

# ML-NIC: Accelerating Machine Learning Inference using Smart Network Interface Cards

Raghav Kapoor<sup>1</sup>, David C. Anastasiu<sup>2</sup> and Sean Choi<sup>1\*</sup>

<sup>1</sup>Cloud Lab, Department of Computer Science and Engineering, Santa Clara University, Santa Clara, CA, United States

<sup>2</sup>Anastasiu Lab, Department of Computer Science and Engineering, Santa Clara University, Santa Clara, CA, United States

Correspondence\*:

Sean Choi, David C. Anastasiu  
sean.choi@scu.edu, danastasiu@scu.edu

## 2 ABSTRACT

3 Low-latency inference for machine learning models is increasingly becoming a necessary  
4 requirement, as these models are used in mission-critical applications such as autonomous  
5 driving, military defense (e.g., target recognition), and network traffic analysis. A widely studied  
6 and used technique to overcome this challenge is to offload some or all parts of the inference  
7 tasks onto specialized hardware such as graphic processing units. More recently, offloading  
8 machine learning inference onto programmable network devices, such as programmable network  
9 interface cards or a programmable switch, is gaining interest from both industry and academia,  
10 especially due to the latency reduction and computational benefits of performing inference directly  
11 on the data plane where the network packets are processed. Yet, current approaches are relatively  
12 limited in scope, and there is a need to develop more general approaches for mapping offloading  
13 machine learning models onto programmable network devices.

14 To fulfill such a need, this work introduces a novel framework, called ML-NIC, for deploying  
15 trained machine learning models onto programmable network devices' data planes. ML-NIC  
16 deploys models directly into the computational cores of the devices to efficiently leverage the  
17 inherent parallelism capabilities of network devices, thus providing huge latency and throughput  
18 gains. Our experiments show that ML-NIC reduced inference latency by at least  $6\times$  on average  
19 and in the 99th percentile and increased throughput by at least  $16x$  with little to no degradation in  
20 model effectiveness compared to the existing CPU solutions. In addition, ML-NIC can provide  
21 tighter guaranteed latency bounds in the presence of other network traffic with shorter tail  
22 latencies. Furthermore, ML-NIC reduces CPU and host server RAM utilization by 6.65% and  
23 320.80 MB. Finally, ML-NIC can handle machine learning models that are  $2.25\times$  larger than the  
24 current state-of-the-art network device offloading approaches.

25 **Keywords:** Machine Learning, SmartNIC, Netronome, Data Plane, Inference

## 1 INTRODUCTION

26 Machine learning (ML) permeates a vast amount of everyday life, from personalized recommendations  
27 to stock market analysis and novel drug synthesis. While the machine learning models created to solve  
28 problems in these various fields are proven to be highly effective, these models often need large amount

29 of time to make predictions (also referred to as model inference) on data instances. Often times, such  
30 limitation becomes a huge limiting factor for deploying ML models for latency critical applications. For  
31 example, applications such as high-frequency trading, military target recognition, pilots traveling at aircraft  
32 speeds (i.e., at least 621 mph) need to perform ML inference with the tight latency budget of at most 50  
33 ms. Making things worse, it is often not computationally and physically feasible to host large, effective  
34 ML models on directly on the devices like military aircraft. Thus, the ML model computations are often  
35 offloaded onto ground stations or edge devices, where data transmission can significantly add to the latency  
36 to execute model inference.

37 While many different types of machine learning accelerators have been developed, such as Graphical  
38 Processing Units (GPU) (Choquette et al., 2018), Field Programmable Gate Arrays (FPGA) (Fowers  
39 et al., 2018; He et al., 2018; Tong et al., 2017), and specialized application-specific integrated circuits  
40 (ASIC) (Chen et al., 2014; Jouppi et al., 2017), their efficiency is lowered by the data transfer time  
41 over the PCIe bus from the host system's network interface card (NIC). To overcome this challenge,  
42 we investigated methods to perform ML inference at the edge of the network to reduce the need and  
43 the overhead of transferring data from the edge to these accelerators. To achieve this, we note that the  
44 emergence of programmable data planes makes network devices (i.e., programmable switches and NICs)  
45 potential candidates for accelerating ML inference, especially given that programmable network devices  
46 have already been shown to be significantly power efficient while also providing high throughput and low  
47 latency in a variety of in-network computing tasks such as caching (Jin et al., 2017), consensus (Dang  
48 et al., 2020), and network monitoring (Kim et al., 2015). However, leveraging programmable data planes  
49 for machine learning inference still is an ongoing area of research with room for improvement.

50 Much of the prior works in this area has shown that network devices with programmable data planes,  
51 primarily programmable switches, demonstrate superior latency performance with minor degradation in  
52 model effectiveness (Zhang et al., 2023). While the line rate performance of programmable switches is  
53 beneficial for model inference, their limitation to match+action logic, memory size, cost and the placement  
54 in the network restrict the feasibility and accuracy of models that can be mapped onto them. For example,  
55 since many modern machine learning algorithms rely on operations such as multiplication during inference,  
56 finite-sized match+action tables cannot support every possible combination of multiplied values. Even  
57 though prior methods have found ways around this, it was not without loss in model effectiveness. And  
58 as machine learning models continue to grow, additional losses in model effectiveness seem likely. In  
59 addition, prior methods focused on general models that may not be used in the real-world for low-latency  
60 applications.

61 To compensate for this limitation, we propose the Smart Network Interface Cards (SmartNICs) as a  
62 viable alternative. SmartNICs possess additional computational resources and several packet processing  
63 accelerators that can be adapted to mimic essential machine learning inference operations, such as  
64 multiplication and logarithm functions, more accurately. Furthermore, SmartNICs are much more cost and  
65 power efficient, are more easier to deploy and test.

66 Therefore, in this paper, we present ML-NIC, a framework for compiling and deploying trained machine  
67 learning models onto SmartNICs by providing intelligent model mapping methods. This current work  
68 mainly focuses on mapping tree-based models onto SmartNICs due to it's wide usage of low-latency  
69 applications, but we have proposals for future work with proposed methods to support inference for other  
70 types machine learning models as well. Compared to many prior works that implement machine learning  
71 algorithms onto programmable network devices, ML-NIC implementation uses more device parallelism in

72 the inference process. Finally, our Python implementation of ML-NIC is made publicly available upon  
73 publication of the manuscript.

74 Our contributions include:

- 75 1. We present an algorithm to extract logic learned by a generic decision tree to facilitate parallelized  
76 feature analysis during inference.
- 77 2. We present a method to map and compile trained decision trees onto a SmartNIC in a manner that  
78 leverages its parallelism capabilities.
- 79 3. We demonstrate our framework's potential for accelerating the inference of decision trees for different  
80 tasks compared to conventional CPU and current state-of-the-art SmartNIC model deployment  
81 strategies.
- 82 4. We created an open-source project that contains all of ML-NIC's implementation and experimentation <sup>1</sup>.

83 The rest of this paper is organized as follows. Section 2 presents some background information on  
84 SmartNICs relevant to our work. Section 3 explains our approach towards deploying machine learning  
85 models onto a SmartNIC. Sections 4, 5, and 6 describe our experimental setup and discuss our results.  
86 Section 7 presents an overview of recent work that utilizes programmable data planes to accelerate the  
87 inference time for various machine learning algorithms on different problems. Section 8 points out future  
88 directions, and Section 9 ends with concluding remarks.

## 2 BACKGROUND

89 In this section, we provide the background for ML-NIC and some underlying motivations.

### 90 2.1 SmartNIC

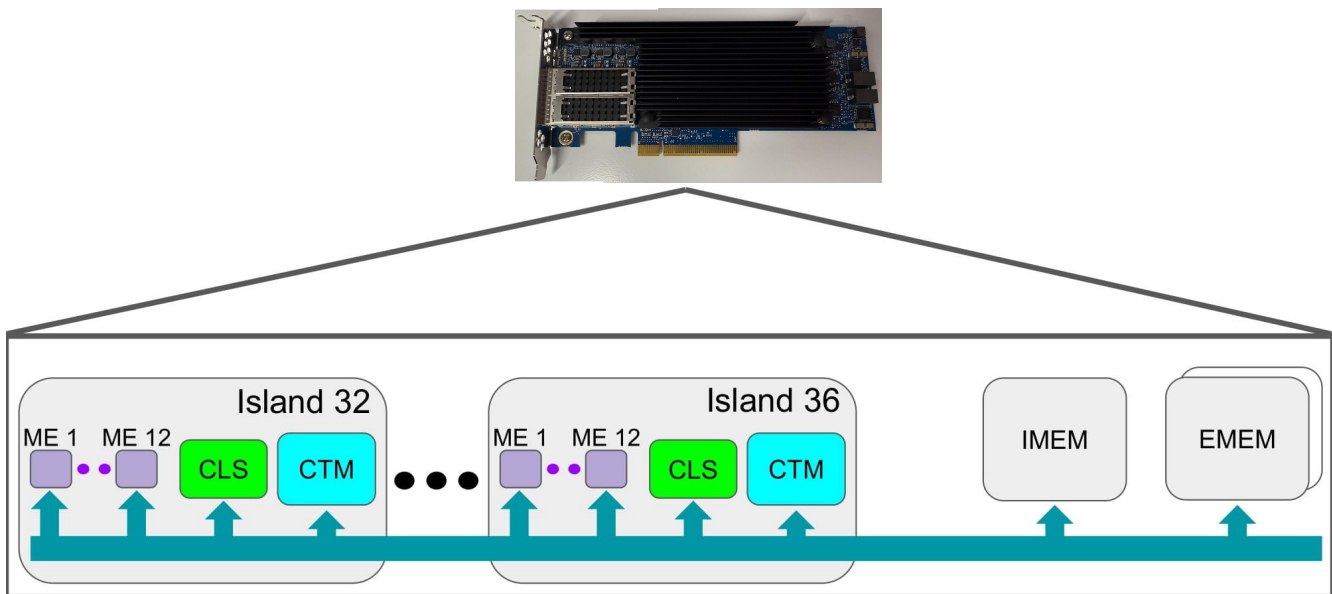
91 Smart Network Interface Cards (SmartNICs) possess additional computational resources and memory  
92 storage compared to traditional network interface cards. These resources enable SmartNICs to perform  
93 deep packet inspection, network function virtualization, and zero-trust security (Netronome Systems, 2024).  
94 As a result, offloading such operations to the SmartNIC frees a host system's CPU from conducting them.  
95 Compared to programmable switches, SmartNICs have more computational resources that can be leveraged.  
96 This motivates our choice to use SmartNICs for machine learning inference, since this process can be quite  
97 computationally intensive. For the rest of the section, we will focus on one particular type of SmartNIC:  
98 ASIC-Based Netronome SmartNICs supporting the NFP4000 architecture.

#### 99 2.1.1 ASIC-Based Netronome SmartNIC

100 ASIC-based SmartNIC represents a type of SmartNIC where the ASIC is custom built to support  
101 programmability using a set of custom languages. One notable example of an ASIC-based SmartNIC is  
102 a Netromome SmartNIC that can be programmed using languaged called P4 and Micro-C. This project  
103 utilizes a specific subset of the Netromome SmartNICs, which support the NFP4000 architecture. To  
104 elaborate further, Netronome SmartNICs support the NFP4000 architecture feature 48 packet processing  
105 cores and 60 programmable flow processing cores (Corigine, 2020). In much of the literature and technical  
106 documentation, the flow-processing cores are referred to as microengines (ME), which we denote as purple  
107 squares in Figure 1. Each microengine acts as an independent 32-bit processor with its own code store and  
108 local memory to run different programs in parallel with the other microengines. The microengines can be

---

<sup>1</sup> The project can be found on <https://github.com/The-Cloud-Lab/ML-NIC>



**Figure 1.** In this figure, we present a high-level overview of the NFP4000 architecture, focusing on the components most pertinent to ML-NIC.

109 programmed using a low-level language like Micro-C, an extended subset of C89, or a high-level language  
 110 like P4. A key difference between Micro-C and P4 is that Micro-C provides the flexibility to program each  
 111 microengine differently, whereas P4 defaults to loading the same program onto all microengines. However,  
 112 both languages lack floating-point number support. We provide a high-level illustration of the NFP4000  
 113 architecture in Figure 1.

114 Each microengine supports 8 threads, where each thread runs the same program and has its own block of  
 115 memory/registers. The following memory in a microengine is evenly partitioned among the 8 threads in a  
 116 microengine:

- 117 • 256, 32-bit General-Purpose Registers - used for general per-packet computations
- 118 • 256, 32-bit Transfer Registers - used for transferring data between memory regions
- 119 • 128, 32-bit Next-Neighbor Registers - used for communicating between neighboring microengines in  
 120 the same island
- 121 • 4kB of Local Memory - used for additional data storage as needed
- 122 • 120 Signal Registers - used to notify threads that a certain hardware event has occurred

123 The partitioning of memory among the 8 threads facilitates fast context switching between them, so they  
 124 can process different packets efficiently (Siracusano et al., 2022).

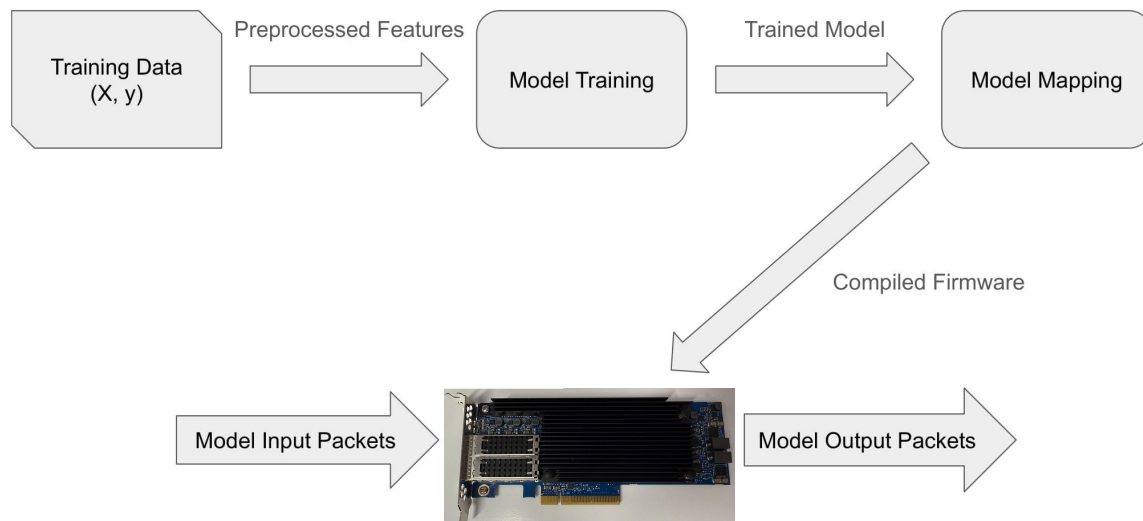
125 The microengines are organized into islands. While these islands can vary in number of microengines and  
 126 specialized functionality, standard islands, shown in Figure 1, contain 12 microengines with two regions of  
 127 memory shared between all the microengines in the island: Cluster Local Scratch (CLS) and Cluster Target  
 128 Memory (CTM). CLS, denoted in green in Figure 1, commonly stores small forwarding tables shared  
 129 between the microengines (Wray, 2014). CTM, denoted in cyan in Figure 1, holds packet headers and  
 130 coordinates between the microengines and other subsystems on the card (Wray, 2014). As CTM is larger  
 131 than CLS, more clock cycles are required to read/write to CTM.

132 Outside of the islands, the Netronome SmartNICs have three additional memory units, as shown in  
 133 Figure 1, shared with all microengines: one Internal Memory Unit (IMEM) and two External Memory  
 134 Units (EMEM) (Langlet, 2019). IMEM is used for storing packet payloads and medium-sized match-action  
 135 tables (Wray, 2014). EMEM is used to store larger match-action tables and other flow statistics (Wray,  
 136 2014). As these three memory units are the largest of those mentioned prior, with EMEM being larger than  
 137 IMEM, they require a greater number of clock cycles to read/write to them. Microengines can access data  
 138 in all these memory regions using the Command Push Pull (CPP) bus, denoted in turquoise in Figure 1.

139 Based on our knowledge of data access times for the different memory regions, hardware signals, and  
 140 transfer registers, we design ML-NIC to efficiently leverage these resources to ensure a high degree of  
 141 performance from a SmartNIC.

### 3 ML-NIC ARCHITECTURE

142 ML-NIC comprises three components: machine learning model training, model mapping, and model  
 143 deployment, as shown in Figure 2. We explain the details of each component below.



**Figure 2.** This figure shows the big picture of the ML-NIC architecture.

#### 144 3.1 Model Training

145 In the machine learning model training component, we consider a labeled dataset  $(X, y)$ , where  
 146  $X \in \mathbb{R}^{m \times n}$  represents our data matrix ( $m$  data points,  $n$  features) and  $y \in \{1 \dots q\}^m$  represents our  
 147 class labels ( $q$  possible classes). We make no assumptions on whether  $X$  consists of only continuous  
 148 features, only categorical features, or a mix. We assume the continuous features are normalized within the  
 149 range  $[0, 1]$ . We do so to simplify the range of numerical representation that the first iteration of ML-NIC  
 150 needs to account for. We find that this assumption is reasonable, since data normalization is a common  
 151 technique in machine learning to prevent certain features from dominating over other features due to  
 152 differences in scaling. For the categorical features, we assume that they are one-hot encoded (Liu, 2017)  
 153 (i.e., a feature with three categories is expanded to three features with values 0 or 1). This dataset is used to  
 154 learn the parameters of a particular machine learning model.

155 In this first iteration of ML-NIC, we choose to focus on decision trees for three reasons. First, we cite  
 156 the relative computational simplicity of tree traversal compared to floating-point operations in hardware  
 157 without a Floating-Point Unit (FPU) like a SmartNIC that was designed for fast computations at network  
 158 line rates.

159 Second, in comparison to other classical supervised machine learning models, such as Naive Bayes,  
 160 k-nearest neighbor, and support vector machine, we find that the decision tree is the more suitable choice  
 161 for offloading. With Naive Bayes, it is known that the algorithm perform poorly when the features used for  
 162 training are not conditionally independent. However, in practice, decision trees can perform well even if the  
 163 features are correlated. For k-nearest neighbor, we find that the required storage of every training instance  
 164 to be an obstacle for offloading onto network devices, especially given the size of modern datasets. Even if  
 165 only a selection of the training set was used to make offloading feasible, this could result in more significant  
 166 performance degradation in certain machine learning problems. With respect to support vector machines  
 167 for multiclass classification, the model may not be optimal for offloading with a large number of features  
 168 and number of classes. For this explanation, we temporarily denote  $m$  to be the number of features,  $q$  as  
 169 the number of classes (greater than 2 for multiclass classification), and  $z$  as the number of support vectors.  
 170 First, assume the support vectors for the support vector machine can be stored and there exists suitable  
 171 means of multiplication and a kernel on an off-the-shelf programmable network device that are at least as  
 172 expensive as a comparison operation. Since off-the-shelf programmable network device are optimized for  
 173 match + action, we expect the device architecture to have an efficient compare operation. Then, we see  
 174 that the support vector machine requires at least  $m \times z \times q$  multiplications, whereas a decision tree would  
 175 require at most  $m$  compare operations for a single inference. Therefore, based on instruction count, the  
 176 decision tree model has a better chance of yielding inference latency reduction when offloaded onto an  
 177 off-the-shelf programmable network device.

178 Third, we note that tree-like machine learning models, such as XGBoost (Chen and Guestrin, 2016),  
 179 are commonly used to learn tasks from structured tabular data over neural networks given faster training  
 180 time, potential performance gain, and model transparency. We use the decision tree model to show that  
 181 our framework can be used in real-world applications and offloading a machine learning model onto a  
 182 SmartNIC can reduce model inference latency significantly. While our focus is currently on decision tree  
 183 inference, we provide a discussion on how our framework can be augmented to account for additional  
 184 machine learning models in Section 6.

185 To train a decision tree, we consider the high-level algorithm outlined in Algorithm 1. Since finding the  
 186 globally-optimal tree structure for a learning task is computationally challenging, locally-optimal heuristic  
 187 algorithms are used such as ID3 (Quinlan, 1986), C4.5 (Quinlan, 1986), and CART (Leo Breiman and  
 188 Olshen, 1984). In practice, metric  $M$  is commonly Information Gain in the case of ID3 and C4.5 or Gini  
 189 Impurity for CART. Formulations for these metrics are provided in Equations 1 and 2. In Equations 1 and  
 190 2, we further define  $C$  as the number of classes,  $P$  as the number of splits,  $S_c$  as the number of examples  
 191 in the training dataset subset with class label  $c$ , and  $S_{p,c}$  as the number of examples in the  $p$ th split of the  
 192 training dataset subset with class label  $c$ .

$$Gini(S) = \sum_{c=1}^C \frac{|S_c|}{|S|} \left( 1 - \frac{|S_c|}{|S|} \right) \quad (1)$$

**Algorithm 1** Decision Tree Training Algorithm

**Require:**  $S$  is a valid subset of the training dataset  $D$ ,  $R$  is a set of stopping criteria for the algorithm,  $M$  is a valid impurity metric to locally optimize,  $split\_node$  is a valid function for a splitting  $S$  at a node

```

1: function TRAIN_TREE( $S, R, split\_node$ )
2:    $c = majority\_label(S)$ 
3:    $tree\_node = Node(label = c)$ 
4:   if  $not\_satisfied(R)$  then
5:      $ms = list()$ 
6:      $splits = split\_node(S)$ 
7:     for each  $split \in splits$  do
8:        $ms.append(M(split))$ 
9:     end for
10:     $best\_split = splits[arg\_optimal(ms, M)]$ 
11:    for each  $s \in best\_split$  do
12:       $tree\_node.insert\_branch(TRAIN\_TREE(s, R, split\_node))$ 
13:    end for
14:  end if
15:  return  $tree\_node$ 
16: end function

```

$$Info(S) = \sum_{c=1}^C \left( -\frac{|S_c|}{|S|} \log_2 \left( \frac{|S_c|}{|S|} \right) \right) - \sum_{p=1}^P \frac{|S_p|}{|S|} \sum_{c=1}^C \left( -\frac{|S_{p,c}|}{|S_p|} \log_2 \left( \frac{|S_{p,c}|}{|S_p|} \right) \right) \quad (2)$$

193 **3.2 Model Mapping**

194 Before discussing the technical details of decision tree mapping onto a SmartNIC, we discuss our mapping  
195 approach at a high-level. To run inference, we find the disjunctive normal form (Roth, 2016) of a decision  
196 tree. In the disjunctive normal form, the logic for assigning a class label to a data instance is expressed as a  
197 disjunction of conjunctions (i.e., (condition 1 and condition 2 and ...) or (condition3 and condition1 and ...) or ...).  
198 Each conjunction in the disjunction (i.e., condition 1 and condition 2 and ...) represents a path from  
199 the root node to a leaf node in a decision tree. We prefer the disjunctive norm form over the typical tree  
200 structure of a decision tree for inference, since it makes executing inference in a parallelized manner more  
201 convenient. To parallelize the inference process from the disjunctive normal form, we take the conditions  
202 from all the conjunctions that correspond to a particular feature, noting which path in the decision tree the  
203 condition corresponds to. To run inference on a data instance then, the conditions for each feature can be  
204 evaluated in parallel, where the result of each feature evaluation yields a set of paths in the decision tree  
205 that are possible for the data instance to take. Then, by aggregating the all possible paths and taking the  
206 intersection among them, a single path can be found. By matching the path to its corresponding class label,  
207 the decision tree inference process is complete.

208 To map a decision tree onto a SmartNIC, we take the output of the machine learning model training  
209 process (i.e., a pickle file) and proceed to generate an implementation of the SmartNIC data plane. Currently,  
210 we support SmartNICs that are programmable in Micro-C, primarily SoC-based Netronome SmartNICs. In  
211 the current iteration of our framework, we consider a trained decision tree classifier  $C$  with  $l$  leaf nodes,  
212 where  $l$  is at most 256. We make no additional assumptions on the number of splits per non-leaf node or  
213 the training algorithm used. Based on the number of leaf nodes  $l$  in model  $C$  and number of features  $n$  in  
214  $X$ , there are three possible scenarios for mapping  $C$  onto the SmartNIC:

- 215 • The model can fit on one island of the SmartNIC. Each island is then programmed with its own set of  
 216 feature computation, result aggregation, and packet collection microengines (i.e., inference for model  
 217  $C$  is run on all the islands)
- 218 • The model can fit on the entire SmartNIC with one feature assigned per feature microengine and one  
 219 packet collection microengine.
- 220 • The model can fit on the entire SmartNIC with multiple features assigned per feature microengine and  
 221 one packet collection microengine.

222 After selecting one of the three above mapping schemes, the next step is to extract the logic (i.e., find  
 223 disjunctive normal form and extract the conditions that match to a particular feature) learned by model  $C$ .  
 224 To do so, we iterate through all the  $n$  features in  $X$  and perform a depth-first search through the decision  
 225 tree. We record the operation for those nodes that run a comparison operation on our feature of interest and  
 226 continue the depth-first search until all the leaf nodes have been reached. Formally, Algorithm 2 illustrates  
 our logic extraction approach.

---

### Algorithm 2 Decision Tree Logic Extraction Algorithm

---

**Require:**  $node$  points to valid node in decision tree,  $ftre$  is feature seen by decision tree during training,  
 $clt$  has enough space to store decision tree logic for  $ftre$

```

1: function GET_LOGIC( $node, ftre, clt$ )
2:   if  $is\_leaf(node)$  then
3:      $clt.insert(node.prediction)$ 
4:   else
5:     if  $node.ftre = ftre$  then
6:        $clt.insert(node.logic)$ 
7:     end if
8:     for each  $child \in node.children$  do
9:        $clt.insert(GET\_LOGIC(child, ftre, clt))$ 
10:    end for
11:  end if
12:  return  $clt$ 
13: end function

```

---

227

228 We also assign the each of microengines on the SmartNIC as one of three types: packet collection, feature  
 229 computation, and result aggregation. The packet collection microengine(s) are programmed to signal the  
 230 CTM packet engine that they are ready to receive packets. Once a packet is received, the packet collection  
 231 microengine(s) will verify that packet is a model input packet, extract the features from the packet payload,  
 232 and asynchronously signal all the feature computation microengines of the inference request and transmit  
 233 the corresponding feature to each via transfer registers.

234 The feature computation microengines are responsible for evaluating the conditions on a feature for  
 235 a given data instance and determine which paths in the decision tree are possible. Since each feature  
 236 computation microengine is responsible for different features and run simulatenously, all the features can  
 237 be evaluated and all the possible paths in the decision tree can be determined in parallel. To implement the  
 238 conditions and determine the possible paths per feature on the SmartNIC, we use Micro-C if-statements  
 239 to evaluate the conditions and update an array of integers to reflect which paths are possible. For the  
 240 update, we treat the array of integers as a single bit string, where most significant bit in the integer at the  
 241 last index in the array corresponds to path 1. We assign paths based on the order in which the nodes are



242 encountered by Algorithm 2. Since the values for comparison in the conditions for evaluating each feature  
243 in the decision tree and the features themselves can be floating-point, and the SmartNIC does have an  
244 FPU, we consider an alternative floating-point representation. We represent floating-point numbers on the  
245 SmartNIC using a fixed-point representation that consists of 16 bits, where the last 13 bits represent the  
246 non-integer portion of a floating-point number. As each feature computation microengine completes its  
247 evaluation of its corresponding feature, they notify the result aggregation microengine(s) of the decision  
248 tree paths that are possible based on the feature they each evaluation.

249 Once all the feature computation microengines finish their evaluation, the result aggregation  
250 microengine(s) finds the intersecting path the decision tree between all the possible paths, matches the path  
251 to the corresponding class label, and sends an asynchronous signal to the packet collection microengine(s)  
252 along with the class label via transfer register(s). Once the packet collection microengine(s) receives the  
253 signal from the result aggregation microengine(s), it edits the original packet payload with the class label  
254 for the data instance and notified the CTM packet engine that the packet needs to be transmitted. Also  
255 note that during the time the feature computation and result aggregation microengines are completing their  
256 tasks, the packet collection microengine(s) are editing the model input packet's header in preparation for  
257 transmission as a model output packet.

258 To program all the packet collection, feature computation, and result aggregation microengines, separate  
259 Micro-C code is written to program each microengine to complete their specific task for model inference,  
260 whereas prior methods often program all the microengines with one piece of P4 code to perform the same  
261 tasks for the model inference and do not fully leverage the parallel operating capacity of the SmartNIC.  
262 Example Micro-C code for packet collection, feature computation, and result aggregation can be found  
263 below in Listings 1, 2, and 3.

### 264 3.3 Model Deployment

265 Once all the Micro-C code files are created, they are all compiled and linked to generate the device  
266 firmware to run on the data plane in the model deployment component. Then, the firmware file output  
267 is loaded onto the SmartNIC. An example of the full process is shown in Listing 4. Each microengine  
268 assumes a specified behavior based on one of the three microengine assignments specified above. The  
269 SmartNIC can now ingress packets with features in the packet payload, run machine learning inference in a  
270 parallelized manner, and egress packets with the classification result as the packet payload.

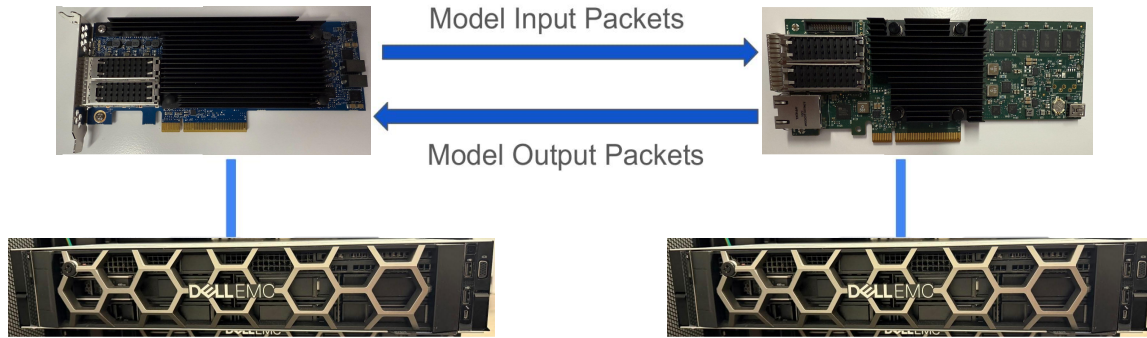
## 4 EXPERIMENTAL SETUP

### 271 4.1 Testbed

272 Our testbed consists of two Dell PowerEdge Rack Servers. Server 1 hosts an NVIDIA Mellanox Bluefield-  
273 2 DPU 25 GbE SmartNIC for packet transmission and data collection. Server 2 hosts a Netronome  
274 AgilioCX 2 × 25 GbE SmartNIC, on which our decision tree models are deployed. Both systems are  
275 directly connected via qsfp cable between the Mellanox and Netronome SmartNICs. We illustrate our setup  
276 in Figure 3.

### 277 4.2 Datasets and Models

278 Our evaluation considers four tasks: land mine detection, satellite image pixel classification, gas sensor  
279 drift compensation, and network traffic classification. The main characteristics of the datasets used for  
280 each task can be found in Table 1. We train a decision tree model for each task using the scikit-learn



**Figure 3.** The testbed setup used for evaluation. The left server hosts the Netronome AgilioCX  $2 \times 25$  GbE. The right server hosts the NVIDIA Mellanox Bluefield-2 DPU 25 GbE.

281 library (Pedregosa et al., 2011). We summarize the hyperparameters used for each tree in Table 2. For  
 282 hyperparameters not explicitly mentioned in the table that can be tuned for the decision tree models (i.e.,  
 283 criterion, splitter, max features, etc.), we resort to the default values provided by scikit-learn.

**Table 1.** Summary of Datasets Used (refer to dataset subsections for class label abbreviations)

Attribute	Mine	Landsat	Gas	CICIDS
# of features	3	36	128	7
# of data instances	338	6435	13910	22887218
# of classes	5	6	6	7
# of training data	270	4435	11128	500000
# of test data	68	2000	2782	6957375

**Table 2.** Summary of decision tree Models Created

Parameter	Mine	Landsat	Gas	CICIDS
# of leaves	114	256	256	89
depth	17	15	24	15
# of nodes	227	511	511	177
min samples leaf	1	1	1	2
min samples split	2	2	2	2
min impurity decrease	0	0	0	0.00001
max leaf nodes	None	256	256	None

#### 284 4.2.1 Dataset Preprocessing

285 As mentioned in Section 3, we assume the continuous features are in the range  $[0, 1]$  and categorical  
 286 features are one-hot-encoded. To achieve this, we apply min-max normalization to scale the continuous  
 287 features of each dataset to range between 0 and 1 using the training set. Test features that lie outside the  
 288 range  $[0, 1]$  after min-max normalization has been applied are clipped to the closest endpoint. We also  
 289 one-hot-encode the categorical features for each dataset based on the values observed from the training set.  
 290 If the categorical features in the test set take on values not observed in the training set, the one-hot-encoded  
 291 feature is represented as a bit string of zeros.

#### 292 4.2.2 Land Mine Detection

293 We use the Land Mines dataset (Yilmaz et al., 2018) for the land mine detection task. The authors  
294 propose three features to classify a mine into five types, Null, Anti-Tank, Anti-Personnel, Booby-Trapped  
295 Anti-Personnel, and M14 Anti-Personnel, with 65 – 71 samples per class. Our motivation for choosing  
296 this dataset is based on the number of features (8 after data preprocessing), where we can evaluate the first  
297 SmartNIC mapping scenario (fitting on one island) as described in Section 3. In later sections, we will  
298 refer to this dataset as Mine.

#### 299 4.2.3 Satellite Image Pixel Classification

300 We use the Statlog (Landsat Satellite) dataset (Srinivasan, 1993) for the satellite image pixel classification  
301 task. The goal of this task is to examine multispectral values from a  $3 \times 3$  neighborhood of a satellite image  
302 and classify the central pixel as one of five classes: Red Soil, Cotton Crop, Grey Soil, Damp Grey Soil,  
303 Soil with Vegetation Stubble, Mixture, or Very Damp Grey Soil. There are 626 – 1533 samples per class.  
304 Our motivation for choosing this dataset is based on the number of features (36), where we can evaluate  
305 the second mapping scenario (fitting on whole SmartNIC, one feature per microengine) as described in  
306 Section 3. In later sections, we will refer to this dataset as Landsat.

#### 307 4.2.4 Gas Sensor Drift Compensation

308 We use the Gas Sensor Array Drift dataset (Rodríguez-Luján et al., 2014) for the gas sensor drift  
309 compensation tasks. This dataset consists of measurements from 16 chemical sensors to identify six gases,  
310 Ammonia, Acetaldehyde, Acetone, Ethylene, Ethanol, and Toluene, with 1508 – 3009 samples per class.  
311 Our motivation for choosing this dataset is based on the number of features (128), where we can evaluate  
312 the third mapping scenario (fitting on whole SmartNIC, multiple features per microengine) as described in  
313 Section 3. In later sections, we will refer to this dataset as Gas.

#### 314 4.2.5 Network Traffic Classification

315 We use the CICIDS2017 dataset (Sharafaldin et al., 2018) for the network traffic classification use case.  
316 This use case's purpose is to identify network flows as benign or malicious (brute force attack, heartbleed  
317 attack, botnet, DoS attack, DDoS attack, web attack, infiltration attack). However, we follow the approach  
318 used by Xavier et al. (2021) to generate the dataset for classifying individual network packets rather  
319 than network flows and the training and tests sets based on the network flows instead of the conventional  
320 stratified 80/20 split used for the above datasets above. In the dataset, the packets were labeled as benign,  
321 DoS GoldenEye, DoS Hulk, DoS Slowhttptest, DoS Slowloris, Web Brute Force, or Port Scan. Each class  
322 has 30059 – 20121944 samples. Our motivation for choosing this dataset is based on its use in the work by  
323 Xavier et al. (2021), which is similar to our approach. Our evaluation on this dataset clearly compares our  
324 approach and Xavier et al. (2021)'s approach. In later sections, we will refer to this dataset as CICIDS.

### 325 4.3 Baselines

326 We compare our approach against the following two baselines. First, we implement a traditional CPU  
327 baseline, which uses socket programming to receive incoming packets, extract the payload, run inference  
328 with the trained scikit-learn decision tree, and build and send a model output packet with the model  
329 prediction in the packet payload.

330 Second, we implement the approach developed by Xavier et al. (2021) using P4-16. Like our approach,  
331 Xavier et al. (2021)'s approach traverses through the scikit-learn decision tree structure, extracts the model's

332 logic, and rebuilds the tree in P4 using Python. Note that the original implementation was solely created  
333 for the CICIDS dataset, and the authors did not provide a method to handle floating-point features. In  
334 evaluating this method on the other datasets, we modified it slightly to use our fixed-point representation of  
335 floating-point numbers. Also, due to limitations with Xavier et al. (2021)'s approach, we could not evaluate  
336 it on larger decision trees, such as those generated with the Landsat and Gas Datasets.

#### 337 4.4 Evaluation Metrics

338 In our experiments, we measure the effectiveness, (average and tail) latency, throughput, and hardware  
339 utilization of ML-NIC against the baselines.

340 For effectiveness, we measured the accuracy, F1 score, recall, and precision metrics on each dataset's  
341 test set. Given that our datasets are for multiclass classification tasks, we take the macro-average (i.e.,  
342 unweighted mean) of the per-class scores for the F1 score, recall, and precision measurements.

343 For latency, we collected the time between model input packet transmission and model output packet  
344 reception on server 1 in microseconds for 1000 packets. In addition to our vanilla latency experiments  
345 (i.e., no CPU load or network link utilization), we conduct latency experiments with background traffic on  
346 the network link and CPU load. We generate random network traffic at different speeds using Tcpreplay  
347 for latency experiments with background traffic to achieve 25%, 50%, and 99% network link utilization.  
348 We use stress-ng for latency experiments with CPU load and generate CPU loads of 25%, 50%, and 99%.  
349 To ensure an apples-to-apples comparison with the CPU baseline, we append zero padding to the model  
350 input packets for our approach and the P4 baseline. Hence, they are the same size as the CPU model input  
351 packets. Note that, when collecting the data for the CPU baseline, we remove the time taken to decode the  
352 data features and encode the model's prediction.

353 For throughput, we use Tcpreplay to loop through the PCAP files containing the test set packets for each  
354 dataset at top speed. Simultaneously, we also run Tshark to filter and collect the model prediction packets  
355 for 60 seconds.

356 Lastly, for hardware utilization, we run Tcpreplay for 60 seconds like we did for the throughput experiment  
357 and measure CPU, server host RAM, and SmartNIC memory utilization. We also measure CPU, server host  
358 RAM, and SmartNIC memory utilization 30 seconds before and after running Tcpreplay for reference. We  
359 do not explicitly measure SmartNIC microengine utilization in this experiment. Instead, we use the results  
360 from our network link utilization and CPU load latency experiments as a proxy for SmartNIC microengine  
361 utilization.

## 5 RESULTS

362 From our experimental results, we demonstrate the following:

- 363 1. Our approach achieves effectiveness scores similar to those of the CPU baseline and identical to the P4  
364 baseline.
- 365 2. Our approach has better a latency guarantee than the CPU and P4 baselines in various network link  
366 utilization and CPU load scenarios.
- 367 3. The throughput of our approach is significantly greater compared to the CPU baseline and on par with  
368 the P4 baseline.
- 369 4. Our approach uses fewer server host resources (i.e., CPU and server RAM) compared to the CPU and  
370 P4 baselines.

## 371 5.1 Effectiveness Scores

372 As shown in Table 3<sup>2</sup>, the effectiveness scores between our approach and the CPU baselines are similar  
 373 with minor degradation. For conciseness, we only show the plots of accuracy and F1 score of the models,  
 374 since the precision and recall results follow a similar pattern.

**Table 3.** Effectiveness Measurements On All Datasets

Dataset	Measure	CPU	Xavier et al. (2021)	ML-NIC
Mine	Accuracy (%)	<b>58.82</b>	57.35	57.35
	F1 Score (%)	<b>58.22</b>	56.92	56.92
Landsat	Accuracy	<b>85.65</b>	N/A	85.55
	F1 Score (%)	<b>83.93</b>	N/A	83.76
Gas	Accuracy (%)	<b>97.41</b>	N/A	95.15
	F1 Score (%)	<b>97.26</b>	N/A	94.78
CICIDS	Accuracy (%)	<b>95.12</b>	<b>95.12</b>	<b>95.12</b>
	F1 Score (%)	<b>47.04</b>	<b>47.04</b>	<b>47.04</b>

375 For the Mine dataset, we note differences of 1.47%, 1.30%, 1.35%, and 1.54% across the accuracy,  
 376 F1 score, precision, and recall metrics. For the Landsat dataset, we note differences of 0.100%, 0.17%,  
 377 0.14%, and 0.20% across the accuracy, F1 score, precision, and recall metrics. For the Gas dataset, we  
 378 note differences of 2.26%, 2.49%, 2.00%, and 2.66% across the accuracy, F1 score, precision, and recall  
 379 metrics. For the CICIDS dataset, our approach and the CPU baseline do not differ in accuracy, F1 score,  
 380 precision, or recall. This is because all the features used to train the scikit-learn decision tree are integers,  
 381 so no quantized representation of features is needed as with the previous three datasets. Our approach  
 382 achieves identical effectiveness scores as the P4 baseline on the Mine and CICIDS datasets.

## 383 5.2 Latency

384 From the latency data we collected, we provide zoomed-in empirical cumulative distribution functions  
 385 (eCDF) for each dataset, network link utilization, and CPU load in Figure 4. We also provide more concrete  
 386 numbers on the 50th, 99th, and 99.9th percentiles across each dataset, link utilization, and CPU load in  
 387 Tables 4 and 5. From our experiments, we make the following observations. We generally see a significant  
 388 gap in the latency measurements between ML-NIC and the CPU baseline and a very small gap between  
 389 ML-NIC and Xavier et al. (2021)'s approach. Specifically, we found that ML-NIC's latency can be at least  
 390 132.62  $\mu$ s faster than the CPU baseline and 1.35  $\mu$ s faster than Xavier et al. (2021)'s approach in the 50th  
 391 percentile. However, there is a significant difference in the tails between ML-NIC and Xavier et al. (2021)'s  
 392 approach, suggesting that ML-NIC has a stronger latency guarantee. Based on the 99.9th percentiles, we  
 393 see that the tail latency of Xavier et al. (2021)'s approach can be at least 1.53 $\times$  larger than ML-NIC's tail.

394 Looking at impact of high network link utilization and CPU load, we observe very minimal fluctuation in  
 395 the eCDFs of the ML-NIC and Xavier et al. (2021)'s approach. This suggests that both approaches are  
 396 robust against high network link utilization and CPU load on these datasets. But, there is a more noticeable  
 397 impact of high network link utilization and CPU load on the CPU baseline. As the network link utilization  
 398 increases, we tend to see more probability mass shift towards the higher latency in the eCDF. With respect  
 399 to the 99.9th latency percentile, we see an increase of at least 1.82 $\times$  from 0% link utilization to 99% link  
 400 utilization. Concerning the increases in CPU load, there is a more significant shift in the eCDF curves

<sup>2</sup> Bold values indicate best values found for a given metric during experiments

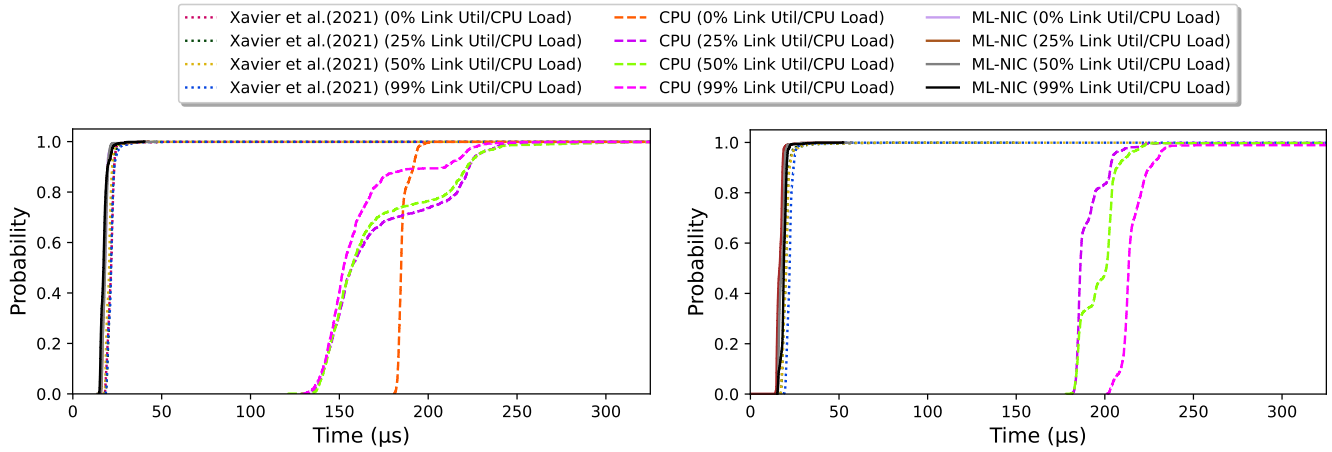


Figure 4a. Mine Tcpreplay

Figure 4b. Mine Stress-Ng

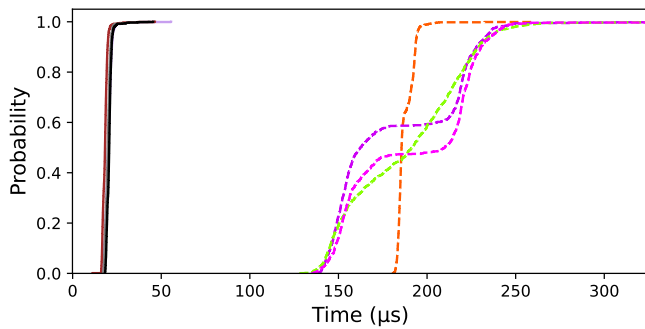


Figure 4c. Landsat Tcpreplay

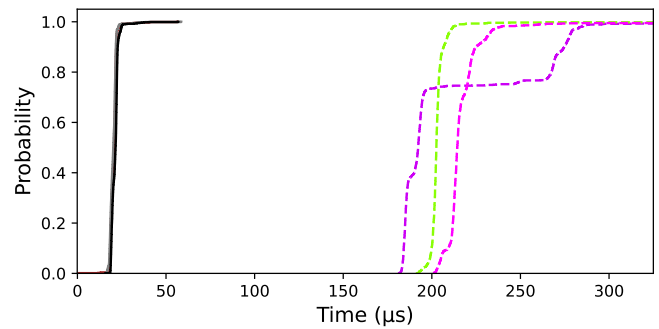


Figure 4d. Landsat Stress-Ng

401 towards higher latencies. Referring to the 99.9th percentile, there is a latency increase of at least  $20.27 \times$   
 402 from 0% CPU load to 99% CPU load.

Table 4. Latency Measurements Using Tcpreplay ( $\mu s$ )

Dataset	Link Util (%)	Percentile	CPU	Xavier et al. (2021)	ML-NIC
Mine	0	50	184.77	21.55	<b>18.00</b>
		99	197.39	25.41	<b>23.60</b>
		99.9	202.94	71.26	<b>35.54</b>
	25	50	156.93	21.30	<b>17.73</b>
		99	243.95	27.35	<b>23.63</b>
		99.9	285.89	65.31	<b>36.11</b>
	50	50	156.96	20.41	<b>18.15</b>
		99	267.83	27.33	<b>21.53</b>
		99.9	364.59	60.37	<b>31.05</b>
99	50	152.32	21.55	<b>17.24</b>	
	99	230.96	31.10	<b>24.45</b>	
	99.9	454.20	70.65	<b>36.48</b>	
Landsat	0	50	185.61	N/A	<b>19.99</b>
		99	200.27	N/A	<b>24.88</b>

		99.9	260.64	N/A	<b>39.63</b>
	25	50	162.39	N/A	<b>18.27</b>
		99	247.78	N/A	<b>22.70</b>
		99.9	365.45	N/A	<b>36.65</b>
	50	50	191.58	N/A	<b>19.82</b>
		99	258.51	N/A	<b>23.87</b>
		99.9	458.64	N/A	<b>42.08</b>
	99	50	209.91	N/A	<b>20.63</b>
		99	247.88	N/A	<b>24.96</b>
		99.9	503.92	N/A	<b>37.40</b>
Gas	0	50	185.99	N/A	<b>18.91</b>
		99	198.86	N/A	<b>26.57</b>
		99.9	287.77	N/A	<b>48.31</b>
	25	50	153.22	N/A	<b>20.60</b>
		99	255.51	N/A	<b>27.00</b>
		99.9	367.69	N/A	<b>41.65</b>
	50	50	159.18	N/A	<b>20.13</b>
		99	251.86	N/A	<b>28.91</b>
		99.9	410.64	N/A	<b>45.20</b>
	99	50	161.33	N/A	<b>22.11</b>
		99	244.74	N/A	<b>25.35</b>
		99.9	523.63	N/A	<b>43.98</b>
CICIDS	0	50	179.30	19.58	<b>17.21</b>
		99	189.71	23.86	<b>21.66</b>
		99.9	198.83	54.91	<b>31.70</b>
	25	50	154.72	20.00	<b>15.97</b>
		99	224.22	26.66	<b>19.76</b>
		99.9	277.53	55.06	<b>29.29</b>
	50	50	155.85	18.94	<b>17.55</b>
		99	225.52	23.11	<b>20.68</b>
		99.9	442.80	55.48	<b>33.41</b>
	99	50	154.81	21.91	<b>15.99</b>
		99	226.99	28.21	<b>20.39</b>
		99.9	483.02	51.85	<b>33.80</b>

**Table 5.** Latency Measurements Using Stress-Ng ( $\mu s$ )

Dataset	CPU Load (%)	Percentile	CPU	Xavier et al. (2021)	ML-NIC
Mine	25	50	186.16	20.16	<b>16.45</b>
		99	221.97	31.82	<b>20.43</b>
		99.9	343.53	74.56	<b>32.35</b>
	50	50	200.85	20.12	<b>17.69</b>
		99	222.64	31.16	<b>20.96</b>
		99.9	535.06	67.87	<b>34.50</b>
	99	50	213.40	22.12	<b>18.84</b>

		99	332.46	27.96	<b>22.66</b>
		99.9	7979.24	76.50	<b>36.86</b>
Landsat	25	50	192.30	N/A	<b>20.69</b>
		99	295.65	N/A	<b>24.48</b>
		99.9	405.26	N/A	<b>41.19</b>
	50	50	202.55	N/A	<b>20.23</b>
		99	215.20	N/A	<b>24.12</b>
		99.9	7582.86	N/A	<b>38.09</b>
	99	50	214.34	N/A	<b>21.54</b>
		99	260.63	N/A	<b>25.80</b>
99.9		9008.55	N/A	<b>40.27</b>	
Gas	25	50	183.90	N/A	<b>19.17</b>
		99	201.89	N/A	<b>27.42</b>
		99.9	469.44	N/A	<b>45.69</b>
	50	50	203.31	N/A	<b>23.01</b>
		99	238.13	N/A	<b>26.94</b>
		99.9	8852.89	N/A	<b>46.95</b>
	99	50	218.32	N/A	<b>19.40</b>
		99	251.67	N/A	<b>25.54</b>
99.9		16182.94	N/A	<b>46.74</b>	
CICIDS	25	50	155.58	20.13	<b>18.17</b>
		99	201.23	23.99	<b>21.31</b>
		99.9	357.74	51.48	<b>31.69</b>
	50	50	167.66	19.88	<b>18.18</b>
		99	184.74	24.39	<b>22.68</b>
		99.9	481.76	50.01	<b>30.33</b>
	99	50	175.56	20.15	<b>17.16</b>
		99	198.38	24.00	<b>21.83</b>
99.9		4030.55	59.39	<b>32.74</b>	

### 403 5.3 Throughput

404 As seen in Figure 5, there is a significant improvement in throughput with our approach compared to the  
 405 CPU baseline. In our approach, we note  $24.80\times$ ,  $19.30\times$ ,  $16.95\times$ , and  $20.11\times$  more packets per minute  
 406 compared to the CPU baseline across the Mine, Landsat, Gas, and CICIDS datasets. Furthermore, our  
 407 approach yields moderately higher throughput than the P4 baseline in the Mine and CICIDS dataset. In our  
 408 approach, we observe  $1.26\times$  and  $1.11\times$  more packets per minute compared to the P4 baseline for the Mine  
 409 and CICIDS datasets.

### 410 5.4 Hardware Utilization

411 From our hardware utilization experiment, we report the minimum, maximum, and average CPU and  
 412 server host RAM utilization in Table 6. We also report the SmartNIC memory utilization as a constant, since  
 413 dynamic memory allocation is not available on the AgilioCX  $2 \times 25$  GbE SmartNIC. Since we are not able  
 414 to directly measure the SmartNIC RAM used for the CPU baselines, we approximate it. Our approximation  
 415 takes an unweighted average of the ratio of size of SmartNIC firmware for the CPU baselines over the size



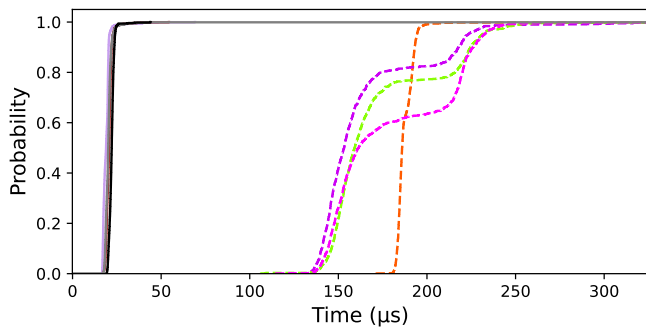


Figure 4e. Gas Tcpreplay

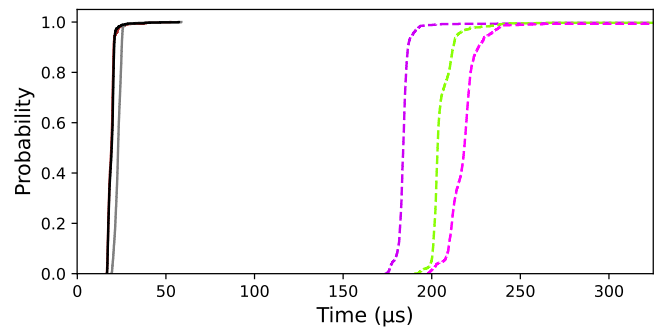


Figure 4f. Gas Stress-Ng

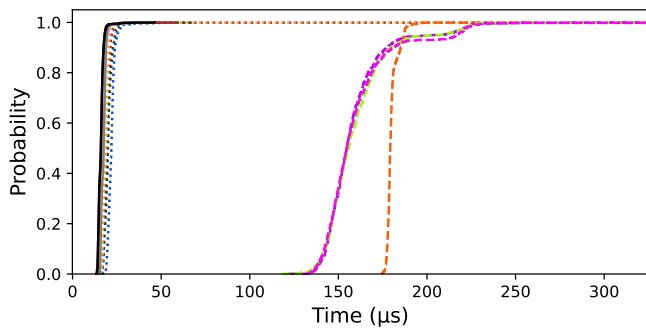


Figure 4g. CICIDS Tcpreplay

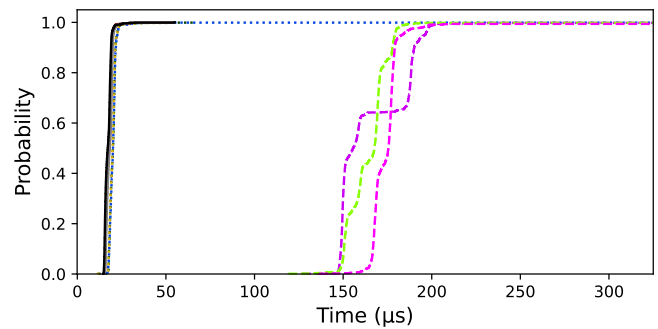


Figure 4h. CICIDS Stress-Ng

Figure 4. Figures 4a, 4c, 4e, and 4g depict the eCDFs for the latency experiments conducted using Tcpreplay to saturate the network link on all the datasets. Figures 4a, 4c, 4e, and 4g depict the eCDFs for the latency experiments conducted using Stress-Ng to generate a CPU load on all the datasets.

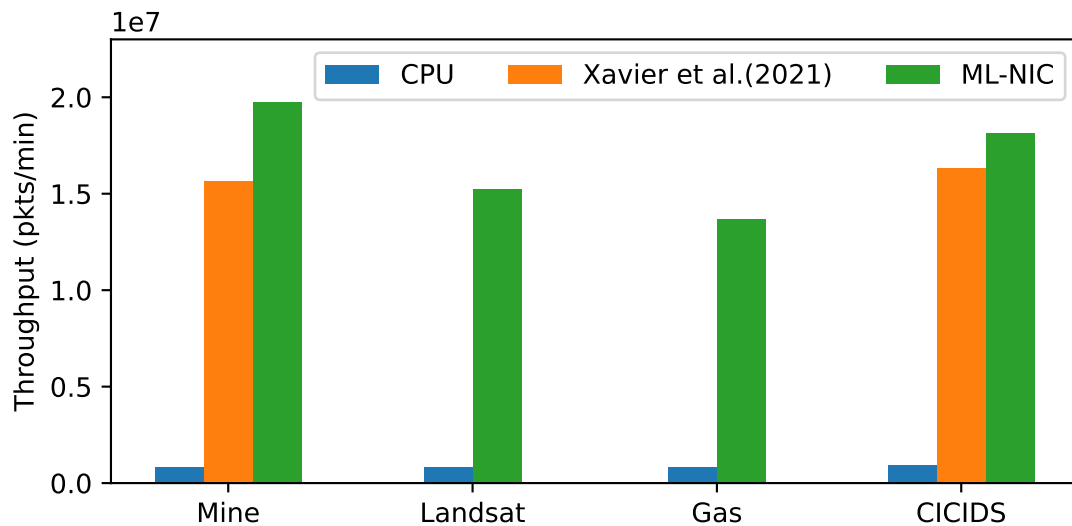


Figure 5. This figure provides a bar graph of the throughputs observed with our approach and two baselines.

416 of the SmartNIC firmware for the P4 baselines and our approach multiplied by the SmartNIC RAM used  
 417 by those models for each of the datasets.

**Table 6.** Hardware Utilization Measurements

Dataset	Measure	CPU	Xavier et al. (2021)	ML-NIC
Mine	Min CPU (%)	<b>0.25</b>	<b>0.25</b>	<b>0.25</b>
	Avg CPU (%)	3.82	0.77	<b>0.52</b>
	Max CPU (%)	7.56	7.75	<b>0.94</b>
	Min Host RAM (MB)	2226.76	2147.76	<b>1908.97</b>
	Avg Host RAM (MB)	2228.43	2150.75	<b>1910.71</b>
	Max Host RAM (MB)	2230.56	2153.44	<b>1911.81</b>
	SmartNIC RAM (MB)	96.15	973.00	<b>21.13</b>
Landsat	Min CPU (%)	<b>0.25</b>	N/A	<b>0.25</b>
	Avg CPU (%)	3.81	N/A	<b>0.51</b>
	Max CPU (%)	7.53	N/A	<b>0.88</b>
	Min Host RAM (MB)	2226.50	N/A	<b>1909.86</b>
	Avg Host RAM (MB)	2229.00	N/A	<b>1910.69</b>
	Max Host RAM (MB)	2231.77	N/A	<b>1911.70</b>
	SmartNIC RAM (MB)	96.15	N/A	<b>21.12</b>
Gas	Min CPU (%)	0.25	N/A	<b>0.19</b>
	Avg CPU (%)	3.78	N/A	<b>0.51</b>
	Max CPU (%)	7.46	N/A	<b>0.88</b>
	Min Host RAM (MB)	2223.59	N/A	<b>1908.86</b>
	Avg Host RAM (MB)	2226.08	N/A	<b>1909.83</b>
	Max Host RAM (MB)	2228.07	N/A	<b>1910.98</b>
	SmartNIC RAM (MB)	96.15	N/A	<b>21.122</b>
CICIDS	Min CPU (%)	<b>0.19</b>	0.25	0.25
	Avg CPU (%)	3.68	0.80	<b>0.52</b>
	Max CPU (%)	7.20	9.09	<b>0.81</b>
	Min Host RAM (MB)	2225.02	2077.43	<b>1907.22</b>
	Avg Host RAM (MB)	2227.38	2078.88	<b>1908.36</b>
	Max Host RAM (MB)	2230.01	2080.04	<b>1909.21</b>
	SmartNIC RAM (MB)	96.15	973.00	<b>21.12</b>

418 From Table 6, we see that ML-NIC consistently uses lower resources for average and maximum host  
419 system's CPU, host system's RAM, and SmartNIC's RAM usage compared to the CPU baseline and  
420 Xavier et al. (2021)'s method. In the case where there the CPU baseline has slightly lower minimum CPU  
421 usage than the ML-NIC on the CICIDS dataset, this measurement likely corresponds to some degree of  
422 randomness in the measurement, since ML-NIC also achieved the same minimum CPU usage on the Gas  
423 dataset. In addition to the CPU baseline having a higher maximum CPU usage than ML-NIC by at least  
424  $7.91\times$ , we also observe that Xavier et al. (2021)'s method can achieve similar or higher levels of maximum  
425 CPU usage. We attribute this to the runtime environment (RTE) server that is running on our host system,  
426 which is needed to run P4 code. We believe this also accounts for slightly higher host system RAM usage.  
427 For the SmartNIC's RAM usage, we observe our proxy for the CPU baseline to be lower than Xavier et al.  
428 (2021)'s method. Since the firmware running on the SmartNIC for the CPU baseline runs as a regular NIC,  
429 the SmartNIC would not require additional memory beyond storing extracting packet headers into local  
430 memory or general purpose registers. Furthermore, the larger SmartNIC RAM usage from Xavier et al.

431 (2021)'s method likely occurs because their approach involves running packet collect, feature computation,  
432 and result aggregation on every microengine. Since ML-NIC distributes these operations across multiple  
433 microengines, the resulting SmartNIC RAM usage would be lower by at least  $46.05\times$ .

## 6 DISCUSSION

### 434 6.1 Generalization

435 The work focuses on converting trained scikit-learn decision trees into Micro-C for deployment onto a  
436 SmartNIC. We focused on the Netronome AgilioCX  $2 \times 25$  GbE. Deployment across different SmartNICs  
437 (assuming Micro-C support) may require significant code changes to accommodate the resources available  
438 on the card compared to the baselines.

### 439 6.2 Benefits

440 Despite their resource constraints compared to a host system's CPU, SmartNICs show potential as  
441 alternative hardware for deploying appropriately sized decision tree models. Deploying the decision tree  
442 model onto the SmartNIC, which brings it closer to the network edge, saves latency time by removing the  
443 need to transfer data over the PCIe bus from the NIC to the CPU without additional hardware. Furthermore,  
444 the lower-level programming used in our approach compared to the P4 baseline allows us to leverage  
445 device parallelism to deploy larger decision trees.

### 446 6.3 Scope

447 Our work only considers deploying decision tree models trained for various tasks onto a SmartNIC.  
448 Improvements in any specific use case are beyond the scope of our work.

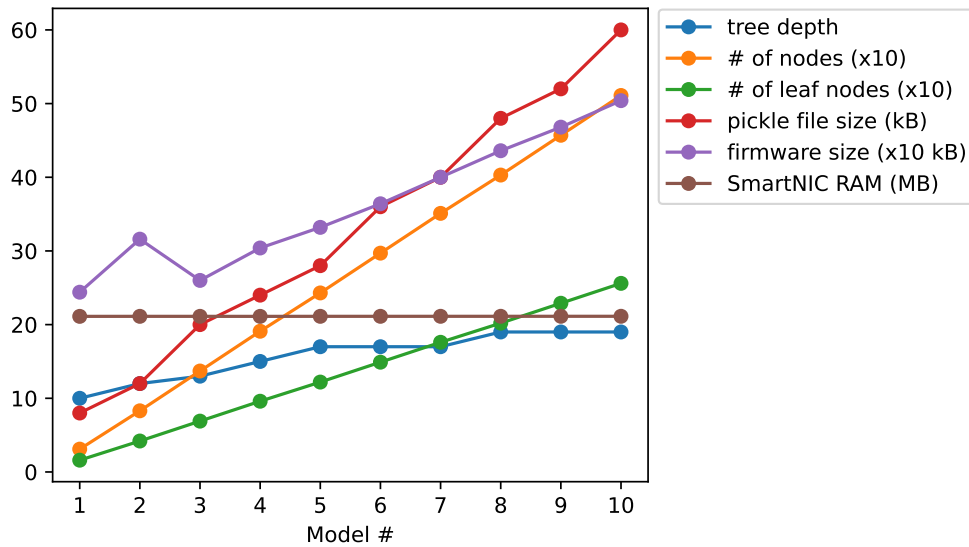
### 449 6.4 Limitations

450 Our method's limitations depend on the SmartNIC's memory, computational, asynchronous I/O, and  
451 data rate constraints. Within the Netronome AgilioCX  $2 \times 25$  GbE, the primary limits are the number  
452 of microengines (60), data rate (25 GbE), and number of hardware signals (15 per thread). While the  
453 constraint on microengines can be mitigated by assigning multiple features to a microengine, the memory  
454 (i.e., number of transfer registers) and asynchronous I/O (i.e., number of hardware signals) constraints limit  
455 the depth of the trained decision tree to 480 leaf nodes. Our current implementation is limited to decision  
456 trees with 256 leaf nodes, since more complex firmware is required to assign a hardware signal to a greater  
457 number of transfer registers.

### 458 6.5 Offloading More Models

#### 459 6.5.1 Decision Tree

460 In addition to the decision tree models we have deployed in this work, we provide some more insights  
461 on other decision trees that our current work can offload onto a SmartNIC via an ablation study. In our  
462 ablation study, we use the CICIDS datasets to construct 10 decision trees with a constraint on the maximum  
463 number of leaf nodes between 16 and 256, based on the leaf node limit we mentioned in Section 3. For  
464 each decision tree, we look at the depth, number of nodes, number of leaf nodes, size of the pickle file,  
465 size of the firmware file, and SmartNIC RAM usage. We present our findings in Figure 6 below. Note  
466 that we scaled some of the measurements by a factor of 10, so the trends in the some of the decision tree  
467 parameters would be more clear.



**Figure 6.** This figure shows how ML-NIC scales with respect to different decision tree parameters.

468 From Figure 6, we observe the following. First, we note the slow inclination, followed by a brief  
 469 declination, then continued inclination in the usage of SmartNIC RAM. We also observe a similar pattern of  
 470 inclination, declination, then inclination again in the trend for the model firmware size. We attribute this to  
 471 how we compiled two instances of the model 1 and 2 per island to maximize our usage of the computation  
 472 resources on the SmartNIC. For the remaining models, we only compiled one instance per island. Since  
 473 models 1 and 2 have two instances compiled per island, more RAM and instruction memory would be  
 474 needed to store the labels for the leaf nodes and decision tree logic. Based on the rate of inclination between  
 475 firmware size and SmartNIC RAM usage, we see that a primary concern for offloading larger models is the  
 476 amount of instruction memory available per microengine. After model 2, we see that the (scaled) trend for  
 477 firmware size grows slower than that of the size of model pickle file and number of nodes by 1.64 and 1.53  
 478 and grows faster than the trend for decision tree depth by 4.07 and number of leaf nodes by 1.30.

## 479 6.5.2 Other Machine Learning Models

480 Besides decision trees, we also consider approaches for executing inference with other machine learning  
 481 models. Since inference for many popular machine learning models relies heavily on the matrix-vector  
 482 multiplication operation, we look into techniques for efficient and effective matrix-vector multiplication  
 483 that can be performed by a SmartNIC. In addition to conducting inference on models such as neural  
 484 network and support vector machines for supervised tasks, we find an implementing a suitable matrix-  
 485 vector multiplication method necessary for unsupervised learning, such as with implementing k-means  
 486 using cosine similarity as the similarity measure instead of Euclidean distance. First, a naive approach  
 487 that we consider is creating a lookup table per weight to match a feature value with the multiplication  
 488 of that feature with the specific weight. In this approach, the feature computation microengines would  
 489 be responsible for doing the multiplication lookup based on the weight and feature value, and the result  
 490 aggregation microengines would sum up the multiplied weight-feature values to obtain the final result.  
 491 However, this approach may not be feasible for problems that require a large number of weights due to  
 492 memory constraints on the SmartNIC.

493 An alternative approach would be to consider using natural logarithm and the exponent function (i.e.,  $e^x$ ).  
 494 Instead of storing a lookup table per feature, two lookup tables can be stored to approximately compute

495 the natural logarithm and exponent of the feature values based on their fixed range (i.e., we assume each  
496 feature is in the range  $[0, 1]$  in Section 3). Then, each feature computation microengine would be operate  
497 on a specific feature by conducting a lookup for the natural logarithm of the feature, taking the sum of the  
498 natural logarithms of the weights and the feature value, and conduct a lookup of the exponent of the sum of  
499 the natural logarithm values. While this approach may resolve issues with the memory constraint, a large  
500 number of features requires multiple lookups to the memory region holding the tables (i.e., CLS or IMEM)  
501 that can congest the CPP bus. So, to avoid this issue, the lookup tables could be replaced with first-order  
502 Taylor approximations of the natural logarithm and exponent functions for a specific number of reference  
503 points in the range  $[0, 1]$ , where the Taylor approximations can be represented using additions and bit shifts.  
504 At the same time though, the use of first-order Taylor approximations can result in more erroneous model  
505 predictions.

506 More recently though, work by Blalock and Gutttag (2021) proposed a novel technique for matrix-matrix  
507 multiplication that used locality-sensitive hashing to determine suitable functions (denote as  $g(A)$ ) that  
508 can be executed efficiently using balanced binary regression trees. Based on their findings and our current  
509 implementation for decision tree inference, we believe their approach to be a more promising direction for  
510 executing matrix-vector multiplication on a SmartNIC.

## 511 6.6 Productionization and Scaling

512 When deploying decision tree models in the real world, we consider factors such as model updates and  
513 scaling. Concerning model updates (i.e., models retrained on larger datasets), we still limit the number of  
514 leaf nodes to 256, which may or may not be helpful as a regularization technique to prevent decision tree  
515 overfitting. In order to deploy a new decision tree, the original decision tree (i.e., the model firmware file)  
516 needs to be unloaded from the SmartNIC, and then the firmware for the new decision tree can be loaded  
517 onto the SmartNIC. This means that SmartNIC would be inactive while unloading the old decision tree and  
518 loading the new decision tree. So, inference requests can not be handled by a SmartNIC during that time.

519 For scaling, we primarily focus on the first model deployment scenario (model fits on an island). We  
520 do not believe much scaling of SmartNIC resources can be done as the entire card is required for one  
521 instantiation of the model. In the first deployment scenario, though, based on the amount of network traffic,  
522 firmware can be developed to set aside some islands for the decision tree model and others for other tasks.  
523 However, like the model update case, old firmware would need to be unloaded and new firmware loaded.  
524 So, inference requests can not be handled by a SmartNIC while unloading old and loading new firmware.

## 525 6.7 Hardware Improvements

526 Considering SmartNICs with Micro-C support that possess additional hardware capabilities, we first note  
527 SmartNICs with additional programmable flow-processing microengines. A SmartNIC can include more  
528 programmable microengines in two ways: additional islands or microengines per island. With additional  
529 microengines per island, we expect a further reduction in latency for all three model deployment scenarios.  
530 This would happen because more models would be able to fit on an island, which would remove the need  
531 for communication between the microengines on different islands. Communication between microengines  
532 on different islands is more expensive than between microengines on the same island. With additional  
533 islands, we expect a further increase in throughput for model deployment scenario one (the model can fit  
534 on one island). More islands mean more instances of the model that can be instantiated on the SmartNIC,  
535 which would allow it to meet more inference requests. In either scenario, we expect larger trees (with

536 respect to the number of features) to be more easily deployed, given that each microengine would be  
537 responsible for analyzing fewer features.

538 Next, we consider a SmartNIC with additional transfer registers per microengine. More transfer registers  
539 means that fewer microengines would be needed to perform result aggregation based on the feature analysis  
540 conducted by the feature computation microengines. With greater availability of microengines, we could  
541 likely deploy larger trees on the SmartNIC.

542 Lastly, another hardware improvement we consider is the addition of one or more FPUs. Adding FPUs  
543 resolves the issue with minor effectiveness degradation that we observe with the current iteration of our  
544 work while likely maintaining similar latency and throughput performance that we've observed in this  
545 work.

## 7 RELATED WORK

546 Recently, several works have leveraged programmable data planes to make aspects of machine learning  
547 more efficient. These works can be split into two categories: model training and model inference. We  
548 focus on the latter. Within model inference, research efforts are focused on leveraging how entire or  
549 portions of the machine learning model inference process can be offloaded onto programmable data planes  
550 while maintaining adequate model performance. These implementations are most commonly conducted on  
551 programmable switches and SmartNICs.

### 552 7.1 Programmable Switch

553 For works that map machine learning models onto programmable switches, we generally observe a focus  
554 on specific models used to address certain tasks. We first note Net2Net (Siracusano and Bifulco, 2018)  
555 that proposed quantizing neural networks into binary neural networks, since they require operations that  
556 are readily available on modern switching chips. Rather than quantizing a trained neural network into a  
557 binary neural network, Qin et al. (2020) directly trained binary neural networks and mapped them onto the  
558 data planes of programmable switches using P4 to handle the network intrusion detection use case. As an  
559 alternative to binary neural network quantization, Dao et al. (2021) used neuron pruning to map neural  
560 networks onto programmable switches for the network intrusion use case. While the above neural network  
561 works solely considered a single programmable switch for deployment, Saquetti et al. (2021) implemented  
562 a neural network neuron distribution method across multiple programmable switches and coordinated  
563 inference of the model between the switches to optimize resource usage. Similarly, JointNIDS (Dao  
564 and Lee, 2022) also employed a distributed neural network inference approach. But, the neural network  
565 intrusion detection models were split into two sequential submodels with overlapping hidden units and  
566 mapped the submodels onto two programmable switches. The authors assigned one programmable switch  
567 to detect major network attacks, while the second handled the more subtle aspects of network traffic  
568 classification. Rather than staying within the constraints of off-the-shelf programmable network devices,  
569 Taurus (Swamy et al., 2022) extended the PISA architecture of programmable switches by adding custom  
570 hardware to support parallelism and additional operations (i.e., multiplication, nonlinear operations) needed  
571 to run neural network inference without any quantization.

572 Regarding tree-based models, pForest (Busse-Grawitz et al., 2019) developed a optimization technique to  
573 map a random forest classifier to a programmable switch in P4 for network flow classification. Furthermore,  
574 this approach adaptively switches out the current classifier with others based on the network flows observed.  
575 In addition to mapping a random forest classifier, Planter (Zheng and Zilberman, 2021) mapped a xgboost

576 and isolation forest classifier to programmable switches using overlapping trees to overcome some of  
577 the inefficiencies observed in pForest (Busse-Grawitz et al., 2019). Similar to pForest (Busse-Grawitz  
578 et al., 2019), SMASH (Kamath and Sivalingam, 2021) also focused on the network flow classification task  
579 and used an improved hash-and-store algorithm with a decision tree model for early flow classification.  
580 Also working with decision trees, pHeavy (Zhang et al., 2021) implemented trained decision trees on the  
581 data plane to reduce the overhead involved with communicating to the control plane in Software-Defined  
582 Networking (SDN) when classifying highly-congested network flows. In contrast to the other tree-based  
583 model offloading approaches that focus on network-related use cases, NetPixel (Siddique et al., 2021)  
584 implemented decision trees on P4 programmable switches to handle image classification. To address some  
585 of the issues with deployment of decision trees and other machine learning algorithms onto programmable  
586 data planes, Mousika (Xie et al., 2022) introduced a teacher-student knowledge distillation approach to  
587 translate machine learning models to binary decision trees, which are more suitable for mapping onto the  
588 data plane.

589 On top of supervised machine learning, the deployment of unsupervised learning algorithms onto  
590 programmable switches has also been explored. Clustreams (Friedman et al., 2021) used a combination  
591 of the quadtree data structure and a match+action table stored in Ternary Content Addressable Memory  
592 (TCAM) to cluster network traffic efficiently. In addition, ACC-Turbo (Alcoz et al., 2022) redesigned the  
593 original Aggregate-based Congestion Control (ACC) approach using online clustering and a scheduling  
594 algorithm to mitigate pulse-wave DDoS attacks.

595 Unlike the works above that focus on a specific machine learning algorithm type (i.e., neural networks,  
596 tree-based models, clustering, etc.), IIsy (Xiong and Zilberman, 2019) introduced mapping schemes for  
597 several machine learning algorithms, such as decision trees, k-means, naive bayes, and support vector  
598 machines, to the data plane using the match-action pipeline in programmable switches. Also, Hong et al.  
599 (2024) developed a feature engineering and model deployment strategy for tree-based models (i.e., decision  
600 trees, random forests, xgboost), k-nearest neighbor, and k-means to handle the high-frequency stock market  
601 trading task.

## 602 7.2 SmartNIC

603 Similar to the works that address deployment of machine learning model inference onto programmable  
604 switches, we see works about model inference onto SmartNICs that also consider neural networks and  
605 decision trees used for particular applications. Using the approach proposed by Net2Net (Siracusano  
606 and Bifulco, 2018), BaNaNa split (Sanvito et al., 2018) accelerated the inference of neural networks by  
607 splitting a neural network at its fully-connected layers, sending all prior layers to the host system's CPU  
608 for inference, and quantizing the fully-connected layers to run portion of the inference on the host system's  
609 SmartNIC. Different from the other works that primarily look into P4 implementations, N3IC (Siracusano  
610 et al., 2022) used Micro-C and P4 to map binary neural networks onto a greater variety of targets (i.e.,  
611 SmartNICs) for traffic analysis use cases. Regarding tree-based models, Xavier et al. (2021) presented  
612 a framework for deploying decision tree models onto SmartNICs in P4. The authors demonstrated that  
613 their framework can achieve high accuracy (above 95%) in a network intrusion detection use case. While  
614 similar to our work, we note that ML-NIC works on a greater variety of use cases outside of network traffic  
615 analysis. Furthermore, ML-NIC's Micro-C implementations can parallelize the model inference process,  
616 which is not possible with P4. While the works on machine learning inference offloading for SmartNICs  
617 do not cover unsupervised learning to our knowledge, they do address traditional reinforcement learning.



618 Opal (Simpson and Pezaros, 2022) implemented online reinforcement learning onto a SmartNIC data plane,  
619 relying on classical reinforcement algorithms such as Sarsa (Sutton, 2018) and avoiding neural networks.

## 8 FUTURE DIRECTIONS

### 620 8.1 Implementing Additional Models

621 As mentioned in Section 6, our next step is to expand our framework to other machine learning algorithms,  
622 such as support vector machines and neural networks. We find that implementing the approximate matrix-  
623 matrix multiplication approach developed by Blalock and Gutttag (2021) to be means of achieving this goal.  
624 In addition to matrix-matrix multiplication, there are floating point operations that are often performed  
625 for various models such as neural networks. Thus, in order to implement such models, future work may  
626 involve, either adding hardware support for these operations or using quantized operations as a default. We  
627 discuss the potential future work in this area in Section 8.2.

### 628 8.2 Improving the Floating-Point Representation

629 While our current approach leverages fixed-point 16-bit features and produces comparable model  
630 effectiveness scores to the baselines, it is possible that the proposed float-representation scheme can  
631 be improved without adding more hardware. One potential future direction to be explored is to look into  
632 the posit representation Gustafson and Yonemoto (2017) as a potential alternative, since it resolves the  
633 issue of NaN quantities observed in the standard floating-point representation. Also, while implementing  
634 the standard floating-point representation on the SmartNIC may seem like a viable solutions, we believe  
635 that the float-pointing representation standard would introduce additional latency due to the additional  
636 computation spent managing mantissa bits, exponent bits, and NaN quantities. In addition, adding any  
637 additional hardware support may cause added cost and energy consumption of SmartNIC, which defeats  
638 the purpose of using the SmartNIC in the first place. Thus, a software / algorithmic based approaches for  
639 performing floating point operations will be a great direction for future work.

### 640 8.3 Automating the Model Deployment

641 Furthermore, despite automating the process of decision tree logic extraction, the process of building the  
642 model mapping still mostly requires the developer to manually allocate cores as one of packet collection,  
643 feature computation, or result aggregation. We think the model mapping component can be made more  
644 efficient with additional code that considers the computation constraints of the SmartNIC and presents  
645 a mapping scheme to remove some of the tedious work in deploying a model onto the SmartNIC. Thus,  
646 a direction for future work may involve building a more sophisticated system that can perform model  
647 compilation, optimization and deployment automatically to the SmartNIC. This work can be further  
648 strengthened by adding a notion of distributed deployment and model inference across multiple SmartNICs  
649 located on multiple server.

### 650 8.4 Utilizing Different Types of SmartNICs

651 While this work primarily utilizes ASIC-based SmartNIC, it is possible to implement similar work on  
652 other ASIC-based and other types of SmartNICs, such as FPGA-based SmartNICs. While the optimizations  
653 we performed in this paper is specific to Netronome SmartNIC, the overall idea of mapping memory into  
654 different SmartNIC region is quite generic. Thus, a potential valuable future work is to perform similar



655 optimization strategies across different types of hardware implementations to understand the similarities  
656 and differences in the effectiveness of the proposed optimization and compilation strategies.

## 9 CONCLUSION

657 Low-latency model inference is a necessity for many time-sensitive machine learning applications. This  
658 paper demonstrates that ML-NIC is a suitable framework for performing machine learning model inference.  
659 Our evaluation of the first iteration of ML-NIC shows that it can deploy larger models than the state-of-the-  
660 art SmartNIC approach, can produce predictions at faster speeds with a minor loss in model effectiveness  
661 compared to the CPU solution, and is robust to high network utilization and CPU loads.

## CONFLICT OF INTEREST STATEMENT

662 The authors declare that the research was conducted in the absence of any commercial or financial  
663 relationships that could be construed as a potential conflict of interest.

## AUTHOR CONTRIBUTIONS

664 SC conceived the presented idea and provided computing resources. DCA and SC supervised the project,  
665 designed experiments, and validated the findings. RK conducted experiments, collected and analyzed  
666 data, and wrote the initial draft of the manuscript. All authors were involved in reviewing and revising the  
667 manuscript.

## DATA AVAILABILITY STATEMENT

668 The datasets analyzed for this study can be found in the University of California, Irvine Machine  
669 Learning Repository [<https://archive.ics.uci.edu/>] and University of New Brunswick Canadian Institute for  
670 Cybersecurity [<https://www.unb.ca/cic/>].

## REFERENCES

- 671 Alcoz, A. G., Strohmeier, M., Lenders, V., and Vanbever, L. (2022). Aggregate-based congestion control  
672 for pulse-wave ddos defense. In *Proceedings of the ACM SIGCOMM 2022 Conference*. 693–706
- 673 Blalock, D. and Gutttag, J. (2021). Multiplying matrices without multiplying. In *International Conference*  
674 *on Machine Learning* (PMLR), 992–1004
- 675 Busse-Grawitz, C., Meier, R., Dietmüller, A., Bühler, T., and Vanbever, L. (2019). pforest: In-network  
676 inference with random forests. *arXiv preprint arXiv:1909.05680*
- 677 Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd*  
678 *acm sigkdd international conference on knowledge discovery and data mining*. 785–794
- 679 Chen, Y., Luo, T., Liu, S., Zhang, S., He, L., Wang, J., et al. (2014). Dadiannao: A machine-learning  
680 supercomputer. In *2014 47th Annual IEEE/ACM International Symposium on Microarchitecture* (IEEE),  
681 609–622
- 682 Choquette, J., Giroux, O., and Foley, D. (2018). Volta: Performance and programmability. *Ieee Micro* 38,  
683 42–52
- 684 Corigine (2020). *Corigine NFP-4000 Flow Processor*. Tech. rep., Corigine
- 685 Dang, H. T., Bressana, P., Wang, H., Lee, K. S., Zilberman, N., Weatherspoon, H., et al. (2020). P4xos:  
686 Consensus as a network service. *IEEE/ACM Transactions on Networking* 28, 1726–1738

- 687 Dao, T.-N., Hoang, V.-P., Ta, C. H., et al. (2021). Development of lightweight and accurate intrusion  
688 detection on programmable data plane. In *2021 International Conference on Advanced Technologies for*  
689 *Communications (ATC)* (IEEE), 99–103
- 690 Dao, T.-N. and Lee, H. (2022). Jointnids: Efficient joint traffic management for on-device network intrusion  
691 detection. *IEEE Transactions on Vehicular Technology* 71, 13254–13265
- 692 Fowers, J., Ovtcharov, K., Papamichael, M., Massengill, T., Liu, M., Lo, D., et al. (2018). A configurable  
693 cloud-scale dnn processor for real-time ai. In *2018 ACM/IEEE 45th Annual International Symposium on*  
694 *Computer Architecture (ISCA)* (IEEE), 1–14
- 695 Friedman, R., Goaz, O., and Rottenstreich, O. (2021). Clustreams: Data plane clustering. In *Proceedings*  
696 *of the ACM SIGCOMM Symposium on SDN Research (SOSR)*. 101–107
- 697 Gustafson, J. L. and Yonemoto, I. T. (2017). Beating floating point at its own game: Posit arithmetic.  
698 *Supercomputing frontiers and innovations* 4, 71–86
- 699 He, Z., Sidler, D., István, Z., and Alonso, G. (2018). A flexible k-means operator for hybrid databases.  
700 In *2018 28th International Conference on Field Programmable Logic and Applications (FPL)* (IEEE),  
701 368–3683
- 702 Hong, X., Zheng, C., Zohren, S., and Zilberman, N. (2024). Accelerating machine learning for trading using  
703 programmable switches (IOS Press), *Frontiers in Artificial Intelligence and Applications*, 3429–3436
- 704 Jin, X., Li, X., Zhang, H., Soulé, R., Lee, J., Foster, N., et al. (2017). Netcache: Balancing key-value stores  
705 with fast in-network caching. In *Proceedings of the 26th Symposium on Operating Systems Principles*.  
706 121–136
- 707 Jouppi, N. P., Young, C., Patil, N., Patterson, D., Agrawal, G., Bajwa, R., et al. (2017). In-datacenter  
708 performance analysis of a tensor processing unit. In *Proceedings of the 44th annual international*  
709 *symposium on computer architecture*. 1–12
- 710 Kamath, R. and Sivalingam, K. M. (2021). Machine learning based flow classification in dcns using p4  
711 switches. In *2021 International Conference on Computer Communications and Networks (ICCCN)*  
712 (IEEE), 1–10
- 713 Kim, C., Sivaraman, A., Katta, N., Bas, A., Dixit, A., Wobker, L. J., et al. (2015). In-band network  
714 telemetry via programmable dataplanes. In *ACM SIGCOMM*. vol. 15, 1–2
- 715 Langlet, J. (2019). Towards machine learning inference in the data plane
- 716 Leo Breiman, C. J. S., Jerome Friedman and Olshen, R. (1984). *Classification and Regression Trees*  
717 (Chapman and Hall/CRC)
- 718 Liu, Y. H. (2017). *Python machine learning by example* (Packt Publishing Ltd)
- 719 Neutronome Systems, I. (2024). *Agilio CX 2x25GbE SmartNIC*. Tech. rep., Neutronome Systems, Inc.
- 720 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn:  
721 Machine learning in python. *the Journal of machine Learning research* 12, 2825–2830
- 722 Qin, Q., Poularakis, K., Leung, K. K., and Tassiulas, L. (2020). Line-speed and scalable intrusion detection  
723 at the network edge via federated learning. In *2020 IFIP Networking Conference (Networking)* (IEEE),  
724 352–360
- 725 Quinlan, J. R. (1986). Induction of decision trees. *Machine learning* 1, 81–106
- 726 Rodríguez-Luján, I., Fonollosa, J., Vergara, A., Homer, M. L., and Huerta, R. (2014). On the calibration  
727 of sensor arrays for pattern recognition using the minimal number of experiments. *Chemometrics and*  
728 *Intelligent Laboratory Systems* 130, 123–134
- 729 [Dataset] Roth, D. (2016). Decision trees
- 730 Sanvito, D., Siracusano, G., and Bifulco, R. (2018). Can the network be the ai accelerator? In *Proceedings*  
731 *of the 2018 Morning Workshop on In-Network Computing*. 20–25

- 732 Saquetti, M., Canofre, R., Lorenzon, A. F., Rossi, F. D., Azambuja, J. R., Cordeiro, W., et al. (2021).  
733 Toward in-network intelligence: Running distributed artificial neural networks in the data plane. *IEEE*  
734 *Communications Letters* 25, 3551–3555
- 735 Sharafaldin, I., Lashkari, A. H., Ghorbani, A. A., et al. (2018). Toward generating a new intrusion detection  
736 dataset and intrusion traffic characterization. *ICISSp* 1, 108–116
- 737 Siddique, H., Neves, M., Kuzniar, C., and Haque, I. (2021). Towards network-accelerated ml-based  
738 distributed computer vision systems. In *2021 IEEE 27th International Conference on Parallel and*  
739 *Distributed Systems (ICPADS)* (IEEE), 122–129
- 740 Simpson, K. A. and Pezaros, D. P. (2022). Revisiting the classics: Online rl in the programmable dataplane.  
741 In *NOMS 2022-2022 IEEE/IFIP Network Operations and Management Symposium* (IEEE), 1–10
- 742 Siracusano, G. and Bifulco, R. (2018). In-network neural networks. *arXiv preprint arXiv:1801.05731*
- 743 Siracusano, G., Galea, S., Sanvito, D., Malekzadeh, M., Antichi, G., Costa, P., et al. (2022). Re-architecting  
744 traffic analysis with neural network interface cards. In *19th USENIX Symposium on Networked Systems*  
745 *Design and Implementation (NSDI 22)*. 513–533
- 746 [Dataset] Srinivasan, A. (1993). Statlog (Landsat Satellite). UCI Machine Learning Repository. DOI:  
747 <https://doi.org/10.24432/C55887>
- 748 Sutton, R. S. (2018). Reinforcement learning: an introduction. *A Bradford Book*
- 749 Swamy, T., Rucker, A., Shahbaz, M., Gaur, I., and Olukotun, K. (2022). Taurus: a data plane architecture  
750 for per-packet ml. In *Proceedings of the 27th ACM International Conference on Architectural Support*  
751 *for Programming Languages and Operating Systems*. 1099–1114
- 752 Tong, D., Qu, Y. R., and Prasanna, V. K. (2017). Accelerating decision tree based traffic classification on  
753 fpga and multicore platforms. *IEEE Transactions on Parallel and Distributed Systems* 28, 3046–3059
- 754 Wray, S. (2014). *The Joy of Micro-C*. Tech. rep., Netronome Systems, Inc.
- 755 Xavier, B. M., Guimarães, R. S., Comarela, G., and Martinello, M. (2021). Programmable switches for  
756 in-networking classification. In *IEEE INFOCOM 2021-IEEE Conference on Computer Communications*  
757 (IEEE), 1–10
- 758 Xie, G., Li, Q., Dong, Y., Duan, G., Jiang, Y., and Duan, J. (2022). Mousika: Enable general in-network  
759 intelligence in programmable switches by knowledge distillation. In *IEEE INFOCOM 2022-IEEE*  
760 *Conference on Computer Communications* (IEEE), 1938–1947
- 761 Xiong, Z. and Zilberman, N. (2019). Do switches dream of machine learning? toward in-network  
762 classification. In *Proceedings of the 18th ACM workshop on hot topics in networks*. 25–33
- 763 Yilmaz, C., Kahraman, H. T., and Söyler, S. (2018). Passive mine detection and classification method  
764 based on hybrid model. *IEEE Access* 6, 47870–47888. doi:10.1109/ACCESS.2018.2866538
- 765 Zhang, B., Kannan, R., Prasanna, V., and Busart, C. (2023). Accelerating gnn-based sar automatic target  
766 recognition on hbm-enabled fpga. In *2023 IEEE High Performance Extreme Computing Conference*  
767 (HPEC) (IEEE), 1–7
- 768 Zhang, X., Cui, L., Tso, F. P., and Jia, W. (2021). pheavy: Predicting heavy flows in the programmable  
769 data plane. *IEEE Transactions on Network and Service Management* 18, 4353–4364
- 770 Zheng, C. and Zilberman, N. (2021). Planter: seeding trees within switches. In *Proceedings of the*  
771 *SIGCOMM '21 Poster and Demo Sessions* (New York, NY, USA: Association for Computing Machinery),  
772 SIGCOMM '21. 12–14. doi:10.1145/3472716.3472846

## APPENDIX

773 **A Packet Collection Implementation**

774 An example implementation for packet collection is provided below. At the start of the main function,  
 775 the microengine signals the CTM packet engine that it is ready to receive a packet. Once the MAC  
 776 block has received and preprocessed a packet, it is transferred to the microengine's corresponding CTM  
 777 buffer and an MU buffer (if needed). The remainder of the code then validates the packet as a model  
 778 input packet, parses the features embedded in the packet payload, and transfers the features to various  
 779 microengines. After transferring the packets, the result aggregation microengine waits to be signaled that  
 780 the final classification is produced. Then, the original model input packet is modified with the classification  
 781 result, some post-processing is performed, and the CTM packet engine is informed that the packet is ready  
 782 to be sent out.

**Listing 1.** Example Packet Collection Code

```

783 #include <nfp.h>
784 #include <nfp/remote_me.h>
785 #include <nfp/mem_bulk.h>
786 ... /* Some more packet imports */
787 #include <net/hdr_ext.h>
788 #include "config.h"
789
790 /*
791 Read packet data from memory in xfer registers as a two byte
792 offset so that the payload get aligned to a word boundary. This
793 makes the header extraction code more efficient.
794 */
795 #define PKT_START_OFF 2
796
797
798 /* A structure used for extracting , the different protocol header */
799 struct pkt_hdr {
800     struct eth_hdr eth;
801     struct ip4_hdr ip4;
802 };
803
804 // Signals for interacting with other MEs
805 SIGNAL local_signal;
806
807 __remote SIGNAL remote_signal1;
808 __remote __xread uint32_t remote_xfer1;
809
810 __remote SIGNAL remote_signal2;
811 __remote __xread uint32_t remote_xfer2;
812
813 ... /* Repeat similar variable for each feature */
814
```

```

815 __visible SIGNAL complete_signal;
816 __visible __xread uint32_t complete_xfer;
817
818 __intrinsic void proc_pkt(__mem40 char *buf_addr, __gpr uint32_t
819   buf_off) {
820
821   // Extract the headers first before the payload
822   __xread uint32_t pkt_buf[17];
823   __xwrite uint32_t data_out[15];
824   __lmem uint32_t src_buf[13]; // just store the headers
825   __gpr struct pkt_hdr eh;
826   __gpr uint32_t csum_prepend;
827   __gpr int src_off = buf_off;
828   __gpr int res;
829   __gpr int len;
830
831   // Extract the packet
832   mem_read32(pkt_buf, buf_addr + buf_off - PKT_START_OFF, sizeof(
833     pkt_buf));
834
835   /* Copy xfer register to a Local Memory buffer for easier
836     extraction */
837   reg_cp(src_buf, pkt_buf, sizeof(src_buf));
838   src_off = PKT_START_OFF;
839
840   /*
841    * Handle the checksum prepend if configured
842    */
843   #ifndef CFG_RX_CSUM_PREPEND
844     /* read the MAC parsing info for CSUM (first 4B are timestamp) */
845     csum_prepend = pkt_csum_read(pkt_buf, PKT_START_OFF + 4);
846     src_off += MAC_PREPEND_BYTES;
847
848     if (NFP_MAC_RX_CSUM_L3_SUM_of(csum_prepend) ==
849         NFP_MAC_RX_CSUM_L3_IPV4_FAIL) {
850       /* L3 checksum is wrong */
851       return;
852     }
853
854     if ((NFP_MAC_RX_CSUM_L4_SUM_of(csum_prepend) ==
855         NFP_MAC_RX_CSUM_L4_TCP_FAIL) ||
856         (NFP_MAC_RX_CSUM_L4_SUM_of(csum_prepend) ==
857         NFP_MAC_RX_CSUM_L4_UDP_FAIL)) {
858       /* L4 checksum is wrong */
859       return;

```

```

860     }
861 #endif
862
863     // Layer 2 Extraction
864     res = he_eth(src_buf, src_off, &eh.eth);
865     len = HE_RES_LEN_of(res);
866     src_off += len;
867
868     // Check if ethernet header is correct
869     if (((__gpr uint16_t *)&eh.eth.dst)[0] == 0x0015
870         && ((__gpr uint16_t *)&eh.eth.dst)[1] == 0x4d13
871         && ((__gpr uint16_t *)&eh.eth.dst)[2] == 0x79ac
872         && net_eth_is_uc_addr(&eh.eth.src)
873         && eh.eth.type == 0x4d49) {
874
875         // Edit Ethernet Header
876         ((__lmem struct eth_hdr *)((__lmem char *)src_buf+10))->dst =
877             eh.eth.src;
878         ((__lmem struct eth_hdr *)((__lmem char *)src_buf+10))->src =
879             eh.eth.dst;
880         ((__lmem struct eth_hdr *)((__lmem char *)src_buf+10))->type
881             = 0x4d4f;
882
883         // Edit IP Header
884         res = he_ip4(src_buf, src_off, &eh.ip4);
885         ((__lmem struct ip4_hdr *) (src_buf+6))->dst = eh.ip4.src;
886         ((__lmem struct ip4_hdr *) (src_buf+6))->src = eh.ip4.dst;
887
888         // Copy header and payload to data_out
889         reg_cp(data_out, src_buf+2, 44);
890         reg_cp(data_out+11, pkt_buf+13, 16);
891
892         // Inform ME 1 that data is ready
893         remote_me_reg_write_signal_remote(data_out+11, 32, 1, 0,
894             __xfer_reg_number(&
895                 remote_xfer1, __nfp_meid
896                 (32, 1)) + (__ctx() * 32)
897             ,
898             sizeof(remote_xfer1), &
899             local_signal);
900
901         // Inform ME 2 that data is ready
902         remote_me_reg_write_signal_remote(data_out+11, 32, 2, 0,

```

```

903                                     __xfer_reg_number(&
904                                     remote_xfer2 , __nfp_meid
905                                     (32, 2)) + (__ctx() * 32)
906                                     ,
907                                     sizeof(remote_xfer2), &
908                                     local_signal);
909
910     ... /* Repeat microengine data transfer for each feature */
911
912     // Wait for ME 9 to finish
913     __wait_for_all(&complete_signal);
914     reg_cp(data_out+11, &complete_xfer , sizeof(complete_xfer));
915     mem_write64(data_out , buf_addr + buf_off + 6, 48);
916 }
917 }
918
919 int main(void) {
920     __mem40 char *pbuf;
921     __xread struct nbi_meta_catamaran nbi_meta;
922     __xread struct nbi_meta_pkt_info *pi = &nbi_meta.pkt_info;
923     __gpr int in_port , pkt_off;
924     __gpr struct pkt_ms_info msi;
925
926     __assign_relative_register(&local_signal , 1);
927     __assign_relative_register(&complete_signal , 10);
928     __implicit_write(&complete_signal , sizeof(complete_signal));
929     __implicit_write(&complete_xfer , sizeof(complete_xfer));
930
931     for (;;) {
932         /* Receive packet */
933         pkt_nbi_recv(&nbi_meta , sizeof(nbi_meta));
934         in_port = MAC_TO_PORT(nbi_meta.port);
935         pbuf = pkt_ctm_ptr40(pi->isl , pi->pnum, 0);
936
937         /* Collect features from pkt */
938         pkt_off = PKT_NBI_OFFSET;
939         proc_pkt(pbuf , pkt_off);
940
941         /* Send packet */
942         // Write the MAC egress CMD and adjust offset and len as
943         needed
944         pkt_off += MAC_PREPEND_BYTES;
945         pkt_mac_egress_cmd_write(pbuf , pkt_off , 1, 1);
946
947         pkt_off -= 4;

```

```

948     msi = pkt_msd_write(pbuf, pkt_off);
949     pkt_nbi_send(pi->isl, pi->pnum, &msi,
950                pi->len - MAC_PREPEND_BYTES + 4,
951                NBI, PORT_TO_TMQ(in_port),
952                nbi_meta.seqr, nbi_meta.seq, PKT_CTM_SIZE_256);
953 }
954
955     return 0;
956 }

```

## 957 B Feature Computation Implementation

958 The code listing below exemplifies how we program a microengine for feature computation. Once the  
959 packet collection microengine signals the microengine and receives a feature to analyze, it executes the  
960 logic learned by the trained decision for that specific feature as a series of if-else statements to determine  
961 which leaf nodes can be reached based on the value taken on by the feature in question. Once the possible  
962 leaf nodes are calculated, the feature computation microengine transmits the possible leaf nodes to the  
963 result aggregation microengine.

### Listing 2. Example Feature Computation Code

```

964 #include <nfp.h>
965 #include <stdint.h>
966 #include <std/reg_utils.h>
967 #include <nfp/me.h>
968 #include <nfp/remote_me.h>
969 #include <nfp/cls.h>
970
971 __visible SIGNAL remote_signal2;
972 __visible __xread uint32_t remote_xfer2;
973
974 __remote SIGNAL result_signal2;
975 __remote __xread uint32_t result_xfer2[4];
976
977 SIGNAL local_signal;
978
979 int main(void) {
980     __gpr uint32_t feature;
981     __gpr uint32_t result_gpr[4];
982     __xwrite uint32_t result_write[4];
983
984     __assign_relative_register(&remote_signal2, 1);
985     __assign_relative_register(&local_signal, 3);
986     __implicit_write(&remote_signal2, sizeof(remote_signal2));
987     __implicit_write(&remote_xfer2, sizeof(remote_xfer2));
988
989     for (;;) {

```



```

990     reg_set(result_gpr , 0x3ffff , sizeof(uint32_t));
991     reg_set(result_gpr+1, 0xffffffff , sizeof(uint32_t));
992     reg_set(result_gpr+2, 0xffffffff , sizeof(uint32_t));
993     reg_set(result_gpr+3, 0xffffffff , sizeof(uint32_t));
994
995     __wait_for_all(&remote_signal2);
996     reg_cp(&feature , &remote_xfer2 , sizeof(feature));
997     feature &= 0xffff;
998
999     // Decision Tree logic
1000    if (!(feature <= 5585))
1001        result_gpr[3] &= ~(1 << 3);
1002
1003    if (!(feature <= 5585))
1004        result_gpr[3] &= ~(1 << 4);
1005
1006    if (!(feature <= 5585))
1007        result_gpr[3] &= ~(1 << 5);
1008
1009    if (!(feature <= 5585))
1010        result_gpr[3] &= ~(1 << 6);
1011
1012    if (!(feature <= 5585 && feature <= 4841))
1013        result_gpr[3] &= ~(1 << 7);
1014
1015    if (!(feature <= 5585 && feature > 4841))
1016        result_gpr[3] &= ~(1 << 8);
1017
1018    if (!(feature <= 5585 && feature <= 3351 && feature <= 1117))
1019        result_gpr[3] &= ~(1 << 9);
1020
1021    ... /* Remaining decision tree feature logic */
1022
1023    // Transfer result to result aggregation microengine
1024    reg_cp(&result_write , &result_gpr , sizeof(result_write));
1025    remote_me_reg_write_signal_remote(&result_write , 32, 9, 0,
1026                                     __xfer_reg_number(&
1027                                                         result_xfer2 , __nfp_meid
1028                                                         (32, 9)) + (__ctx() * 32)
1029                                     ,
1030                                     sizeof(result_xfer2) , &
1031                                     local_signal);
1032 }
1033
1034 return 0;

```

1035 }

### 1036 C Result Aggregation Implementation

1037 Below, we provide an example implementation for result aggregation. Once the microengine receives all  
 1038 possible leaf nodes based on the values taken on by each feature from the feature computation microengines,  
 1039 it takes the single possible leaf node predicted by all the feature computation microengines and searches  
 1040 for the corresponding class label. The result aggregation microengine transmits the label to the packet  
 1041 collection microengine to complete the inference process.

#### Listing 3. Example Result Aggregation Code

```

1042 #include <nfp.h>
1043 #include <stdint.h>
1044 #include <std/reg_utils.h>
1045 #include <nfp/me.h>
1046 #include <nfp/remote_me.h>
1047 #include <nfp/cls.h>
1048
1049 __visible SIGNAL result_signal1;
1050 __visible __xread uint32_t result_xfer1[4];
1051
1052 ... /* Repeat similar variables for each feature */
1053
1054 // For transfer back to packet collection microengine
1055 __remote SIGNAL complete_signal;
1056 __remote __xread uint32_t complete_xfer;
1057
1058 SIGNAL local_signal;
1059
1060 int main(void) {
1061     __gpr uint32_t results_gpr1[4];
1062     __gpr uint32_t results_gpr2[4];
1063     ... /* Repeat similar variables for each feature */
1064     __xwrite uint32_t final_result_write;
1065
1066     __lmem uint8_t path_class[] = {1, 4, 4, 1, 4, 1, 1, 4, 5, 1, 5,
1067     1, 5, 5, 3, 3, 4, 1, 3, ...};
1068
1069     __assign_relative_register(&result_signal1, 2);
1070     .../* Repeat similar function call for each feature */
1071     __assign_relative_register(&local_signal, 10);
1072
1073     __implicit_write(&result_signal1, sizeof(result_signal1));
1074     __implicit_write(result_xfer1, sizeof(result_xfer1));
1075     .../* Repeat similar function calls for each feature */
1076

```

```

1077     for (;;) {
1078         __wait_for_all(&result_signal1 , &result_signal2 , ...);
1079
1080         // Collect all the possible paths
1081         reg_cp(results_gpr1 , result_xfer1 , sizeof(results_gpr1));
1082         ... /* Repeat similar function call for each feature */
1083
1084         // Compute final path
1085         results_gpr8[0] = results_gpr1[0] & ...;
1086         ... /* Repeat similar computation for each 32-bit word */
1087
1088         // Calculate leaf node index
1089         if (results_gpr8[0] != 0)
1090             results_gpr8[0] = ffs(results_gpr8[0]) + 96;
1091         else if (results_gpr8[1] != 0)
1092             results_gpr8[0] = ffs(results_gpr8[1]) + 64;
1093         ... /* Repeat similar computation for each 32-bit word */
1094         else
1095             results_gpr8[0] = ffs(results_gpr8[3]);
1096
1097         // Get class prediction
1098         results_gpr8[0] = (uint32_t)(path_class[results_gpr8[0]]);
1099         reg_cp(&final_result_write , results_gpr8 , sizeof(
1100             final_result_write));
1101
1102         // Send final result to packet collection core
1103         remote_me_reg_write_signal_remote(&final_result_write , 32, 0,
1104             0,
1105
1106                                     __xfer_reg_number(&
1107                                         complete_xfer , __nfp_meid
1108                                         (32, 0)) + (__ctx() * 32)
1109                                     ,
1110                                     sizeof(complete_xfer) , &
1111                                     local_signal);
1112     }
1113     return 0;
1114 }

```

## 1115 D Model Deployment Implementation

1116 Once all the Micro-C code files have been created, the following provides an example of deploying the  
1117 mapped decision tree onto the SmartNIC. Each nfcc command processes a Micro-C code file and builds  
1118 the corresponding .list file. This process is repeated for each microengine intended for use. Then the code  
1119 is linked, where we explicitly specify which microengine will run which code file to get the final firmware  
1120 file. Lastly, we load the firmware onto the SmartNIC and start running the firmware.

**Listing 4.** Micro-C Code Compilation onto SmartNIC example

```

1121 # --- Building blm.list
1122 nfas -t -W3 -R -lm 0 -C -chip nfp-4xxx-b0 -DBLM_CUSTOM_CONFIG -
1123   DNFP_LIB_ANY_NFAS_VERSION -I. -I/c_packetprocessing/microc/blocks -
1124   I/c_packetprocessing/microc/include -I/c_packetprocessing/microc/
1125   lib -I/components/standardlibrary/include -I/components/
1126   standardlibrary/microcode/include -I/components/standardlibrary/
1127   microcode/src -DBLM_CUSTOM_CONFIG -DSINGLE_NBI -DPKT_NBI_OFFSET= -
1128   DBLM_BLK_EMEM_TYPE=emem -DNBII=8 -DBLM_INSTANCE_ID=0 -
1129   DBLM_INIT_EMU_RINGS -I. -I/dtree/mine/micro_c -I/c_packetprocessing
1130   /microc/blocks/blm/ -I/c_packetprocessing/microc/blocks/blm/_h -I/
1131   c_packetprocessing/microc/blocks/blm/_uc -o blm.list /
1132   c_packetprocessing/microc/blocks/blm/blm_main.uc
1133
1134 ### Island 32 ####
1135 # --- Building pkt_collect.list
1136 nfcc -W3 -chip nfp-4xxx-b0 -Qspill=7 -Qnn_mode=1 -Qno_decl_volatile -
1137   single_dram_signal -Qnctx_mode=8 -FI config.h -DBLM_CUSTOM_CONFIG
1138   -I. -I/dtree/mine/micro_c -I/c_packetprocessing/microc/include -I/
1139   c_packetprocessing/microc/lib -I/c_packetprocessing/microc/blocks/
1140   blm -I/c_packetprocessing/microc/blocks/blm/_h -I/components/
1141   standardlibrary/include -I/components/standardlibrary/microcode/
1142   include -Fepkt_collect_32.list /dtree/mine/micro_c/island32/
1143   pkt_collect_32.c /c_packetprocessing/microc/lib/nfp/libnfp.c /
1144   c_packetprocessing/microc/lib/std/libstd.c /c_packetprocessing/
1145   microc/lib/pkt/libpkt.c /c_packetprocessing/microc/lib/net/libnet.c
1146   /components/standardlibrary/microc/src/rtl.c
1147
1148 # --- Build feature1.list
1149 nfcc -W3 -chip nfp-4xxx-b0 -Qspill=7 -Qnn_mode=1 -Qno_decl_volatile -
1150   single_dram_signal -Qnctx_mode=8 -FI config.h -DBLM_CUSTOM_CONFIG
1151   -I. -I/dtree/mine/micro_c -I/c_packetprocessing/microc/include -I/
1152   c_packetprocessing/microc/lib -I/c_packetprocessing/microc/blocks/
1153   blm -I/c_packetprocessing/microc/blocks/blm/_h -I/components/
1154   standardlibrary/include -I/components/standardlibrary/microcode/
1155   include -Fefeature1_32.list /dtree/mine/micro_c/island32/
1156   feature1_32.c /c_packetprocessing/microc/lib/nfp/libnfp.c /
1157   c_packetprocessing/microc/lib/std/libstd.c /components/
1158   standardlibrary/microc/src/rtl.c
1159
1160 # --- Build feature2.list

```

```

1161 nfcc -W3 -chip nfp-4xxx-b0 -Qspill=7 -Qnn_mode=1 -Qno_decl_volatile -
1162   single_dram_signal -Qnctx_mode=8 -FI config.h -DBLM_CUSTOM_CONFIG
1163   -I. -I/dtree/mine/micro_c -I/c_packetprocessing/microc/include -I/
1164   c_packetprocessing/microc/lib -I/c_packetprocessing/microc/blocks/
1165   blm -I/c_packetprocessing/microc/blocks/blm/_h -I/components/
1166   standardlibrary/include -I/components/standardlibrary/microcode/
1167   include -Fefeature2_32.list /dtree/mine/micro_c/island32/
1168   feature2_32.c /c_packetprocessing/microc/lib/nfp/libnfp.c /
1169   c_packetprocessing/microc/lib/std/libstd.c /components/
1170   standardlibrary/microc/src/rtl.c
1171
1172 ... ### Repeat similar commands for remaining features ###
1173
1174 # --- Build result_collect.list
1175 nfcc -W3 -chip nfp-4xxx-b0 -Qspill=7 -Qnn_mode=1 -Qno_decl_volatile -
1176   single_dram_signal -Qnctx_mode=8 -FI config.h -DBLM_CUSTOM_CONFIG
1177   -I. -I/dtree/mine/micro_c -I/c_packetprocessing/microc/include -I/
1178   c_packetprocessing/microc/lib -I/c_packetprocessing/microc/blocks/
1179   blm -I/c_packetprocessing/microc/blocks/blm/_h -I/components/
1180   standardlibrary/include -I/components/standardlibrary/microcode/
1181   include -Feresult_collect_32.list /dtree/mine/micro_c/island32/
1182   result_collect_32.c /c_packetprocessing/microc/lib/nfp/libnfp.c /
1183   c_packetprocessing/microc/lib/std/libstd.c /components/
1184   standardlibrary/microc/src/rtl.c
1185
1186 ... ### Repeat similar commands for Islands 33-36 ###
1187
1188 # --- Link code
1189 nfld -chip nfp-4xxx-b0 -mip -rtsyms -o model.fw -map model.map -u i32
1190   .me0 -l pkt_collect_32.list -u i32.me1 -l feature1_32.list -u i32.
1191   me2 -l feature2_32.list ... -u ila0.me0 -l blm.list -i i8 -e /
1192   components/standardlibrary/picocode/nfp6000/catamaran/catamaran.
1193   npfw
1194
1195 ### Load Firmware onto SmartNIC ###
1196 nfp-nffw load --no-start model.fw
1197 nfp-nffw start

```