Investigating Resource-efficient Neutron/Gamma Classification ML Models Targeting eFPGAs

Jyothisraj Johnson a,1 Billy Boxer^b Tarun Prakash^a Carl Grace^a Peter Sorensen^a Mani Tripathi^b

^aLawrence Berkeley National Laboratory (LBNL), Berkeley, CA 94720-8099, USA

^b University of California, Davis, Department of Physics and Astronomy, Davis, CA 95616-5270, USA

E-mail: jyothisrajjohnson@lbl.gov

ABSTRACT: There has been considerable interest and resulting progress in implementing machine learning (ML) models in hardware over the last several years from the particle and nuclear physics communities. A big driver has been the release of the Python package, hls4ml, which has enabled porting models specified and trained using Python ML libraries to register transfer level (RTL) code. So far, the primary end targets have been commercial FPGAs or synthesized custom blocks on ASICs. However, recent developments in open-source embedded FPGA (eFPGA) frameworks now provide an alternate, more flexible pathway for implementing ML models in hardware. These customized eFPGA fabrics can be integrated as part of an overall chip design. In general, the decision between a fully custom, eFPGA, or commercial FPGA ML implementation will depend on the details of the end-use application. In this work, we explored the parameter space for eFPGA implementations of fully-connected neural network (fcNN) and boosted decision tree (BDT) models using the task of neutron/gamma classification with a specific focus on resource efficiency. We used data collected using an AmBe sealed source incident on Stilbene, which was optically coupled to an OnSemi J-series SiPM to generate training and test data for this study. We investigated relevant input features and the effects of bit-resolution and sampling rate as well as trade-offs in hyperparameters for both ML architectures while tracking total resource usage. The performance metric used to track model performance was the calculated neutron efficiency at a gamma leakage of 10^{-3} . The results of the study will be used to aid the specification of an eFPGA fabric, which will be integrated as part of a test chip.

KEYWORDS: Scintillators, Si-PMTs, Trigger concepts and systems

¹Corresponding author.

Contents

1	Introduction	1
2	The Python to RTL Workflow	2
3	Method	7
	3.1 Resource Efficient BDT Models	14
	3.2 Resource Efficient fcNN Models	16
4	Results	17
	4.1 BDT Results	17
	4.2 fcNN Results	20
5	Conclusion	22

1 Introduction

For a variety of radiation detection applications including neutron radiography [1] [2], associated particle imaging (API) [3] and neutron scatter cameras [4] [5] [6], there is a demonstrated need for compact front-end electronics that are easily scalable to perform individual read out of hundreds to thousands of silicon photomultipliers (SiPMs). Such electronics need to minimize power consumption in order to enable portable deployment. In addition, a requirement for the electronics in the majority of these target applications is an efficient method for discriminating neutron energy depositions from gamma depositions. Ideally, such functionality is implemented at the channel level. The added complexity for this task is that gamma dominated operating environments necessitate very low false positive rates while still maintaining high neutron efficiency.

This task of neutron/gamma classification is fundamentally multi-dimensional. Generally, specific plastic or organic crystalline scintillators with demonstrated inherent pulse shape discrimination (PSD) capability are used. This capability derives from the differences of a scintillator's singlet/triplet state activation with regards to the type of incident ionizing radiation [7]. However, other factors such as output scintillation wavelength and expected light yield also affect the overall discrimination power [8]. In addition, the time distribution of scintillation light produced is convolved with the SiPM single photon response to produce an output electronic current pulse, introducing yet another dimension. The classification task itself is typically performed by directly or indirectly using differences in the waveform tail between neutron and gamma energy depositions. The standard discriminator used is a ratio between partial and total integration of the waveform. However, several other methods have been utilized, mapping the discrimination problem into different parameter spaces covering both time and frequency domains. Examples include but are not

limited to [9] [10] [11] [12]. The majority of such methods do not have straightforward real-time compatible implementations.

For this work, we use Stilbene [13], a widely utilized organic scintillator, and the OnSemi J-series 60035 SiPM [14] to benchmark performance. While the SiPM single photon response is primarily set by details of the chosen SiPM, it is additionally shaped by any series resistance in the signal path and by the bandwidth of front-end electronics used [15]. Although the details of specific front-end topologies are outside the scope of this work, we nevertheless indirectly explore the effects of front-end bandwidth in the study.

Application specific integrated circuits (ASICs) provide a straightforward path to integrating individual read out of large numbers of channels, each with PSD capability, while minimizing power consumption. We have previously investigated one method of performing PSD using only analog circuits when targeting ASICs [16]. Briefly, the method utilizes the mapping of a traditional ratio of partial over total integrated charge to the difference between an event's scaled total integration value and corresponding partial integration value. A significant benefit results in that PSD capability increases at larger energy depositions rather than remaining flat as a function of energy past a certain energy deposition threshold. This is explicitly shown in Figure 8 of [16]. For an ASIC with a large number of integrated channels, the straight-forward physical circuit implementation requires multiple stages of circuitry at the channel level. In total, this would result in significant area utilization without spending considerable effort on details of circuit implementation. More realistically, a significant area reduction would likely require a mixed-signal solution with the subtraction implemented digitally. However, if already moving towards a mixed-signal solution, we can go beyond a simple translation of the required subtraction operation to a synthesized digital block on-chip. A novel strategy can be used to shift the classification task to a pipelined machine learning (ML) model implemented on an custom embedded field programmable gate array (eFGPA) fabric, which is then integrated on-chip.

This strategy has become viable in the last few years and is explained in more detail in Section 2. In Section 3, we provide a detailed overview of the study conducted in this work to explore the feasibility of an actual implementation, including setup and methodology. We introduce the test bed used for collecting data, provide an overview of the relevant hyperparameters for both types of ML architectures investigated in this study, and detail the metric used for evaluating model performance. In Section 4, we look at the results of the study individually for both ML model architectures. Finally, in Section 5, we summarize the results and its significance for an eventual eFPGA fabric implementation.

2 The Python to RTL Workflow

The implementation of neutron/gamma classification ML models on commercial FPGAs is wellestablished. Previous examples of work done include but are not limited to [17] [18] [19] [20] [21] [22]. Shifting from commercial FPGAs to eFPGA fabric targets is enabled by recent developments in both open-source embedded FPGA (eFPGA) frameworks [23] [24] and in porting machine learning (ML) models trained using Python ML libraries to register transfer level (RTL) code [25] [26]. The combined advances in both areas now provide the basis for a viable pathway to a real-time, lowpower and compact neutron/gamma classification method using ML models directly on a custom ASIC that integrates all required circuitry. Depending on the exact input features used, analog circuits or digital signal processing algorithms can be used to extract input feature values for a given event. The diagrams in Figure 1 present a generic system signal progression diagram comparing both scenarios. The exact input features required and a decision on if the ability to redefine the input features used is desired will guide such a decision. A re-programmable eFPGA fabric will need to be sufficiently specified to enable such an ability. In either case, the list of circuitry includes necessary front-end electronics for direct SiPM read out, ADCs for either waveform or individual input feature digitization and TDCs for capturing timing information in addition to dedicated digital blocks and/or an appropriately specified eFPGA fabric. A more detailed look at differences between using a waveform digitizer or analog circuits is given in Section 3 using a comparative look at the performance of the standalone charge comparison PSD method. Regardless, it is clear that the specifications of the ADC in both scenarios will directly set the best possible input feature value quantization. In other words, the ADC specification can play a significant role on achievable model performance. To account for this, we chose to fold in the ADC bit-resolution and sampling rate as parameters in our analysis to understand ADC requirements for a future test chip.





In this work, we focus on the first step towards a final integration of a customized open-source eFPGA fabric into an overall chip design for the task of neutron/gamma classification: establishing feasibility. Utilizing an open-source eFPGA framework allows for an FPGA fabric to be integrated

as part of an overall chip design without relying on expensive commercial eFPGA vendors. These frameworks allow for customization of the quantity and placement (location) of the various primitive resources (also known as tiles) that make up the fabric, chosen from those available in the open-source framework [23]. This allows for area-efficient solutions with the eFPGA fabric sized for a specific target application. Unlike a commercial FPGA, the availability of specific resource blocks, including memory resources, and their capabilities will depend on both the framework used and the technology node selected.

Achievable area-optimization for the eFPGA fabric will be limited by achievable routing densities and resource capabilities of the typically older process nodes used for the mixed-signal integrated chips targeting the applications under consideration. Explicitly, an eFPGA fabric implemented in 28nm versus 180nm will have very different capabilities for the same arrangement of primitive resource tiles defined using register transfer level (RTL) code. As a result, the exact latencies, resource usage, required maximum clock frequency, etc will be different for the same ML model depending on the process node selected for implementation. Due to these considerations, we used an older AMD/Xilinx Artix 7-series device (xc7a35ticsg324-1L) on an Arty development board as an initial stand-in for this exploration study. This commercial FPGA target has a total of 90 DSP blocks, 41.6k flip-flops, 20.8k LUTs and the equivalent of 100 18Kb RAM blocks. Compared to the top of the line for the 7-series FPGAs, this is an approximately 20x reduction in resources available. Compared against the top of the line UltraScale + family devices, this represents an approximately 60x reduction in available flip-flops and LUTs, 100x reduction for DSPs and 35x reduction for memory resources. The 7-series family devices are fabricated in a 28nm process node, representing a realistic upper bound comparison point on achievable performance and resource capabilities for the expected highest end node targeted for an actual custom eFPGA fabric.

The exact resource usage numbers and other specifications achieved for the stand-in commercial FPGA target will not directly translate for an eFPGA fabric target due to the open source frameworks considered. However, it is expected that the general trends for the models will hold and exact numbers will be off by some scaling factor. We aim to use the results of this study to inform the specification of resources needed for an actual eFPGA fabric to be implemented on a test chip.

In this work, we explored the design parameter space for two basic ML algorithms: boosted decision tree (BDT) and fully connected neural network (fcNN) models. The costs associated with fabricating chips scale proportionally with area for multi-project wafer (MPW) runs. As a result, we specifically investigated feasibility for resource-efficient implementations of these basic ML models, which will translate to an area-efficient implementation of an appropriately sized custom eFPGA fabric to run these models on a test chip.

The traditional FPGA implementation workflow uses RTL code, which specifies the intended behavioral functionality of a digital block, for design entry. From there, (vendor-specific) logic synthesis and place and route (PnR) software tools map the RTL code to resources such as look-up tables (LUTs), flip-flops (FFs), block random access memory (BRAM), etc on a commercial FPGA target die. In the last two decades, commercial high level synthesis (HLS) software has matured and been adopted to allow for bypassing the requirement to write RTL code directly. These tools use C/C++/SystemC/MATLAB code for design entry. The initial code provides an untimed, functional description of the intended digital block behavior and is then automatically mapped to a data flow graph (DFG), from which operations are scheduled based on capabilities of the target process node.

The final output of the HLS software is (timed) RTL code.

In order to generally implement ML algorithms in hardware, ML models defined using Python ML libraries need to be mapped to synthesizable C/C++/SystemC/MATLAB code. This last bridge was crossed in the past decade with the recent development of the **hls4ml** Python package. A detailed explanation of the **hls4ml** workflow is given in [25]. This package can directly interface with commercial HLS tools such as AMD/Xilinx Vivado HLS and Intel/Altera Quartus HLS to perform high level synthesis on its output C++ code, which functionally describes the ML model.

Originally, support for BDTs was integrated directly into **hls4ml** along with the support for NN-based ML architectures [26]. Since then, the BDT workflow has been moved into **conifer**, a complementary Python package released by the same collaboration. The packages, **hls4ml** and **conifer** respectively allow specification of quantization levels for the biases/weights or thresholds of the fcNN/BDT model and provides support for pipelining operations. Because of this, the ML model can accept new inputs after waiting a set number of clock cycles (initiation interval), potentially before the final prediction for the previous inference is available (achievable latency). This capability is crucial for replacing PSD circuitry of a large number of individual channels with a single pipelined ML model implemented on an eFPGA fabric. The **hls4ml** package is flexible enough that it provides several user-configurable settings, which can be set as needed to explore trade-offs in latency, throughput, power and resource usage for the classification task at hand. In comparison, **conifer** does not provide as extensive configurability, but the design parameter space for BDTs is also simpler [26].

Implementing ML models in hardware requires consideration of several new factors that are generally not necessary for a software implementation. The quantization of the input features plays a significant role in lowering overall resource usage as keeping native 32-bit resolution leads to larger resource usage for not much better performance. For our end-use target application(s), as mentioned, we anticipate having on-chip ADCs. The specifications of the ADC's bit resolution and full scale range will set the upper bound on the input feature quantization. The process of assigning ADC counts a physical value allows for selecting the base unit, which plays a non-trivial role in the required number of bits per input feature. For example, using millivolts (mV) or Volts (V) for storing the value of a pulse's peak amplitude is straightforward and interchangeable in software. However, both **hls4ml** and **conifer** use fixed point format to represent values for input feature. Thus, the number of integer and fractional bits explicitly need to be set. An additional bit will be used for the sign. Figure 2 explicitly demonstrates the difference in input word lengths required to exactly represent a value over a 2V full-scale range depending on if mV or V is the base unit. We expand on this point and the assumptions used in this study in Section 3.

To explore the full parameter space, we had to pick specific Python libraries to specify and train the ML models. While both **scikit-learn** and **XGBoost** have Python libraries allowing for the specification and training of BDTs, it is generally accepted that **XGBoost's** version results in both better performance and faster training. Furthermore, preliminary checks of resource usage for similarly set hyperparameter values between BDT models specified using the two libraries indicate very similar resource usage values. For these reasons, we used **XGBoost's** BDT model for this work. For fcNNs, several Python libraries are also available. However, at present, only **Keras** supports a complementary Python library for performing quantization aware training (**QKeras**), which is also supported by **hls4ml**. Taking into consideration the importance of this step in the overall aim of



Figure 2: The relationship between ADC bit resolution and the number of bits required in fixed point format to represent values in a base unit of V or mV over a 2V full-scale range. The top figure shows the mapping for mV and the bottom, for V.

resource-efficient ML eFPGA implementations, we have used (Q)Keras for specifying the fcNN models.

The overall requirements on neutron/gamma classification must also take into account the expected range of operating environments for the end-use applications. More specifically, the expected range of both gamma and neutron rates encountered will directly affect the absolute efficiency requirement on the PSD circuitry. We can tolerate larger gamma leakage values if the gamma rates will be less than or equal to neutron rates in the range of expected operating environments. However, in gamma dominated operating environments, much lower gamma leakage values are required. In this work, we assume operation in a gamma-dominated environment and



Figure 3: Diagram of the bench top test bed used to acquire the waveform dataset used in this work. Figure is replicated from [16].

define the performance metric used to assess performance of trained machine learning (ML) models as the neutron efficiency at a gamma leakage of 10^{-3} . This is not an absolute lower bound on expected gamma leakage requirements across all potential end-use applications. Rather, this is due to a statistical limitation from the size of the training data set used for the study as explained further in Section 3.

3 Method

The neutron and gamma waveform data set used in this work was acquired with the same benchtop test bed described in detail in [16]. We briefly summarize the test bed again here. Figure 3, replicated from [16], provides a system diagram of the setup. A custom SiPM array fanout board was used to read out selected channels of an 8x8 array of OnSemi J-series 60035 SiPMs. Each SOUT channel is connected in series with a 50 Ω resistor, which provides current-to-voltage conversion and appropriate impedance matching with the characteristic impedance of the coaxial connectors that the signal is routed to. Importantly, it also dominates the single photon response. The additional series resistance increases the decay constant associated with the recharging of the SiPM's single photon avalanche diode (SPAD) microcells after receiving a hit. The waveform tails extend over a time period > 1 μ s as a result. This fact should be kept in mind when interpreting and extrapolating any of the results of this study. Each FOUT signal, already a voltage output, is connected to a balun transformer that provides the correct interface to the coaxial connector that brings the signal off the board. Signals are digitized using a commercial 14-bit, 250 MS/s CAEN DT5725 digitizer with data collected using a computer interface. Specifically, the digitized waveforms were collected for a 6mm cube of Stilbene optically coupled to a single SiPM in the array with AmBe used as a mixed neutron/gamma emitter.

In total, 80k gamma and 49k neutron waveforms were selected from a full 30 million event data set with an energy cut selecting only events with > 90 keVee energy depositions. The goal of this study was to demonstrate feasibility of ML model deployments on an eFPGA, not necessarily



Figure 4: The neutron and gamma events used for this study are plotted as a function of PSD ratio vs. total deposition energy. A partial integration window of 128ns and total integration window of 1500ns is used. An energy cut at 90 keVee is included. Afterwards, upper and lower threshold cuts were defined to exclude events in between the neutron and gamma bands to maintain purity of the truth data set.

to improve classification efficiency at low energy depositions. The vast majority of events in the full data set correspond to the 60 keV Am-241 X-ray emission. Truth labels were assigned by performing a PSD ratio cut, with total integration set to 1500ns and partial integration set to 128ns, in the standard PSD ratio vs energy parameter space. Figure 4 shows these cuts and the resulting distribution of gamma and neutron events. The inversion of the neutron and gamma bands in the figure result from an inversion in the definition of PSD used. Instead of a tail-over-total, we performed a partial-over-total ratio as defined in Equation 3.1 with the start of the integration defined as t_0 , the partial integration ending at t_1 and the total at t_2 :

$$PSD = \frac{\int_{t_0}^{t_1} Waveform(t) dt}{\int_{t_0}^{t_2} Waveform(t) dt}.$$
(3.1)

Further details for this reasoning are found in [16]. The 130k total events in the dataset confines the gamma leakage metric to 10^{-3} . Going a magnitude lower means that the obtained neutron efficiency values will depend on the ability to classify only a few events from the full dataset, once split into training and truth/validation subsets. This would mean an added level of statistical uncertainty, which does not map to a true physical limit.

Although the original waveforms were digitized with 14-bit resolution over 2V full-scale and centered at 0V at a sampling rate of 250 MS/s, as mentioned, we wanted to understand the effect of ADC bit-resolution and sampling rate on the performance of trained BDT and fcNN models. To do so, we can perform re-digitization, digital filtering and decimation of the waveform samples in



Figure 5: The frequency response of the 2nd-order Butterworth filters used to re-filter the waveforms in software. Corner frequencies were specified corresponding to a 1x, 2x, 4x, 8x, and 16x downscaling of the actual 125MHz anti-aliasing filter used by the CAEN digitizer.

software in order to emulate a range of bit-resolutions, front-end bandwidths and effective sampling rates. In total, we considered bit-resolutions from eight to the original 14 bits and sampling rates between 15.625 MS/s (decimation by a factor of 16) to the original 250 MS/s. In addition, re-filtering of the waveforms without downsampling was also performed to take a closer look at the expected degradation in performance of the benchmark charge comparison PSD method. It is important to note that the measured thermal noise of the test bed was 180μ V, which limited effective resolution to around 12 bits.

The CAEN digitizer uses a second-order, 125 MHz anti-aliasing filter. We emulated the filter in software using a second-order Butterworth filter implemented with the **scipy** Python package. The critical frequency (corner frequency) of the filter was varied to be half the effective targeted sampling rate to match the physical digitizer's Nyquist configuration. Figure 5 shows a combined Bode plot of the various filter responses for the different targeted sampling rates. Five effective sampling rates were emulated by downsampling (every Nth sample kept) the original waveform samples by N, where N is the ratio of the original 250 MS/s sampling rate divided by the targeted sampling rate. In total, the parameter space encompassed by these two variables includes 35 unique pairs of effective bit-resolution and sampling rate. Figure 6 shows a representative FOUT and SOUT neutron waveform at the original 14 bit resolution and at a selection of other bit-resolutions down to 6 bits. No software re-filtering of the waveforms was performed for this figure.

In Section 2, a distinction was made on the degeneracy in physically implementing the extraction of input features. The two options discussed include using dedicated analog circuits and digitizing the final values or using a waveform digitizer and calculating input features digitally in a dedicated digital signal processing block. We can look at the performance of the standard charge comparison (PSD ratio) method under both scenarios. For both possibilities, 128ns was used for the partial



Figure 6: Example FOUT and SOUT waveforms re-digitized in software to various bit-resolutions without any applied software re-filtering. Each sample corresponds to the original 250 MS/s rate.

integration window and 1024ns was used for the total integration window. The actual number of samples corresponding to these time windows varies depending on downsampling for the second option.

For the analog circuit solution, the signal should be considered analog (having "infinite" resolution) until the final value is digitized. To emulate this signal chain, we re-filter the acquired waveforms from the dataset but do not downsample nor re-digitize them. This is reflected in Figure 7. The effective bandwidth (filter corner frequency) of the re-filtered waveforms should be interpreted as an effective front-end bandwidth. The bit-resolution, in comparison, still maps to the ADC. After partial and total integration values are calculated for these waveforms, the values are re-digitzed to the specified bit-resolution. In re-digitizing the final integration values, the nearest power of 2 that encompasses the numerical range of values is used to define the full-scale range. This same full-scale range is kept regardless of final targeted bit-resolution.

For the waveform digitizer option, we performed re-digitization as well as downsampling to match the corner frequency of the software filter to emulate the expected signal chain. Examples for the FOUT and SOUT waveforms after this procedure are shown in Figure 8. Partial and total integration values are calculated on these resulting waveforms. Because we can define extra bits to store the final integration value, no re-digitization of the final integration values is needed for



Figure 7: Example FOUT and SOUT waveforms re-filtered using a second order Butterworth filter implemented with the Python **scipy** package at various corner frequencies and 14 bit resolution. Downsampling of the waveforms is not performed for this figure. This is used in part to emulate the analog integration option for calculating input features.

emulating this second option.

We can plot the calculated partial and total integration values for both possible solutions. Figures 9, 10 plot the calculated partial and total integration values assuming an analog integration and digitization of final values, with the former figure showing an ideal 14-bit resolution on the final values and the latter assuming 10-bit resolution. The bottom panel for both figures show a zoomed out view of all events from the dataset used for the study and the top panel shows a zoomed in view of the lowest energy events above the 90 keVee equivalent cut. At 10-bit resolution, a significant degradation of the band separation is observed at the lowest energies compared to 14-bits, where band separation of the original waveforms is maintained across different effective bandwidths. It should be noted that a slight degradation at the lowest considered bandwidth is visible. On the other hand, Figures 11, 12 plot the results for the second option. We see that even at 14-bit resolution, band separation and widths worsen with decreasing effective bandwidth. At 8-bits resolution, a significant degradation of possible band separation is observed even at the highest bandwidth considered.



Figure 8: Example FOUT and SOUT waveforms re-filtered using a second order Butterworth filter implemented with the Python **scipy** package at various corner frequencies and 14 bit resolution. Downsampling of the waveforms is performed for this figure to equal the Nyquist rate equivalent for the software filter's corner frequency.

In general, it is clear that the standard charge comparison PSD metric is highly sensitive to the real-time implementation strategy and requires > 10 bit resolution and O(10) - O(100) MS/s for either possible implementation to maintain the ideal performance shown possible in [16]. In that paper, it was demonstrated that gamma leakages of < 10^{-5} can be achieved while maintaining neutron efficiencies of 99.9%. Maintaining that performance equates to the ability to maintain as close to the ideal band separation and widths at the lowest energy deposition events considered in this work.

As mentioned, for this ML model study, we are limited by dataset statistics to set gamma leakages at 10^{-3} but this should be treated as an upper bound as neutron efficiency can be maintained at lower leakage values. With a larger dataset, we can certainly set the performance metric at lower gamma leakage values. The results shown in Section 4 for both BDT and fcNN ML architectures should be considered from this point of view.

For the feasibility study, we only consider the waveform digitizer based approach to input feature calculation. For ML models with multiple input features, a waveform digitizer and calculation of



Figure 9: SOUT total vs. partial area is shown, both in units of $V^*\mu s$, for various effective bandwidths based on re-filtering, but not downsampling waveforms for various filter corner frequencies. The figure explicitly shows the results for an ideal 14-bit resolution for the final calculated input feature values. Degradation of neutron/gamma band widths is not observed across all filter corners but band separation slightly decreases as the filter corner value decreases. The gamma band is above the neutron band.

values for each feature digitally is, in general, less complex than designing analog circuits for each feature and offers a more flexible path towards re-definition, if required. BDTs and fcNNs have different model hyperparameters and, thus, considerations for resource efficiency. As a result, we independently constructed evaluation pipelines for these two ML architectures to account for such differences. These are discussed in the proceeding subsections. Regardless of the architecture, the same four input features were calculated for each waveform, for each of the 35 sampling rate, bitresolution pairs. The four input features are the SOUT total integration value (integration window of 1024ns), SOUT partial integration value (integration window of 128ns), FOUT peak amplitude, and SOUT peak amplitude. The exact integration windows used were selected to maintain compatibility with the various corner frequencies used for re-filtering of waveforms. The total integration window was changed from that used for labeling truth data in an attempt to decouple any potential biases in model training. Specifically, the partial integration window was selected to ensure that the window encompassed a sufficient area past the waveform peak amplitude to perform PSD even at the lowest filter corner value. The total integration window was set large enough to integrate on the full tail of the waveform. The input features were settled on after a dedicated importance study of various options. The ease for either analog or digital calculation of the features was kept in mind.

The final study parameter is the base unit used for mapping ADC bit resolution to fixed point format. While it seems that using V as the base unit is the obvious choice to minimize input word



Figure 10: SOUT total vs. partial area is shown, both in units of $V^*\mu s$, for various effective bandwidths based on re-filtering, but not downsampling waveforms for various filter corner frequencies. The figure explicitly shows the results for a non-ideal 10-bit resolution for the final calculated input feature values. Degradation of neutron/gamma band widths is observed across all filter corners, with band separation decreasing as filter corner value decreases. The gamma band is above the neutron band.

length, the range of the input feature values and differences in ranges between the different input features plays a significant role in performance of fcNN models. For the BDT models, differences in model performance was also observed depending on choice of input feature value scale. For this reason, the choice was made to use mV as the base unit for both ML model architectures despite the larger number of input bits required. For the area input features used, an area unit of $V^*\mu s$ was assumed so that range of values matched to an order of magnitude with the amplitude features. A more thorough investigation studying the effect of input feature value scale and range can be undertaken but is outside the scope of this paper. In addition, to provide optimal performance capability for the fcNN models, the input features were normalized to mean 0 and set to a standard deviation of 1. The scaler values to accomplish this can be extracted from the trained fcNN model to keep the same normalization for new input feature values. This will need to be implemented as a pre-processing step for an actual real-time execution of fcNN ML model inferences. Despite this, we kept the input word lengths the same across ML algorithms so that resource usage estimates can more closely correspond to differences from the two ML architectures.

3.1 Resource Efficient BDT Models

For BDTs, the **XGBoost** Python package provides a considerable amount of user optimization in finalizing a specific BDT model implementation. For this paper, we focused on the parameters that have the most effect on the final resource usage of the eFPGA and assumed fixed values for the



Figure 11: SOUT total vs. partial area is shown, both in units of $V^*\mu s$, for various effective bandwidths based on re-filtering and downsampling waveforms for various corner frequencies and their corresponding Nyquist equivalent sampling rates. The figure explicitly shows the results for an ideal 14-bit resolution for the waveform samples. Degradation of neutron/gamma band separation and widths are observed at lower filter corner frequencies. However the effects are not as drastic compared to 8 bit waveform sample resolution in Figure 12. The gamma band is above the neutron band.

other parameters. The main parameters that contribute to resource usage for the **XGBoost** BDT model are max_depth and num_rounds. The scaling of FPGA resources with these two parameters has been studied and is defined in [26] and replicated here:

$$\mathbf{r} = k_0 \cdot n_b + k_1 \cdot n_b \cdot 2^d \tag{3.2}$$

The terms k_0 and k_1 are unknown constants to be fitted for a given target, n_b is the number of boosting rounds (num_rounds) and d is the max tree depth (max_depth). Interested readers can refer to [26] for plots that visually show the scaling resulting from this equation. Although the **conifer** implementation of the BDT allows selection of the quantization level for the input feature values, we assumed this was primarily determined by the ADC bit-resolution and appropriate scaling of the feature values as discussed. Similarly, the quantization level of the thresholds and the output raw scores were also selected respectively to account for the range of threshold values in the trained model consistent with the ranges of input feature values and to maintain overall performance of the model. Importantly, these parameters do not improve the performance of the BDT model past the baseline already set by the other parameters.

In total, a four-dimensional parameter space to characterize performance on the neutron/gamma classification task remains: bit-resolution and sampling rate for the waveform digitizer (from which



Figure 12: SOUT total vs. partial area is shown, both in units of $V^*\mu s$, for various effective bandwidths based on re-filtering and downsampling waveforms for various corner frequencies and their corresponding Nyquist equivalent sampling rates. The figure explicitly shows the results for a non-ideal 8 bit waveform sample resolution. Significant degradation in neutron/gamma band widths are observed at lower filter corner frequencies with band separation degraded at all evaluated frequencies. The gamma band is above the neutron band.

the input features are calculated) and the max depth and number of boosting rounds for the BDT model. Due to the exponential increase of resources with max depth, a fixed value of 3 was assigned after a preliminary look at model performance for max depth = 2 showed that significantly larger values for number of boosting rounds (> 200) are required to obtain similar performance to setting max depth = 3. In other words, the decrease in resource usage from using a max depth value of 2 is overwhelmed by the linear increase in resource usage from the number of boosting rounds value. As a result, we only focused on a final three-dimensional parameter space for this study.

3.2 Resource Efficient fcNN Models

Much like **XGBoost**, **QKeras** for fcNNs has a large number of user configurable parameters. However, many more of the parameters have a non-trivial effect on the total resource usage of a trained model. In addition, unlike BDTs, there is no overall straightforward equation to describe the full parameter space's relationship to total resource usage. This is partly because the logic synthesis tool can optimize the exact physical implementation of the circuit. This can obfuscate resource scaling relationships due to degeneracies in implementing certain operations with DSP or LUT resources. Optimally, the decision is left primarily to the logic synthesis software unless specified as a directive/pragma by the user. We can obtain expected scaling of resources if iterating on a single parameter while keeping all others constant, as should reasonably be expected. However, to obtain a resource-efficient final model, a defined and ordered procedure for selecting values for the many contributing hyperparameters need to be defined. In total, there are nine considered hyperparameters that can affect the total resource usage: input feature bit resolution, output probability bit resolution, bit resolution for weights, biases (in theory, this can be specified per hidden layer), bit resolution of the activation function, total number of (and size of each) hidden layers, and sparsity value target for pruning aware training. A quantized tanh activation function is assumed for the hidden layer(s) and a sigmoid activation function for the output layer.

The order in which values are selected for these hyperparameters does matter. In general, the extent to which a model can be pruned/sparsified (i.e. the sparsity value target) without performance degredation partly depends on overall classification task complexity. Focusing on the neutron/gamma classification task at hand, it will depend strongly on the starting model size (number of hidden layers and size of each layer) as well as the input feature quantization level. So, we selected an ordered approach to defining values for the various hyperparameters.

An initial look was conducted by using the unmodified waveforms recorded with the test bed (14 bits and 250 MS/s). The purpose was to narrow down the range of values for the overall hyperparameter space by finding lower bounds. From these initial results, a single hidden layer of size 8 was sufficient to ensure nearly unity neutron efficiency even with bit resolution for the weights and biases set to 6 bits and for the activation function set to 10 bits. The number of epochs was set to 120 after the standard optimization against rate of loss function minimization. A sparsity target value of up to 50% was found to still maintain neutron efficiency above 90% in the initial look. These starting values for the hyperparameters provide anchor points when introducing input feature bit resolution and shaping. Instead of a large number of possible values for the nine parameters, we only accept a small upwards perturbation from these anchor values. For example, number of hidden layers was only allowed to vary between 1 and 2 and the size of the hidden layers could be either 8 or 16 for each hidden layer. Similarly, the number of bits for the weights and biases was allowed to be either 6 or 8, and for the activation function, either 10 or 12 if needed. Sparsity values were allowed to vary between 10% to 70% and the number of epochs for training was fixed at 120. Unlike the number of boosting rounds parameter for BDTs, the number of epochs does not directly increase the model size.

4 Results

The results are considered independently for the two considered ML architectures.

4.1 BDT Results

We calculated the neutron efficiency value at a gamma leakage of 10^{-3} for each of 35 pairs of bit-resolution and re-filtered corner frequency values (with downsampling) at several n_b (number of boosting rounds) values. At all bit-resolutions and corner frequencies, it is possible to achieve > 95% efficiency assuming an appropriate selection for n_b . Even with $n_b = 50$, we demonstrate that an ADC with only 8-bits resolution and Nyquist sampling rates of approximately 15 MS/s can still allow us to maintain excellent neutron efficiency at very low gamma leakage values in Figure 13. More generally, such a large parameter space for ADC specification where performance is maintained allows the final determination of design specs to be set by other system requirements for any specific target application.



Figure 13: The mean neutron efficiency and $\pm 1\sigma$ values are shown as a function of the n_b value for bit-resolutions between 8 to 11 bits for BDT models. Each panel respectively shows the results for various evaluated filter corner frequencies. The efficiencies for nbits = 9 to 11 are plotted offset from the x value for nbits = 8, for visual clarity in the figure.

One downside to selecting lower n_b values, however, is that the range of statistical variations possible for the final BDT model start to become significant. This can be true despite the possibility to still achieve very good neutron efficiency values. A possible origin to these variations is that the information content in the input features dataset is more than can be taken advantage of when the n_b value is too low. In this case, depending on initialization of thresholds during training, the model can achieve sub-optimal neutron efficiency values and, in general, display a high variance on final model performance.

To take this into account, a total of 20 trials were used as a compromise between required computation time and having sufficient sample size to calculate $\pm 1\sigma$ neutron efficiency values for each n_b value across filter corners. For each trial, we kept track of final neutron efficiency values achieved across a varying number of iterations of BDT model training and kept the model with the highest resulting value. The number of iterations varied between 200 for all n_b values less than 50 to 100 at $n_b = 80$. The results are shown in Figure 13. We observe significant variance on final neutron efficiency values are large compared to $n_b = 40$. This trend gets progressively worse at even lower values.

If we look individually at each bit-resolution performance at the various filter corner frequencies, we can observe some interesting differences across neutron efficiency results for $n_b = 10$ to 40. When we constrain the n_b value to be low, the information content (discrimination power) of the input features dataset can also be insufficient. We see this manifest as lower resulting neutron efficiency values for the smaller filter corner frequencies. The standard procedure for software implementations is to simply increase the value used. For the case of targeting a heavily resource constrained hardware implementation, to do so is not as straightforward, since resource usage also increases.

Another important point is that neutron efficiency values will plateau after reaching some critical n_b value and enter a diminishing returns regime. In this regime, for a hardware implementation, the resource requirements to implement the model increase but we only achieve marginally better performance. This behavior is easily observed in Figure 13. For this classification task, the plateau value is between 40 to 60, depending on filter corner frequency and bit resolution. In general, the value will also depend on both task complexity and the discrimination power of the selected/available input features.

Figure 14 histograms the number of flip-flops (nFFs) and number of look-up tables (nLUTs) for various bit resolution and corner frequency values as reported by Vivado post physical implementation. The selected n_b value for the figure is 50 but the trends seen can be generalized. This value was selected to convey the required resources and resulting performance for a model with a neutron efficiency that has just plateaued. Models with a significantly lower fraction of required resources are possible if the neutron efficiency target is lowered. It is important to note that the large statistical variations of the neutron efficiency values. These are primarily set by the hyperparameters that define the BDT model. However, depending on the details of the final trained model, the synthesizer can optimize out small parts of the model which can manifest as small fluctuations of these values.

For testing purposes, a streamlined UART-based module to read/write to/from and execute ML inferences was implemented as a wrapper to the model implementation so that input words can be written to the FPGA dev board and results of an inference execution can be read back to verify performance on an actual Artix 7-series FPGA. The final nFF and nLUT values in the figure include the resources needed for implementing this module. A target clock frequency of 100MHz was used as a constraint. The required number of flip-flops is less than 3% of the total available on the target device. The required number of LUTs across the parameter space is less than 25% of the total available.

While an initial increase in nFFs with bit resolution is observed, the numbers level off at higher bit resolutions. This is not the case for nLUTs. The required nLUTs increase with bit resolution and for each bit resolution, with filter corner frequency. No digital signal processor (DSP) nor block RAM (BRAM) resources are required for a BDT model implementation. The breakdown of the hardware implementation of BDT models and the expected mapping to resource usage is given in [26]. The minute deviations from the expected trends are most likely a result of optimizations made during logic synthesis and physical implementation by the Vivado synthesis tool. Across all bit resolution and sampling rate pairs, initiation intervals of 1 clock cycle were achieved. Similarly, latencies between 3 to 4 clock cycles of were achieved. The clock period for the study was fixed at 10ns but in general can be lower, which would result in potentially lower absolute time values for initial intervals and latencies.



Figure 14: A detailed look at the nFFs and nLUTs required for a BDT model with $n_b = 50$ across all evaluated bit resolution and filter corner frequency values. In addition, the returned latency and initiation interval values are given across all evaluated pairs.

4.2 fcNN Results

We also calculated the neutron efficiency values at a gamma leakage of 10^{-3} for each of the 35 pairs of bit-resolution and re-filtered corner frequency values (with downsampling) for trained fcNN models. Because of the larger parameter space, no exact comparable parameter to n_b for fcNNs is possible. A single hidden layer of size 16, with the number of bits used for the weights and biases equal to 8 and equal to 12 for the activation function, is sufficient to maintain high neutron efficiency values across all 35 considered pairs. This is true even at large sparsity target percentage values when performing pruning and quantization aware training. However, similar neutron efficiency values can also be obtained with a smaller hidden layer size and lower sparsity target values. We have used the larger size hidden layer model in Figure 15, which plots the observed statistical variance of fcNN models across a range of sparsity targets from 10% to 70%. As expected, with higher sparsity targets, the performance of the models at lower bit resolutions decreases, although overall neutron efficiency is still quite high. Significant statistical variance is not observed, unlike the BDT models. Some variance is seen at the largest sparsity values considered. Unlike the BDT models, the fcNN models were not able to maintain the same level of neutron efficiency performance at 8 bits across all filter corner values even at a minimal sparsity value target of 10%. Again, a total of 20 trials were used as a compromise between required computation time and having sufficient sample size to calculate $\pm 1\sigma$ values for each sparsity target value. Unlike for the BDT models, only one training iteration was performed per trial because the number of epochs is set sufficiently high with no significant increase in resource usages.

Figure 16 histograms the number of flip-flops (nFFs) and number of look-up tables (nLUTs) for various bit resolution and corner frequency values as reported by Vivado post physical implementation at a sparsity target value of 30% for the above considered model. Like with the BDT model results, this value was chosen as representative of the edge of the plateau region. However,



Figure 15: The mean neutron efficiency and $\pm 1\sigma$ values are shown as a function of the sparsity value target for bit-resolutions between 8 to 11 bits for fcNN models. Each panel respectively shows the results for various evaluated filter corner frequencies. The efficiencies for nbits = 9 to 11 are plotted offset from the x value for nbits = 8, for visual clarity in the figure. The fcNN models all have one hidden layer of size 16 with the number of bits used for the weights and biases equal to 8 and equal to 12 for the activation function.

in this case, the decrease in neutron efficiency is minimal. The resources required to implement a fcNN model can be significantly reduced if the neutron efficiency target value is lowered, as is the case for the BDT models. We also see small fluctuations in the figure as a result of the synthesizer possibly optimizing out small parts of the model.

The required number of flip-flops is less than 3% of the total available on the target device. The required number of LUTs across the parameter space is less than 10% of the total available. Unlike BDT models, BRAM and DSP resources are required to store the activation function in look-up table format and to perform the required multiplications when executing an inference. The number of DSPs are not explicitly plotted but they vary between two to 11 across the hyperparameter values defined for the figure. This is at maximum, 12% of the total available. Similarly, BRAM usage is reported as 0.5 of a single 18 Kb block for all models. This is the smallest unit reported and the true usage of the full block is expected to be different across the evaluated parameter space. An identical clock frequency constraint of 100MHz is specified. In addition, the same UART-based wrapper is used to read/write to/from and execute ML inferences to verify performance. Thus, resource usage numbers reported include this module as well and should be an identical constant number.

Unlike the BDT models, resource usage for nFFs and nLUTs are relatively constant across the hyperparameter space specified. The exact reason can be extrapolated to how a hardware



Figure 16: A detailed look at the nFFs and nLUTs required for fcNN models across all evaluated bit resolution and filter corner frequency values. The sparsity value target is set to 30% and all fcNN models have one hidden layer of size 16. The number of bits used for the weights and biases is equal to 8 and equal to 12 for the activation function. In addition, the returned latency and initiation interval values are given across all evaluated pairs.

implementation is approached for the two ML architectures. Regardless, resources will scale with sparsity value target much like they do with n_b for BDTs. The exact nFFs required for fcNN models are higher than for BDTs and the opposite is true for nLUTs. However, this needs to be considered against the additional DSPs and memory resource tiles required to implement fcNN models. The breakdown of the hardware implementation of fcNN models and general trends on resource usage is given in [27]. Across all bit resolution and sampling rate pairs, initiation intervals between 2 to 3 clock cycles were achieved. Similarly, latencies between 8 and 9 clock cycles were achieved. The BDT models clearly offer lower initiation intervals and latencies compared to fcNN models. This is expected from the constraints in gating the streaming of data between input to hidden and hidden to output layers for a given inference execution.

5 Conclusion

We have investigated the feasibility of deploying simple BDT and fcNN ML models on eFPGAs by using neutron/gamma classification as a case study and a stand-in commercial Artix 7-series FPGA target. We compared the study results to the golden standard of using the charge comparison method. The data set for the study was acquired with a fanout board where the choice of a 50 Ω resistor in series with the SiPM dominates the single photon response shape. In this regime, we were able to show that an ADC with only 8-bits resolution and Nyquist sampling rates of approximately 15 MS/s can still allow us to maintain excellent neutron efficiency at very low gamma leakage values when using either multiple input feature ML architecture. We compare this against the standalone standard charge comparison metric. With the waveform digitizer option, a >= 10-bit, O(100) MS/s

ADC has to be used to maintain comparable performance. With the analog integration option, a similar resolution ADC is required and sampling rate will be set by expected event rate requirements since only the final integration values are digitized.

More generally, we have shown that a low power, pipelined ML-aided solution is possible for real-time PSD functionality. An appropriately sized eFPGA fabric can be specified and synthesized to directly integrate on a single chip that includes direct SiPM readout front-ends and other required circuitry such as ADCs and TDCs. At 180nm, we expect the required area to be several mm by several mm to accommodate a properly sized eFPGA to implement the BDT models considered here. For fcNN models, the area required for most resource types will be smaller but additional area will be required for the necessary RAM blocks. If scaling down to 28nm, we expect the eFPGA fabric area requirements to decrease significantly.

Acknowledgments

This work was supported in part by the U.S. Department of Energy under Contract No. DE-AC02-05CH11231. In addition, support was provided for collaborators at the University of California, Davis by award DE-NA0003996.

References

- P. Hausladen, M.A. Blackston, E. Brubaker, D. Chichester, P. Marleau and R.J. Newby, *Fast-neutron coded-aperture imaging of special nuclear material configurations*, Tech. Rep. Idaho National Lab.(INL), Idaho Falls, ID (United States) (2012).
- [2] P. Hausladen, J. Newby, F. Liang and M. Blackston, The deployable fast-neutron coded-aperture imager: Demonstration of locating one or more sources in three dimensions, URL http://info. ornl. gov/sites/publications/files/Pub46191. pdf (2013).
- [3] M.R. Heath, B. Canion, L. Fabris, I. Garishvili, A. Glenn, J.U. Hausladen et al., *Development of a portable pixelated fast-neutron imaging panel*, *IEEE Transactions on Nuclear Science* 69 (2021) 1352.
- [4] J. Braverman, J. Brennan, E. Brubaker, B. Cabrera-Palmer, S. Czyz, P. Marleau et al., Single-volume neutron scatter camera for high-efficiency neutron imaging and spectroscopy, arXiv preprint arXiv:1802.05261 (2018).
- [5] J. Balajthy, J. Brown, E. Brubaker, B. Cabrera-Palmer, J. Cates, B. Goldblum et al., Characterization of a sipm-based monolithic neutron scatter camera using dark counts, arXiv preprint arXiv:2309.15280 (2023).
- [6] A. Galindo-Tellez, K. Keefe, E. Adamek, E. Brubaker, B. Crow, R. Dorrill et al., *Design and calibration of an optically segmented single volume scatter camera for neutron imaging, Journal of Instrumentation* 16 (2021) P04013.
- [7] F.D. Brooks, R.W. Pringle and B.L. Funt, Pulse shape discrimination in a plastic scintillator, IRE Transactions on Nuclear Science 7 (1960) 35.
- [8] M. Sénoville, F. Delaunay, M. Pârlog, N. Achouri and N. Orr, Neutron-γ discrimination with organic scintillators: Intrinsic pulse shape and light yield contributions, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 971 (2020) 164080.

- [9] X. Luo, Y. Wang, J. Yang, G. Liu, C. Lin, Q. Hu et al., Neutron/gamma discrimination employing the power spectrum analysis of the signal from the liquid scintillator bc501a, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 717 (2013) 44.
- [10] S. Marrone, D. Cano-Ott, N. Colonna, C. Domingo, F. Gramegna, E. Gonzalez et al., Pulse shape analysis of liquid scintillators for neutron studies, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 490 (2002) 299.
- [11] M. Nakhostin and P. Walker, Application of digital zero-crossing technique for neutron-gamma discrimination in liquid organic scintillation detectors, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 621 (2010) 498.
- [12] M. Aspinall, B. D'Mellow, R. Mackin, M. Joyce, N. Hawkes, D. Thomas et al., Verification of the digital discrimination of neutrons and γ rays using pulse gradient analysis by digital measurement of time of flight, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 583 (2007) 432.
- [13] "Inradoptics scintinel stilbene single crystals, datasheet." https://www.inradoptics.com/pdfs/ datasheets/InradOptics_Datasheet_Stilbene_Final.pdf.
- [14] OnSemi, Silicon Photomultipliers (SiPM), High PDE and Timing Resolution Sensors in a TSV Package; J-Series SiPM Sensors, 8, 2021.
- [15] F. Acerbi and S. Gundacker, Understanding and simulating sipms, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 926 (2019) 16.
- [16] B. Boxer, B. Godfrey, C. Grace, J. Johnson, R. Khandwala and M. Tripathi, *Studies in pulse shape discrimination for an optimized asic design, Journal of Instrumentation* 18 (2023) P01020.
- [17] C. Fu, A. Di Fulvio, S. Clarke, D. Wentzloff, S. Pozzi and H. Kim, *Artificial neural network algorithms for pulse shape discrimination and recovery of piled-up pulses in organic scintillators, Annals of nuclear energy* **120** (2018) 410.
- [18] N.M. Michels, A.J. Jinia, S.D. Clarke, H.-S. Kim, S.A. Pozzi and D.D. Wentzloff, *Real-time classification of radiation pulses with piled-up recovery using an fpga-based artificial neural network*, IEEE Access (2023).
- [19] M. Astrain, M. Ruiz, A.V. Stephen, R. Sarwar, A. Carpeño, S. Esquembri et al., *Real-time implementation of the neutron/gamma discrimination in an fpga-based daq mtca platform using a convolutional neural network, IEEE Transactions on Nuclear Science* 68 (2021) 2173.
- [20] A.D. Kaplan, B. Blair, C. Chen, A. Glenn, J. Ruz and R. Wurtz, A neutron-gamma pulse shape discrimination method based on pure and mixed sources, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 919 (2019) 36.
- [21] W. Zhang, W. Tongyu, B. Zheng, L. Shiping, Y. Zhang and Y. Zejie, A real-time neutron-gamma discriminator based on the support vector machine method for the time-of-flight neutron spectrometer, Plasma Science and Technology 20 (2018) 045601.
- [22] S. Yoon, C. Lee, B.-H. Won, S.-B. Hong, H. Seo and H.-D. Kim, Fast neutron-gamma discrimination in organic scintillators via convolution neural network, Journal of the Korean Physical Society 80 (2022) 427.

- [23] D. Koch, N. Dao, B. Healy, J. Yu and A. Attwood, *Fabulous: An embedded fpga framework*, in *The* 2021 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, pp. 45–56, 2021.
- [24] X. Tang, E. Giacomin, A. Alacchi, B. Chauviere and P.-E. Gaillardon, Openfpga: An opensource framework enabling rapid prototyping of customizable fpgas, in 2019 29th International Conference on Field Programmable Logic and Applications (FPL), pp. 367–374, IEEE, 2019.
- [25] F. Fahim, B. Hawks, C. Herwig, J. Hirschauer, S. Jindariani, N. Tran et al., *hls4ml: An open-source codesign workflow to empower scientific low-power machine learning devices, arXiv preprint arXiv:2103.05579* (2021).
- [26] S. Summers, G.D. Guglielmo, J. Duarte, P. Harris, D. Hoang, S. Jindariani et al., *Fast inference of boosted decision trees in fpgas for particle physics*, *Journal of Instrumentation* 15 (2020) P05026.
- [27] J. Duarte, S. Han, P. Harris, S. Jindariani, E. Kreinar, B. Kreis et al., Fast inference of deep neural networks in fpgas for particle physics, Journal of Instrumentation 13 (2018) P07027.