy and Robustness:**
  - Utilize underexplored machine learning techniques to enhance the model's performance.
- **Streamline Deployment and Usability:**
  - Simplify the setup, usage, and deployment of the model.

### **B.2 Design Objectives**

- **Improve Model Accuracy:**
  - Implement k-fold cross-validation to enhance model robustness.
  - Perform extensive hyperparameter tuning using advanced optimization tools like Optuna or Hyperopt.
- **Develop a Command-Line Interface (CLI):**
  - Create a user-friendly CLI to simplify interactions with the model, covering functionalities from data preprocessing to prediction.
- **Containerize the Application:**
  - Package the model and its dependencies into a Docker container to streamline setup and ensure consistent deployment across different environments.
- **Expand OS Detection Capabilities:**
  - Extend the model's functionality to identify specific OS builds and versions.

### **B.3 Design Specifications and Constraints**

- The model will retain at least the same level of accuracy in classifying the original test set. (Functional constraint).
- The model will use K-fold cross validation to improve accuracy (Functional constraint)
- The model will be containerized to allow for ease of use on other platforms (Functional constraint)
- A user will be able to interact with the model using a command line tool (Functional constraint)

### **B.4 Codes and Standards**

- **IEEE P7001:** Transparency in AI—document the decision-making process of the AI.

## **Section C. Scope of Work**



The project scope defines the boundaries of the project encompassing the key objectives, timeline, milestones and deliverables. It clearly defines the responsibility of the team and the process by which the proposed work will be verified and approved. A clear scope helps to facilitate understanding of the project, reduce ambiguities and risk, and manage expectations. In addition to stating the responsibilities of the team, it should also explicitly state those tasks which fall *outside* of the team's responsibilities. *Explicit bounds* on the project timeline, available funds, and promised deliverables should be clearly stated. These boundaries help to avoid *scope creep*, or changes to the scope of the project without any control. This section also defines the project approach, the development methodology used in developing the solution, such as waterfall or agile (shall be chosen in concert with the faculty advisor and/or project sponsor). Good communication with the project sponsor and faculty advisor is the most effective way to stay within scope and make sure all objectives and deliverables are met on time and on budget.

**Project Timeline:**

- **November - December 2024:**
  - **Finalize team roles, project objectives, and deliverables.**
  - **Conduct initial research and begin drafting the preliminary design report.**
  - **Receive the new dataset from the sponsor and perform initial data preprocessing.**
- **January - February 2025:**
  - **Model development with k-fold cross-validation.**
  - **Extensive hyperparameter tuning using advanced optimization techniques.**
- **March 2025:**
  - **CLI application development and integration.**
  - **Containerization of the application and system integration.**
- **April 2025:**
  - **Documentation and user guides development.**
  - **Capstone EXPO preparation, including final reports and presentations.**

**We will be following Agile methodology for development, so we can ensure each component is thoroughly tested and improved. This allows for flexibility and continuous feedback from the project sponsor and faculty advisor.**

### **C.1 Deliverables**

The project deliverables are those things that the project team is responsible for providing to the project sponsor. They are the things that are to be produced or provided as a result of the engineering design process. Some deliverables might include a specific number of alternative designs, required analyses to prove the design meets specifications, detailed machine drawings, functional diagrams or schematics, required computer code, flow charts, user manuals, desktop models, and functioning prototypes. A design “proof of concept” is not specific and should be more clearly defined. Academic deliverables include the team contract, project proposal,

preliminary design report, fall poster and presentation, final design report, and Capstone EXPO poster and presentation. Provide a bulleted list of all agreed upon project deliverables.

In order to mitigate risks associated with the completion and delivery of the project deliverables, provide an outline of the most potentially disruptive, foreseeable obstacles. Some important issues to discuss with the design team, sponsor, and faculty advisor include the following:

- What deliverables require access to campus? Which/how many students regularly access campus and are physically available to complete tasks?
- What work can be done remotely? What resources might be needed in order to ensure that remote work can be completed effectively (e.g. software licenses, shared drives/folders, etc.)?
- What deliverables require ordering from third-party vendors? Will any components potentially required extended lead times? What can the team do in order to mitigate potential supply chain disruptions?

**The project will produce the following deliverables:**

- 1. Enhanced Machine Learning Model:**
  - A trained model capable of detecting specific OS minor versions and builds using the new dataset provided by the sponsor.
  - Implementation of k-fold cross-validation to improve model robustness.
  - Extensive hyperparameter tuning, including exploration of boost types, max\_depth, gamma, learning rate, eta, and others.
- 2. Unified Command-Line Interface (CLI) Application:**
  - A cohesive CLI that integrates all scripts into a user-friendly application.
  - Improved usability with clear documentation and help messages.
  - Error handling, logging, and configuration management features.
- 3. Containerized Application:**
  - A Docker container of the OsirisML application for easy deployment across various operating systems, not limited to Debian Linux.
  - Dockerfile and instructions for building and running the containerized application.

**Potential Risks and Mitigation:**

- **Data Quality and Availability:**
  - *Risk:* The new dataset may have limitations or require additional preprocessing.
  - *Mitigation:* Close collaboration with the sponsor to understand the dataset and adjust preprocessing steps accordingly.
- **Computational Resource Constraints:**

- **Risk:** Hyperparameter tuning and k-fold cross-validation may be computationally intensive.
- **Mitigation:** Utilize university computational resources or cloud services.
- **Time Management:**
  - **Risk:** Potential delays due to unforeseen technical challenges.
  - **Mitigation:** Regular progress reviews and agile adjustments to the project plan.

## C.2 Milestones

Milestones are major project phases or tasks that need to be completed in order to ensure the project deliverables. They may include, among other things, completion of calculations, the development of a computational model, completion of an analysis, set-up of an experiment, completion of data acquisition, purchasing of hardware, assembly of a prototype, completion of testing procedures, development of required code, completion of wiring, post processing, etc.

A good rule of thumb is to break the project down into tasks of no larger than 2-3 weeks in length. These can be individual or group tasks. Breaking down the project into tasks/milestones gives the team and the advisor/sponsor a realistic understanding of what can be done in the allotted time. In an agile development approach, later tasks are expected to be adjusted (or changed) as the team works with the earlier developed tasks.

The amount of time it will take to accomplish each milestone and the approximate date that each milestone will be completed should be considered. Do not underestimate the time that it takes to write and prepare major reports and presentation materials. All deliverables and milestones should be included in the project timeline found in Appendix 1. Provide a summary table of all project milestones including required times and completion dates here.

**Note:** While the project scope, deliverable, and milestones are not intended to change throughout the project, this section should be revisited between major reports to ensure that it still accurately reflects the expectations and requirements of the project team, client, and faculty advisor. Any changes to the project scope, deliverable, and milestones should be thoroughly discussed and mutually agreed upon by all parties. Any changes to this section should be documented and justified in detail.

### Project Milestones Timeline:

#### Finalize Team Roles, Project Objectives, and Deliverables

- **Tasks:**
  - Finalize team roles, project objectives, scope, and deliverables.
  - Submit the project proposal to VCU.
  - Prepare and submit the fall design poster.
- **Expected Completion: November 15, 2024**

---

## 2. Resource Procurement and Research

- **Tasks:**
  - **Resource Procurement:**
    - Contact VCU IT services to procure a high-performance virtual machine (VM) with at least 128 GB of RAM, accessible by all team members.
      - If VCU IT is unable to provide, research other alternatives such as AWS EC2 or GCP.
  - **Research:**
    - Conduct in-depth research on:
      - Advanced machine learning techniques for OS version detection.
      - Feature extraction methods relevant to OS fingerprinting and distinguishing specific OS versions.
      - K-Fold Cross-Validation and its impact on improving model accuracy and robustness.
      - Hyperparameter optimization techniques, including tools like GridSearchCV, RandomizedSearchCV, and Optuna.
      - Boosting algorithms and their variants to understand how different boost types affect model performance.
    - Begin drafting the **Preliminary Design Report**.
- **Expected Completion: November 29, 2024** (2 weeks)

---

## 3. Getting Model to Run Locally

- **Tasks:**
  - Set up development environments on procured machines.
  - Adjust existing scripts to accommodate the new dataset structure.
  - Ensure the initial model runs successfully on local machines using a subset of the data.
  - Begin initial data preprocessing and feature extraction.
  - Continue working on the **Preliminary Design Report** and complete it.
- **Expected Completion: December 13, 2024** (2 weeks)

---

## 4. Data Preprocessing and Feature Extraction

- **Tasks:**
  - Perform comprehensive data preprocessing and feature extraction, including:
    - Adjusting preprocessing scripts to fully accommodate the new dataset structure.

- Extracting features relevant to specific OS version detection.
    - Cleaning and normalizing the data to prepare it for model training.
    - Validate the processed data and ensure it's ready for model development.
  - **Expected Completion: December 27, 2024** (2 weeks)
- 

## 5. Model Development with K-Fold Cross-Validation

- **Tasks:**
    - Implement **k-fold cross-validation** in the training pipeline to enhance model robustness.
    - Train initial machine learning models using the preprocessed dataset.
    - Analyze model performance metrics to identify areas for improvement.
    - Document initial findings and update the design report accordingly.
  - **Expected Completion: January 31, 2025** (5 weeks)
- 

## 6. Extensive Hyperparameter Tuning

- **Tasks:**
    - Perform comprehensive hyperparameter tuning using advanced optimization techniques.
    - Explore and optimize parameters such as:
      - **Boost Types** (e.g., gbtrees, gblines, dart).
      - **Max Depth.**
      - **Gamma.**
      - **Learning Rate (Eta).**
      - **Min Child Weight.**
      - **Subsample and Colsample\_bytree.**
    - Utilize tools like **Optuna** or **Hyperopt** for efficient hyperparameter search.
    - Record the optimal hyperparameters and their impact on model performance.
    - Update documentation with tuning results and insights.
  - **Expected Completion: February 28, 2025** (4 weeks)
- 

## 7. CLI Application Development

- **Tasks:**
  - Refactor existing scripts into a modular and maintainable codebase.
  - Develop a unified **Command-Line Interface (CLI)** using Python libraries such as **argparse** or **click**.
  - Implement user-friendly features including:

- Clear help messages and usage instructions.
    - Input validation and informative error messages.
    - Logging mechanisms for tracking application processes.
  - Ensure the CLI covers all functionalities, from data preprocessing to model training and testing.
  - Begin integrating the CLI with the developed models.
  - **Expected Completion: March 21, 2025** (3 weeks)
- 

## 8. Containerization of the Application and System Integration

- **Tasks:**
    - Create a **Dockerfile** to containerize the OsirisML application.
    - Ensure all dependencies and environment configurations are included.
    - Integrate all components into a single, containerized application.
    - Perform extensive testing, including:
      - **Unit Tests** for individual components.
      - **Integration Tests** to ensure modules work together seamlessly.
      - **System Tests** to validate the entire application's functionality.
    - Test the containerized application on different operating systems to verify compatibility.
    - Optimize the Docker image for minimal size and efficient performance.
    - Provide documentation on how to build and run the Docker container.
  - **Expected Completion: April 11, 2025** (3 weeks)
- 

## 9. Documentation and User Guides

- **Tasks:**
    - Prepare comprehensive technical documentation, including:
      - Installation and setup instructions.
      - Detailed usage guides for all application features.
      - Developer documentation for future maintenance and enhancements.
    - Develop user manuals with examples and best practices.
    - Update the design report with finalized designs, methodologies, and findings.
  - **Expected Completion: April 18, 2025** (1 week)
- 

## 10. Capstone EXPO Preparation

- **Tasks:**
  - Prepare for the Capstone EXPO by:

- Creating the poster presentation.
  - Writing the final project report.
  - Developing presentation materials and rehearsing the presentation.
- Seek feedback from the project sponsor and faculty advisor to refine deliverables.
- **Expected Completion: April 25, 2025** (1 week)

### C.3 Resources

Resources needed for project completion should be listed at the proposal stage. These resources can either be purchased within the Project Budget, or provided by the project sponsor. Some examples are: hardware such as HPCs or servers, software such as IDEs, data analysis platforms or version control systems. Access to cloud computing services may also be necessary to scale certain procedures. Additionally, databases containing operational data for testing, as well as libraries or APIs relevant to predictive analytics and machine learning may be required.

#### Hardware:

- **Computational Resources:**
  - **High-Performance Computing (HPC) Clusters:**
    - Access to university-provided HPC clusters or servers for computationally intensive tasks such as extensive hyperparameter tuning and k-fold cross-validation on large datasets.
    - Necessary during the **Extensive Hyperparameter Tuning** phase and when training models with k-fold cross-validation to handle increased computational load.

#### Software:

- **Programming Languages and Development Tools:**
  - **Python 3.x:**
    - The primary programming language for all development tasks.
    - Utilized extensively across all project phases.
  - **Integrated Development Environments (IDEs):**
    - **PyCharm**, or **Visual Studio Code** for code development and debugging.
  - **Version Control Systems:**
    - **Git** for version control.
    - **GitHub** for code hosting, collaboration, and issue tracking.

#### Machine Learning and Data Processing Libraries:

- **XGBoost:**
  - For model training, evaluation, and implementing advanced boosting algorithms.

- Central to the **Model Development** and **Hyperparameter Tuning** phases.
- **Scikit-learn:**
  - For implementing k-fold cross-validation, initial model development, and basic machine learning tasks.
  - Used during the **Model Development with K-Fold Cross-Validation** phase.
- **Optuna or Hyperopt:**
  - Advanced hyperparameter optimization frameworks to efficiently search for optimal model parameters.
  - Employed during the **Extensive Hyperparameter Tuning** phase.
- **Pandas and NumPy:**
  - For data manipulation, preprocessing, and analysis.
  - Essential during the **Data Preprocessing** and throughout model development.
- **Scapy or Tshark (if needed):**
  - For advanced packet analysis and feature extraction from PCAP files.
  - May be required during **Dataset Acquisition and Data Preprocessing** if additional feature extraction is necessary.

#### **Command-Line Interface (CLI) Development Libraries:**

- **argparse, click, or typer:**
  - For building a user-friendly CLI application.
  - Integral to the **CLI Application Development** phase.

#### **Containerization Tools:**

- **Docker:**
  - For containerizing the application to ensure consistent deployment across different environments and operating systems.
  - Central to the **Containerization of the Application and System Integration** phase.

#### **Testing Frameworks:**

- **PyTest:**
  - For writing and running unit tests.
  - Employed during the **Containerization and System Integration** phase.

#### **Data Resources:**

- **Labeled Dataset Provided by the Sponsor:**
  - A new dataset containing network traffic data with specific OS version labels for training and testing the machine learning models.



- Central to the **Dataset Acquisition and Data Preprocessing** and subsequent modeling phases.

### **Compute Services:**

- **Cloud Computing Platforms (if required):**
  - **Amazon Web Services (AWS) EC2, Google Cloud Platform (GCP), or Microsoft Azure:**
    - For scalable computing resources, especially during extensive hyperparameter tuning and model training.
    - Backup option if university HPC resources are insufficient or unavailable.

## **Section D. Concept Generation**

The primary objective of these concepts is to improve the accuracy of the existing Machine Learning model used to identify operating systems through network feature extraction. By exploring different avenues for feature selection, dimensionality reduction, and model optimization, we aim to enhance the model's ability to analyze and classify network traffic.

### **Concept 1: Feature Importance Analysis with XGBoost**

Use XGBoost to identify the most important features contributing to the model's accuracy. By focusing on high-impact features, we can streamline the model and potentially reduce noise caused by irrelevant or redundant features.

How it addresses the problem:

- Pinpoints features that significantly affect classification accuracy.
- Helps reduce overfitting by limiting the feature space.

Pros:

- XGBoost is well-documented and has built-in support for feature importance scoring.
- Can be directly integrated into the existing model workflow.

Cons:

- Requires fine-tuning of hyperparameters for optimal performance.

Potential Risks:

- Important features may still interact in ways that XGBoost alone can't capture.

## **Concept 2: K-Fold Cross Validation with PCA Variable Selection**

Integrate Principal Component Analysis to reduce dimensionality, followed by k-fold cross validation to validate the model on multiple splits of the data. This ensures robust performance and helps evaluate how well the model generalizes.

How it addresses the problem:

- Reduces dimensionality, simplifying the model and reducing computational overhead.
- Provides more robust and less variance-prone accuracy estimates.
- Efficiently uses all available data by training and evaluating on the entire set..

Pros:

- PCA highlights key latent variables that may improve model interpretability.

Cons:

- PCA can sometimes remove variables that have nonlinear contributions to the output.

Potential Risks:

- Risk of underfitting if too much information is lost during dimensionality reduction.

## **Concept 3: Model Tuning with CatBoost**

Implement CatBoost as an alternative to the existing model, taking advantage of its ability to handle categorical features natively and reduce overfitting through advanced regularization techniques.

How it addresses the problem:

- CatBoost can capture complex interactions in data that other algorithms might miss.
- CatBoost is especially effective with noisy datasets, common in network traffic analysis.

Pros:

- Faster training time compared to some other ensemble methods.
- Supports GPU acceleration for large datasets.

Cons:

- Original project was built with XGBoost rather than CatBoost, so there is a learning curve for team members and sponsors unfamiliar with CatBoost.

Potential Risks:

- CatBoost may require significant computational resources for tuning.



## Section E. Concept Evaluation and Selection

Evaluated three concepts using the following criteria:

### Feature Importance Analysis with XGBoost

- By optimizing model features using Feature Importance Analysis we are able to pinpoint issues which prevent the model from obtaining a high accuracy by filtering features which may cause excess bloat and memory usage.
- This workflow can be directly added to the process as it is a feature shipped alongside XGBoost.

### K-Fold Cross Validation with PCA Variable Selection

- PCA Variable Selection reduces computational complexity and mitigates overfitting by focusing on the most informative features.
- K-Fold Cross Validation ensures more reliable performance estimates and stronger generalization by testing the model on multiple data splits.

### Model Tuning with CatBoost

- Projected to performance post training which can be measured by both testing the model against datasets or referencing the confusion matrix.
- Reduces overfitting and improves model generalization, leading to more reliable performance.
- CatBoost has highly optimized default parameters, often providing strong results with minimal tuning.
- CatBoost is not often used as it is found to be slower than LightGBM, but with its implementation of gradient boosting can improve performance accuracy with a large amount of features

Weight	Feature Importance Analysis with XGBoost	K-Fold Cross Validation with PCA Variable Selection	Model Tuning with CatBoost
Accuracy Improvement	0.6	0.2	0.2
Compute Usage Improvement	0.4	0.4	0.2
Overfitting Reduction	0.2	0.3	0.5
Total	1.1	0.9	0.9

## Section F. Design Methodology

### F.1 System Architecture

The system consists of two main components:

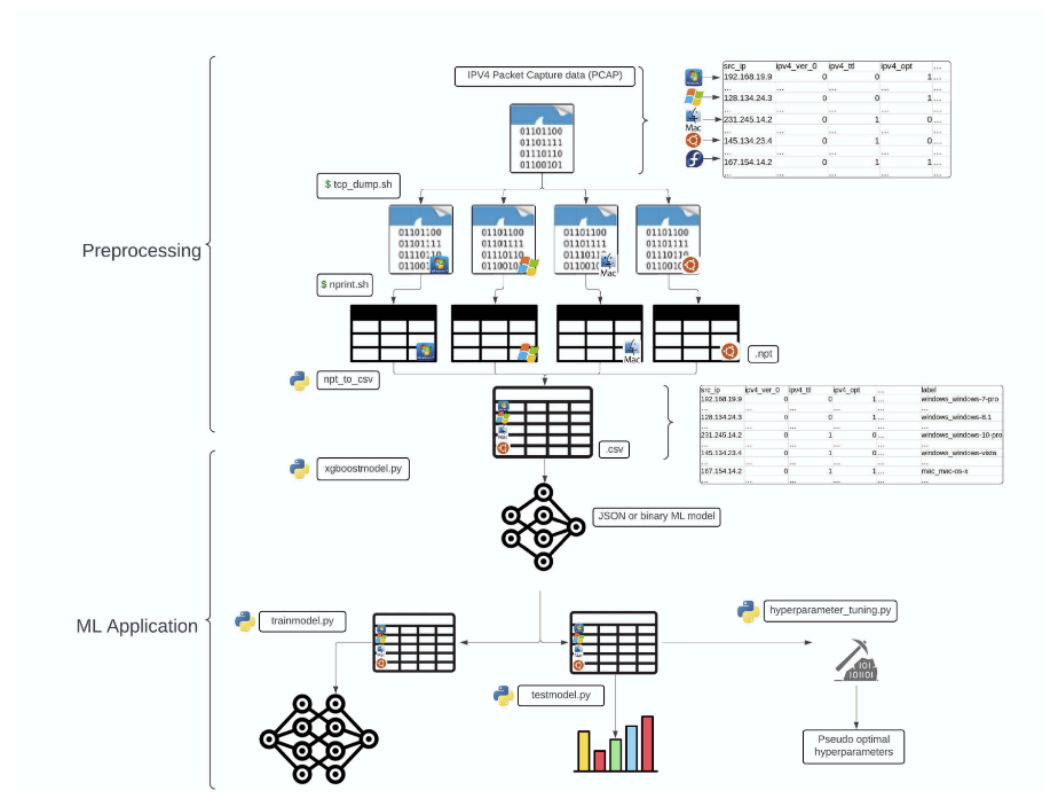
#### Data Preprocessing Pipeline

- Extracts data from packet-capture files using bash
- Divides data into data equivalent to that required from nPrint by operating system and labels the data.
- Converts labeled data to csv format for training the model.
- Performs k-fold cross validation on data for proper randomization.

#### Classification Engine

- Trains the XGBoost model and extracts features from the data by extracting csv data from preprocessing pipeline.
- Produces performance metrics and confusion matrix for training and testing data.
- Performs classification of network traffic to operating system identification.

### F.2 Data Flow



(Credit to original authors of the OsirisML paper Spencer Ekeroth, Jeremy Neale, and Jae Sung Kim for this data flow diagram)

## **F.3 Validation Procedure**

### **Testing Dataset**

- Create or obtain testing data with diverse operating systems and traffic patterns. Alternatively use segments of the CISC-IDS 2017 intrusion dataset.
- Test the model against randomized data to gauge true effectiveness of its classification capabilities and capabilities in identifying versioned operating system builds.
- Validate against known OS configurations which are prone to appearing in a wide variety of system configurations.

### **Performance Metrics**

- Confusion Matrix Accuracy Output
- Utilization of System Resource Improvement
- Training & Model Response Time

### **Deployment Testing**

- Testing containerized lightweight models on different operating systems.
- Testing model against system specification constraints to find true compute requirements.
- Testing limits the model's ability to handle larger datasets in an efficient manner.

## **Section G. Results and Design Details**

### **G.1 Initial Implementation**

- Complete initial lightweight model implementation for testing at decreased time requirement.
- Configuration and development of data preprocessing pipeline and randomization via k-fold cross validation.

### **G.2 Model Development**

- Begin mining and testing for hyperparameters
- Retrain model on current and requested data in randomized format for performance testing with a goal of ~10%+ accuracy improvement and the ability to distinguish between versions of operating systems (e.g. Ubuntu 20.04 vs 24.04)
- Configure and test model's features to prevent overfitting by XGBoost
- Provide a command line interface for interaction with the model

### **G.3 Resource Usage**

- Reduce Overall Memory Footprint: ~100-200gb
- Optimize model to run in a containerized environment and operating systems other than Debian.

## **Section H. Societal Impacts of Design**

### **H.1 Public Health, Safety, and Welfare**

- Enhances military capability by automating the enumeration phase of offensive or defensive cybersecurity and network mapping.
- Easily identifies systems in these environments which may be susceptible to cyber attacks.
- Allows military forces to obtain a mapping of the network around them to perform action as needed.

### **H.2 Societal Impacts**

- Risks misuse for surveillance or malicious purposes.

### **H.3 Political/Regulatory Impacts**

- Does not notify any group which the model could be used on that collection of data and classification of their systems is taking place. This could eventually lead to a regulatory issue if used by anyone outside of the intended environment.
- Implications for network security standards and passive identification of network traffic.

### **H.4 Economic Impacts**

- Reduces costs associated with alternative methods of in-field enumeration phases.
- Enables efficient resource allocation in military environments.
- Significantly reduces required time to map and understand environments with network traffic and could possibly decrease required payroll.



## Section I. Cost Analysis

Provide a simple cost analysis of the project that includes a list of all expenditures related to the project. If an experimental test set-up or prototype was developed, provide a Bill of Materials that includes part numbers, vendor names, unit costs, quantity, total costs, delivery times, dates received, etc. Do not forget to include all manufacturing costs incurred throughout the completion of the project. If the design is expected to become a commercial product, provide a production cost estimate including fixed capital, raw materials, manufacturing (including tooling and/or casting), and labor costs to produce and package the device. Note that this type of detailed cost analysis may be listed as a project deliverable.

**Note:** The Preliminary Design Report should include all costs incurred to date. It is expected that this section will be expanded and updated between the preliminary and final design reports.

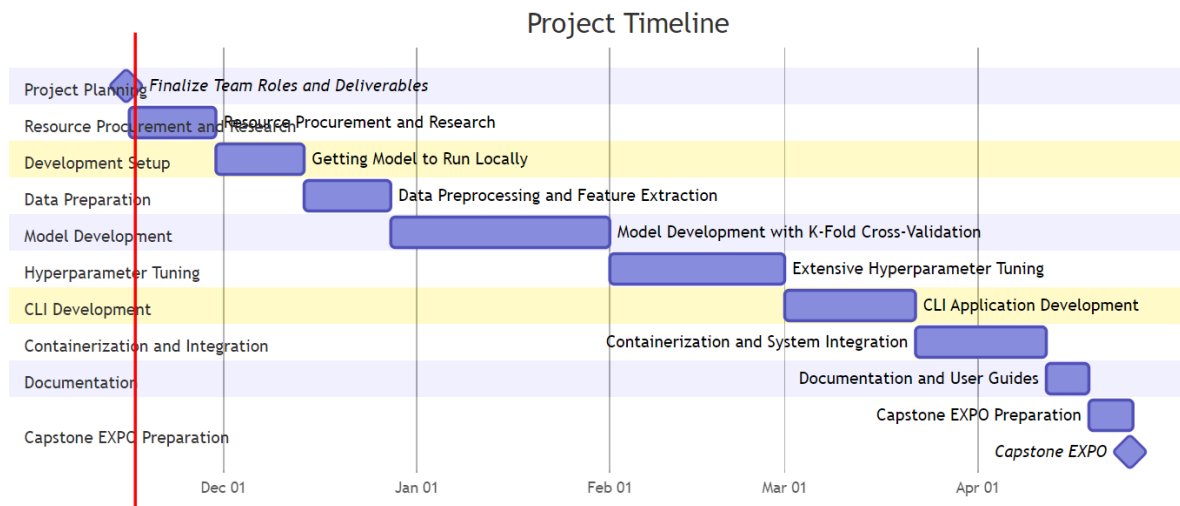
## **Section J. Conclusions and Recommendations**

Use this section to summarize the story of how the design team arrived at the final design. Focus on the evolution of the design through the use of the engineering design process including lessons learned, obstacles overcome, and triumphs of the final design. Revisit the primary project goals and objectives. Provide a brief summary of the final design details and features paramount to the function of the design in meeting these goals and objectives.

A discussion may be included to discuss how the design could be further advanced or improved in the future. If applicable, summarize any questions or curiosities that the final results/design of this effort bring to mind or leave unanswered. If this project might continue on as a future (continuation) senior design project, detail the major milestones that have been completed to date and include any suggested testing plans, relevant machine drawings, electrical schematics, developed computer code, etc. All relevant information should be included in this section such that future researchers could pick up the project and advance the work in as seamless a manner as possible. Documents such as drawings, schematics, and codes could be referenced here and included in one or more appendix. If digital files are critical for future work, they should be saved on a thumb drive, external hard drive, cloud, etc. and left in the hands of the project advisor and/or client.

## Appendix 1: Project Timeline

Provide a Gantt chart of similarly composed visual timeline showing the start and end dates of all completed tasks and how they are grouped together, overlapped, and linked together. Include all senior design requirements including design reports and Expo materials (i.e. Abstract, Poster, and Presentation). All major milestones should be included in the timeline.



## Appendix 2: Team Contract (i.e. Team Organization)

Copy and paste the content from the completed Team Contract here starting with Step 1 of the Team Contract and including all content following the ‘Contents’ list.

### Step 1: Get to Know One Another. Gather Basic Information.

**Task:** This initial time together is important to form a strong team dynamic and get to know each other more as people outside of class time. Consider ways to develop positive working relationships with others, while remaining open and personal. Learn each other’s strengths and discuss good/bad team experiences. This is also a good opportunity to start to better understand each other’s communication and working styles.

<b><i>Team Member Name</i></b>	<b><i>Strengths each member bring to the group</i></b>	<b><i>Other Info</i></b>	<b><i>Contact Info</i></b>
<i>William Lagos</i>	<i>Technical skills, flexibility, problem solving, teamwork</i>		<i>lagoswj@vcu.edu</i>
<i>Sam Soltanian</i>	<i>Always trying to learn, problem-solving flexible</i>	<i>I enjoy learning from other people and in turn offering new knowledge</i>	<i>soltanians@vcu.edu</i>
<i>Christopher Castro</i>	<i>Friendly, curious, bias for action, hard worker</i>	<i>Unfamiliar with cybersecurity topics / PCAP data, but eager to fill in knowledge gaps.</i>	<i>castrocm5@vcu.edu</i>
<i>Ryan Collette</i>	<i>Always want to learn, willing to adapt to complete goals, enjoy working with a team.</i>	<i>Lack of in depth cyber-security concepts, but surface level knowledge and very interested in learning.</i>	<i>colletterc@vcu.edu</i>

<b><i>Other Stakeholders</i></b>	<b><i>Notes</i></b>	<b><i>Contact Info</i></b>
<i>Project Sponsor</i>	<i>Meeting Friday's @ 12:30PM @ Teams</i>	<i>Jae Sung Kim jaesung.kim6.civ@army.mil</i>
<i>Staff Mentor</i>	<i>Initial Meeting Wed, Sep. 4 @ Zoom</i>	<i>Irfan Ahmed iahmed3@vcu.edu</i>

## Step 2: Team Culture. Clarify the Group's Purpose and Culture Goals.

**Task:** Discuss how each team member wants to be treated to encourage them to make valuable contributions to the group and how each team member would like to feel recognized for their efforts. Discuss how the team will foster an environment where each team member feels they are accountable for their actions and the way they contribute to the project. These are your Culture Goals (left column). How do the students demonstrate these culture goals? These are your Actions (middle column). Finally, how do students deviate from the team's culture goals? What are ways that other team members can notice when that culture goal is no longer being honored in team dynamics? These are your Warning Signs (right column).

**Resources:** More information and an example Team Culture can be found in the Biodesign Student Guide "Intentional Teamwork" page ([webpage](#) | [PDF](#))

<i><b>Culture Goals</b></i>	<i><b>Actions</b></i>	<i><b>Warning Signs</b></i>
<i>1: Punctuality</i>	<ul style="list-style-type: none"><li>- <i>Students will arrive to meetings on time</i></li><li>- <i>Students need to communicate any needs to cancel or reschedule meetings</i></li></ul>	<ul style="list-style-type: none"><li>- <i>Student misses with no explanation first meeting, warning is granted</i></li><li>- <i>Student misses meetings with no explanation afterwards – issue is brought up with faculty advisor</i></li></ul>
<i>2: Informing the group of any delays in completing assignments</i>	<ul style="list-style-type: none"><li>- <i>Stay up to date with each other's project responsibilities</i></li><li>- <i>Set reasonable deadlines and note when an extension is needed</i></li></ul>	<ul style="list-style-type: none"><li>- <i>Student shows up for weekly meeting with no considerable work done</i></li></ul>
<i>3: Be open to critical feedback and receptive</i>	<ul style="list-style-type: none"><li>- <i>Be receptive to ideas that potentially differ from your own.</i></li><li>- <i>Be open to alternative design options that you may be unfamiliar with.</i></li><li>- <i>Respond to feedback in a mature and normal manner.</i></li></ul>	<ul style="list-style-type: none"><li>- <i>Students are unreceptive/immature about feedback.</i></li></ul>

### Step 3: Time Commitments, Meeting Structure, and Communication

**Task:** Discuss the anticipated time commitments for the group project. Consider the following questions (don't answer these questions in the box below):

- What are reasonable time commitments for everyone to invest in this project?
- What other activities and commitments do group members have in their lives?
- How will we communicate with each other?
- When will we meet as a team? Where will we meet? How Often?
- Who will run the meetings? Will there be an assigned team leader or scribe? Does that position rotate or will the same person take on that role for the duration of the project?

**Required:** How often you will meet with your faculty advisor, where you will meet, and how the meetings will be conducted. Who arranges these meetings?

See examples below.

<i>Meeting Participants</i>	<i>Frequency Dates and Times / Locations</i>	<i>Meeting Goals Responsible Party</i>
<i>Students Only</i>	<i>As Needed, in-person or on Discord based on requirements per week. Scheduled based on discussion every Thursday.</i>	<i>Actively work on project (Group will document these meetings by taking photos of whiteboards, physical prototypes, etc, then post on Discord and update Capstone Report )</i>
<i>Students Only</i>	<i>Every Thursday during the seminar time block</i>	<i>Update group on weekly challenges and accomplishments, scheduling as needed meetings for other time slots during the week.</i>
<i>Students + Faculty advisor</i>	<i>Every Friday at 12:30 PM on Zoom</i>	<i>Update faculty advisor and get answers to our questions, updates following agile methodology (what we've accomplished, what we're blocked on, what we're doing next)</i>
<i>Project Sponsor</i>	<i>Every Friday at 12:30 PM on Zoom</i>	<i>Update project sponsor and make sure we are on the right track</i>

## Step 4: Determine Individual Roles and Responsibilities

**Task:** As part of the Capstone Team experience, each member will take on a leadership role, *in addition to* contributing to the overall weekly action items for the project. Some common leadership roles for Capstone projects are listed below. Other roles may be assigned with approval of your faculty advisor as deemed fit for the project. For the entirety of the project, you should communicate progress to your advisor specifically with regard to your role.

- **Before meeting with your team**, take some time to ask yourself: what is my “natural” role in this group (strengths)? How can I use this experience to help me grow and develop more?
- **As a group**, discuss the various tasks needed for the project and role preferences. Then assign roles in the table on the next page. Try to create a team dynamic that is fair and equitable, while promoting the strengths of each member.

### Communication Leaders

**Suggested:** Assign a team member to be the primary contact for the client/sponsor. This person will schedule meetings, send updates, and ensure deliverables are met.

**Suggested:** Assign a team member to be the primary contact for faculty advisor. This person will schedule meetings, send updates, and ensure deliverables are met.

### Common Leadership Roles for Capstone

1. **Project Manager:** Manages all tasks; develops overall schedule for project; writes agendas and runs meetings; reviews and monitors individual action items; creates an environment where team members are respected, take risks and feel safe expressing their ideas.  
**Required:** On Edusourced, under the Team tab, make sure that this student is assigned the Project Manager role. This is required so that Capstone program staff can easily identify a single contact person, especially for items like Purchasing and Receiving project supplies.
2. **Logistics Manager:** coordinates all internal and external interactions; lead in establishing contact within and outside of organization, following up on communication of commitments, obtaining information for the team; documents meeting minutes; manages facility and resource usage.
3. **Financial Manager:** researches/benchmarks technical purchases and acquisitions; conducts pricing analysis and budget justifications on proposed purchases; carries out team purchase requests; monitors team budget.
4. **Systems Engineer:** analyzes Client initial design specification and leads establishment of product specifications; monitors, coordinates and manages integration of sub-systems in the prototype; develops and recommends system architecture and manages product interfaces.
5. **Test Engineer:** oversees experimental design, test plan, procedures and data analysis; acquires data acquisition equipment and any necessary software; establishes test protocols and schedules; oversees statistical analysis of results; leads presentation of experimental finding and resulting recommendations.

<b><i>Team Member</i></b>	<b><i>Role(s)</i></b>	<b><i>Responsibilities</i></b>
William Lagos	Systems Engineer	Establish project specifications, monitor, coordinate and manages integration of subsystems in the prototype, and develop/recommend systems for use in product
Christopher Castro	Project Manager	Write agendas and run meetings, review and monitor individual action items. Come up with major stories and set up agile board, organize action items.
Sam Soltanian	Logistics Manager	Communicate with Faculty Advisor & Advisor and communicate between team and faculty/sponsor. Make sure everyone is on the same page.
Ryan Collette	Test Engineer	Collaborate and plan experimental design, plan testing frameworks, scheduling testing, statistical analysis of training results or confidence weights.

## Step 5: Agree to the above team contract

*Team Member: Sam S*

*Signature: Sam Soltanian*

*Team Member: Ryan C*

*Signature: Ryan Collette*

*Team Member: Chris C*

*Signature: Christopher Castro*

*Team Member: William L*

*Signature: William Lagos*



### **Appendix 3: [Insert Appendix Title]**

Note that additional appendices may be added as needed. Appendices are used for supplementary material considered or used in the design process but not necessary for understanding the fundamental design or results. Lengthy mathematical derivations, ancillary results (e.g. data sets, plots), and detailed mechanical drawings are examples of items that might be placed in an appendix. Multiple appendices may be used to delineate topics and can be labeled using letters or numbers. Each appendix should start on a new page. Reference each appendix and the information it contains in the main text of the report where appropriate.

**Note:** Delete this page if no additional appendices are included.

## References

Provide a numbered list of all references in order of appearance using APA citation format. The reference page should begin on a new page as shown here.

- [1] VCU Writing Center. (2021, September 8). *APA Citation: A guide to formatting in APA style*. Retrieved September 2, 2024. <https://writing.vcu.edu/student-resources/apa-citations/>
- [2] Teach Engineering. *Engineering Design Process*. TeachEngineering.org. Retrieved September 2, 2024. <https://www.teachengineering.org/populartopics/designprocess>