Automatic heterogeneous quantization of deep neural networks for low-latency inference on the edge for particle detectors

Claudionor N. Coelho Jr. Palo Alto Networks (California, USA)

Aki Kuusela, Shan Li, and Hao Zhuang Google LLC (California, USA)

Thea Aarrestad,* Vladimir Loncar,[†] Maurizio Pierini, Adrian Alan Pol, and Sioni Summers European Organization for Nuclear Research (CERN) (Geneva, Switzerland)

> Jennifer Ngadiuba California Institute of Technology (Caltech) (California, USA) (Dated: June 22, 2021)

Although the quest for more accurate solutions is pushing deep learning research towards larger and more complex algorithms, edge devices demand efficient inference and therefore reduction in model size, latency and energy consumption. One technique to limit model size is quantization, which implies using fewer bits to represent weights and biases. Such an approach usually results in a decline in performance. Here, we introduce a method for designing optimally heterogeneously quantized versions of deep neural network models for minimum-energy, high-accuracy, nanosecond inference and fully automated deployment on chip. With a per-layer, per-parameter type automatic quantization procedure, sampling from a wide range of quantizers, model energy consumption and size are minimized while high accuracy is maintained. This is crucial for the event selection procedure in proton–proton collisions at the CERN Large Hadron Collider, where resources are strictly limited and a latency of $\mathcal{O}(1) \mu$ s is required. Nanosecond inference and a resource consumption reduced by a factor of 50 when implemented on field-programmable gate array hardware are achieved.



FIG. I. An ultra-compressed deep neural network for particle identification on a Xilinx FPGA.

^{*} E-mail: thea.aarrestad@cern.ch

[†] Also at Institute of Physics Belgrade, Serbia.

I. INTRODUCTION

With edge computing, real-time inference of deep neural networks (DNNs) on custom hardware has become increasingly relevant. Smartphone companies are incorporating Artificial Intelligence (AI) chips in their design for ondevice inference to improve user experience and tighten data security, and the autonomous vehicle industry is turning to application-specific integrated circuits (ASICs) to keep the latency low. While the typical acceptable latency for real-time inference in applications like those above is $\mathcal{O}(1)$ ms [1, 2], other applications may require sub-microsecond inference. For instance, high-frequency trading machine learning (ML) algorithms are running on field-programmable gate arrays (FPGAs) to make decisions within nanoseconds [3]. At the extreme inference spectrum end of both the low-latency (as in high-frequency trading) and limited-area (as in smartphone applications) is the processing of data from proton-proton collisions at the Large Hadron Collider (LHC) at CERN [4]. In the particle detectors around the LHC ring, tens of terabytes of data per second are produced from collisions occurring every 25 ns. This extremely large data rate is reduced by a real-time event filter processing system – the trigger – which decides whether each discrete collision event should be kept for further analysis or be discarded. Data is buffered close to the detector while the processing occurs, with a maximum latency of $\mathcal{O}(1)$ μ s to make the trigger decision. High selection accuracy in the trigger is crucial in order to keep only the most interesting events while keeping the output bandwidth low, reducing the event rate from 40 MHz to 100 kHz. In 2027 the LHC will be upgraded from its current state, capable of producing up to one billion proton-proton collisions per second, to the so-called High Luminosity-LHC (HL-LHC) [5]. This will involve increasing the number of proton collisions occurring every second by a factor of five to seven, ultimately resulting in a total amount of accumulated data one order of magnitude higher than what is possible with the current collider. With this extreme increase, ML solutions are being explored as fast approximations of the algorithms currently in use to minimize the latency and maximize the precision of tasks that can be performed.

Hardware used for real-time inference in particle detectors usually has limited computational capacity due to size constraints. Incorporating resource-intensive models without a loss in performance poses a great challenge. In recent years many developments aimed at providing efficient inference from the algorithmic point of view. This includes compact network design [6–10], weight and filter pruning [11, 12] or quantization. In post-training quantization [13–17] the pre-trained model parameters are translated into lower precision equivalents. However, this process is, by definition, lossy and sacrifices model performance. Therefore, solutions to do quantization-aware training have been suggested [18–27]. In these, a fixed numerical representation is adopted for the whole model, and the model training is performed enforcing this constraint during weight optimization. More recently [28–31], it is argued that some layers may be more accommodating for aggressive quantization, whereas others may require more expensive arithmetic. This suggests that per-layer heterogeneous quantization is the optimal way to achieve higher accuracy at low resource cost, however, might require further specialization of the hardware resources.

In this paper, we introduce a novel workflow for finding the optimal heterogeneous quantization per layer and per parameter type for a given model, and deploying that model on FPGA hardware. Through minimal code changes, the model footprint is minimized while retaining high accuracy, and then translated into low-latency firmware. This paper makes the following contributions:

- We have implemented a range of quantization methods in a common library, which provide a broad base from which optimal quantizations can easily be sampled;
- We introduce a novel method for finding the optimal heterogeneous quantization for a given model, resulting in minimum area or minimum power DNNs while maintaining high accuracy;
- We have made these methods available online in easy-to-use libraries, called *QKeras* and *Au*to*QKeras*¹, where simple drop-in replacement of Keras [32] layers makes it straightforward for users to transform Keras models to their equivalent deep heterogeneously quantized versions, which are trained quantization aware. Using *AutoQKeras*, a user can trade-off accuracy by model size reduction (e.g. area or energy);
- We have added support for quantized QKeras models in the library, hls4ml [13], which converts these pre-trained quantized models into highly-parallel FPGA firmware for ultra low-latency inference.

To demonstrate the significant practical advantages of these tools for high-energy physics and other inference on the edge applications:

- We conduct an experiment consisting of classifying events in an extreme environment, namely the triggering of proton-proton collisions at the CERN LHC, where resources are limited and a maximum latency of $\mathcal{O}(1)\mu s$ is imposed;
- We show that inference within 60 ns and a reduction of the model resource consumption by a factor of 50 can be achieved through automatic heterogeneous quantization, while maintaining similar accuracy (within 3% of the floating point model accuracy);
- We show that the original floating point model accuracy can be maintained for homogeneously quantized DNNs down to a bit width of six while reducing resource consumption up to 75 % through quantization-aware training with QKeras.

The proposed pipeline provides a novel, automatic end-toend flow for deploying ultra low latency, low-area DNNs

¹ https://github.com/google/qkeras

on chip. This will be crucial for the deployment of ML models on FPGAs in particle detectors and other fields with extreme inference and low-power requirements.

The remainder of the paper is organized as follows. In Section II we discuss previous work related to model quantization and model compression with a focus on particle detectors. In Section V we uncover the novel library for training ultra low-latency optimally heterogeneously quantized DNNs, *QKeras*. Section VI describes the procedure of automatic quantization for optimizing model size and accuracy simultaneously. Finally, in Sections VII we deploy these optimally quantized QKeras models on FPGAs and evaluate their performance.

II. MOTIVATION

The hardware triggering system in a particle detector at the CERN LHC is one of the most extreme environments one can imagine deploying DNNs. Latency is restricted to $\mathcal{O}(1)\mu s$, governed by the frequency of particle collisions and the amount of on-detector buffers. The system consists of a limited amount of FPGA resources, all of which are located in underground caverns 50-100 meters below the ground surface, working on thousands of different tasks in parallel. Due to the high number of tasks being performed, limited cooling capabilities, limited space in the cavern, and the limited number of processors, algorithms must be kept as resource-economic as possible. In order to minimize the latency and maximize the precision of tasks that can be performed in the hardware trigger, ML solutions are being explored as fast approximations of the algorithms currently in use. To simplify the implementation of these, a general library for converting pre-trained ML models into FPGA or application-specific integrated circuits (ASIC) firmware has been developed, hls4ml [13]. The package comprises a library of optimized C++ code for common network layers, which can be synthesized through a high-level synthesis (HLS) tool. Converters are provided for multiple model formats, like TensorFlow [33], Keras [32], PyTorch [34] and ONNX [35].

Although other libraries for the translation of ML models to FPGA firmware exist, as summarized in Refs. [36– 39], hls4ml targets extreme low-latency inference in order to stay within the strict constraints of $\mathcal{O}(1)\mu s$ imposed by the hardware trigger systems. In addition, the unique aspect of hls4ml is the support for multiple HLS-vendor backends like Xilinx Vivado HLS, Intel Quartus HLS [40] and Mentor Catapult HLS [41], all of which are in use at the LHC experiments. The Vivado HLS backend is the most advanced and therefore the one used in this paper.

The hls4ml conversion process maps the user-provided neural network model into a given vendor-specific abstraction (like Vivado HLS), with easy-to-use handles to tune performance. The hls4ml NN architecture is introduced in [13]. A model-specific, layer-unrolled architecture is used in order to produce ultra low latency, resource efficient inference engines for particle physics. Computation for each NN layer is carried out in distinct hardware elements of the target device, which allows for high computational throughput through the laver pipeline, as well as fine-grained configuration of each layer (including quantization). A simple handle, named "Reuse Factor" enables users to control the parallelization of the computation, again at a per-layer level. In the fully parallel model, using a Reuse Factor of 1, each individual multiplication of the NN layers is carried out on different resources (whether FPGA DSPs or LUTs). With a Reuse Factor greater than 1, multiplication elements are reused sequentially to reduce the resource cost, at the expense of latency and throughput. This simple handle enables rapid design space exploration as well as configurability to target specific constraints in available resources, latency, and throughput.

In addition, the data access at the NN input and output, as well as the data movement between NN layers, can be configured to be fully parallel or fully serial. The former option is used to target ultra low latency, high throughput inference in the real-time processing of particle physics experiments, while the latter can be used to fit larger NN models within the available FPGA resources when ultra low latency is not as much of a constraint.

The hls4ml library is implemented as a Python package to facilitate ease of use for non-experts, as well as consistency with other popular Deep Learning libraries. The first step in the conversion into FPGA firmware consists of translating a given model into an internal representation of the network graph. During this conversion, user-specified optimization configurations are attached to the model, such as the choice of quantization and parallelisation. The internal representation is written out into an HLS project, assigning the appropriate layers of the target NN and the user configuration. This HLS project can then be synthesized with the FPGA vendor tools, generating an IP core that can be used in the target application. Many commonly used NN layers are supported: Dense; Convolution; BatchNormalization; and several Activation layers. In addition, domain specific layers can be easily added, one example being compressed distance-weighted graph networks [42].

In hls4ml, the precision used to represent weights, biases, activations, and other components are configurable through post-training quantization, replacing the floating point values by lower precision fixed-point ones. This allows compression of the model size, but to some extent sacrifices accuracy. Recently, support for binary and ternary precision DNNs [43] trained quantization-aware has been included in the library. This greatly reduces the model size, but requiring such an extremely lowprecision of each parameter type sacrifices accuracy and generalization.

As demonstrated in Refs. [28–31], mixed-precision quantization, i.e. keeping some layers at higher precision and some at lower precision, is a promising approach to achieve smaller models with high accuracy. However, finding the optimal heterogeneous quantization per layer and per parameter type, hereby referred to as quantization config*uration*, is extremely challenging, with the search space increasing exponentially with the number of layers in a model [30]. A solution for finding the mixed quantization configuration that yields best generalization/accuracy using the Hessian spectrum is proposed in Ref. [30]. For ML applications in hardware triggering systems, the resources one has at disposal, as well as the minimum tolerable model accuracy, are usually known. Finding the best model for a given task is, therefore, a fine compromise between the desired model compression and accuracy with respect to the floating point based model. Both factors must be considered when tuning quantization. The goal of this work is hence to provide a method for finding the optimal mixed-precision configuration for a given model, accounting for both the desired model size and accuracy when optimizing the architecture, and to transform these into highly parallel firmware for ultra low-latency inference on chip.

III. RELATED WORK

Closely related to the work presented here are the FINN [44] and FINN-R [45] frameworks from Xilinx Research Labs, which aim to deploy quantized neural networks on Xilinx FPGAs. The same group have also developed a library for quantization-aware training, BRE-VITAS [46], based on PyTorch model formats. The LogicNets design flow [47], also from Xilinx Research Labs, allows for the training of quantized DNNs that map to highly efficient Xilinx FPGA implementations. A comparison between the approach presented here and Logic-Nets is provided in Section VII. The FP-DNN [48] framework takes TensorFlow [33]-described DNNs as input and maps them onto FPGAs. The open-source alternative DNNWeaver [49] automatically generates accelerator Verilog code using optimized templates. Other frameworks focusing on the mapping of convolutional architectures onto efficient hardware design include Snowflake [50], fpgaConvNet [51–53] and Ref. [54]. For other work on FPGA DNN inference, we refer to the recent surveys at Refs. [36–39, 55]. TensorFlow Lite [56] is a set of tools for on-device inference with low latency and small binary sizes, targeting mobile, embedded and internet of things (IoT) devices. Currently, TensorFlow Lite supports deployment on Android and iOS devices, embedded Linux, and microcontrollers.

Our approach differs from those above with its emphasis on being a multi-backend tool, embracing a fully onchip design to target the microsecond latency imposed in physics experiments. The hls4ml library is completely open-source, and aims to provide domain scientists with easy-to-use software for deploying highly efficient ML algorithms on hardware.

In HAQ [28], a hardware-aware automated framework for quantization is introduced. The automization procedure consists of computing the curvature of the weight space of a layer, assuming a low curvature will require a lower bit-precision for the weights. Our approach differs from HAQ by combining reduced bit-precision with *filter* or neuron unit tuning, where the number of filters or neurons can be automatically tuned during the scan. In this case, the problem becomes highly non-linear, and we therefore take advantage of an AutoML-type of approach. A Bayesian optimization or randomized search is performed to find a solution that encompasses the precision used for the weights and activations, and the number of units or filters of the layer.

IV. PARTICLE IDENTIFICATION IN THE HARDWARE TRIGGER

A crucial task performed by the trigger system that could be greatly improved by a ML algorithm, both in terms of latency and accuracy, is the identification and classification of particles coming from each proton-proton collision. In this paper, we analyze the publicly available dataset introduced in Ref. [13, 57]. Here, a dataset [58] for the discrimination of *jets*, a collimated spray of particles, stemming from the decay and/or hadronization of five different particles was presented. It consists of quark (q), gluon (g), W boson, Z boson, and top (t) jets, each represented by 16 physics-motivated high-level features. In Ref. [13], this data set was used to train a DNN for deployment on a Xilinx FPGA. This model was compressed through post-training quantization in order to further reduce the FPGA resource consumption and provides a baseline to measure the benefits of quantization-aware training with heterogeneous quantization, over post-training quantization.

Adopting the same architecture as in Ref. [13], we use a fully-connected neural network consisting of three hidden layers (64, 32, and 32 nodes, respectively) with ReLU activation functions, shown in Fig. II. The output layer has five nodes, yielding a probability for each of the five classes through a Softmax activation function. The model definition in TensorFlow Keras is given in Listing 1.

As in [13], the weights of this network are homogeneously quantized post-training to a fixed pos yielding the best compromise between accuracy, latency, and resource consumption. This is found to be a fixed point precision, or *bit width*, of 14 bits with 6 integer bits. in the following referred to as $\langle 14, 6 \rangle$. We refer to this configuration as the baseline full model (BF). We then train a second pruned version of the BF model, hereby referred to as *baseline pruned* (BP). This model has 70%of its weights set to zero through an iterative process where small weights are removed using the TensorFlow Pruning API [59], following what was done in Ref. [13]. This reduces the model size and resource consumption significantly, as all zero-multiplications are excluded during the firmware implementation. We then create one heterogeneously quantized version of the BP model, where each layer is quantized independently post-training to yield



FIG. II. Model architecture for the fully-connected NN architecture under study. The numbers in brackets are the precisions used for each layer, quoted as $\langle B, I \rangle$, where B is the precision in bits and I the number of integer bits. When different precision is used for weights and biases, the quantization is listed as w and b, respectively. These have been obtained using the per-layer, per-parameter type automatic quantization procedure described in Section VI.

TABLE I. Per-layer quantization for the different baseline models (quantized post-training). When different precision is used for weights and biases, the quantization is listed as w and b, respectively.

Model	Precision							
	Dense	ReLU	Dense	ReLU	Dense	ReLU	Dense	$\operatorname{Softmax}$
BF/BP	$\langle 14, 6 \rangle$	$\langle 14, 6 \rangle$	$\langle 14, 6 \rangle$	$\langle 14, 6 \rangle$	$\langle 14, 6 \rangle$	$\langle 14, 6 \rangle$	$\langle 14, 6 \rangle$	$\langle 14, 6 \rangle$
\mathbf{BH}	w: $\langle 8, 3 \rangle$ b: $\langle 4, 2 \rangle$	$\langle 13, 7 \rangle$	$\langle 7, 2 \rangle$	$\langle 10, 5 \rangle$	$\langle 5, 2 \rangle$	$\langle 8, 4 \rangle$	w: $\langle 7, 3 \rangle$ b: $\langle 4, 1 \rangle$	$\langle 16, 6 \rangle$

Listing 1.	TensorFlow	Keras	model	definition.
TICOTIC I.	TOTOOLT TO W	TTOTOD	mout	CICILITICICICIT.

from tensorflow.keras.layers import Input
from tensorflow.keras.layers import Dense, Activation
from tensorflow.keras.layers import BatchNormalization
$\mathbf{x} = \text{Input}((16))$
$\mathbf{x} = \text{Dense}(64)(\mathbf{x})$
x = BatchNormalization()(x)
$\mathbf{x} = \operatorname{Activation}(\operatorname{'relu'})(\mathbf{x})$
$\mathbf{x} = \text{Dense}(32)(\mathbf{x})$
$\mathbf{x} = \text{BatchNormalization}()(\mathbf{x})$
$\mathbf{x} = \operatorname{Activation}(\operatorname{'relu'})(\mathbf{x})$
$\mathbf{x} = \text{Dense}(32)(\mathbf{x})$
x = BatchNormalization()(x)
x = Activation('relu')(x)
$\mathbf{x} = \text{Dense}(5)(\mathbf{x})$
x = Activation('softmax')(x)

the highest accuracy possible at the lowest resource cost. We start with an initial configuration of the model quantization using a wide bit-width, then iteratively reduce

the bit-width until reaching a threshold in accuracy loss relative to the initial floating-point model, evaluated on the training set. We iterate over the model in layer order, finding the appropriate precision for weights, biases, and output of a given layer, before moving to the next. We apply a more strict threshold in accuracy for earlier layers, since each round of precision reduction only degrades the accuracy. In this case we restrict to a 1% accuracy difference in the first layer, loosening to 2% for the final layer. This model is referred to as the *baseline heterogeneous* (BH) model. A summary of the per-layer quantizations for the baselines is given in Table I. From Ref. [13], we know that a post-training quantization of this model results in a degradation in model accuracy. The smaller the model footprint is made through post-training quantization, the larger the accuracy degradation becomes. To overcome this, we develop a novel library that, through minimal code changes, allows us to create deep heterogeneously quantized versions of Keras model, trained quantizationaware. In addition, as the amount of available resources on chip is known in advance, we want to find the optimal model for a given use-case allowing a trade-off between

model accuracy and resource consumption. We therefore design a method for performing automatic quantization, minimizing model area while maximizing accuracy simultaneously through a novel loss function. These solutions, simple heterogeneous quantization-aware training and automatic quantization are explained in the following sections.

V. QKERAS: OBTAINING OPTIMAL HETEROGENEOUS QUANTIZATION

Keras [32] is a high-level API designed for building and training deep learning models. It is used for fast prototyping, advanced research, and production. To simplify the procedure of quantizing Keras models, we introduce QKeras [60]: A quantization extension to Keras that provides a drop-in replacement for layers performing arithmetic operations. This allows for efficient training of quantized versions of Keras models.

QKeras is designed using Keras' design principle, i.e. being user-friendly, modular, extensible, and *minimally intrusive* to Keras native functionality. The code is based on the work of Refs. [18, 22], but provides a significant extension to these. This includes providing a richer set of layers (for instance including ternary and stochastic ternary quantization), extending the functionality by providing functions to aid the estimation of model area and energy consumption, allowing for simple conversion between non-quantized and quantized networks, and providing a method for performing automatic quantization. Importantly, the library is written in such a way that all the QKeras layers maintain a true drop-in replacement for Keras ones so that minimal code changes are necessary, greatly simplifying the quantization process. During quantization, QKeras uses the straight-through estimator (STE) [19], where the forward pass applies the quantization functions, but the backward pass assumes the quantization as the identity function to make the gradient differentiable.

For the model in Listing 1, creating a deep quantized version requires just a few code changes. An example conversion is shown in Listing. 2. The necessary code modifications consist of typing Q in front of the original Keras data manipulation layer name and specifying the quantization type, i.e. the kernel_quantizer and bias_quantizer parameters in a QDense layer. We change only data manipulation layers that perform some form of computation that may change the data input type and create variables (trainable or not). Data transport layers, namely layers performing some form of change of data ordering, without modifying the data itself, remain the same, e.g. Flatten. When quantizers are not specified, no quantization is applied to the layer and it behaves as the un-quantized Keras layer².

from tensorflow.keras.layers import Input, Activation from qkeras import quantized_bits from qkeras import QDense, QActivation from qkeras import QBatchNormalization

$\mathbf{x} = \text{Input}((16))$
$\mathbf{x} = \mathbf{QDense}(64,$
kernel_quantizer = $quantized_bits(6,0,alpha=1)$,
$bias_quantizer = quantized_bits(6,0,alpha=1))(x)$
$\mathbf{x} = \text{QBatchNormalization}()(\mathbf{x})$
$\mathbf{x} = \mathbf{Q}\mathbf{A}\mathbf{ctivation}(\mathbf{'quantized_relu}(6,0)')(\mathbf{x})$
$\mathbf{x} = \text{QDense}(32,$
kernel_quantizer = $quantized_bits(6,0,alpha=1)$,
$bias_quantizer = quantized_bits(6,0,alpha=1))(x)$
$\mathbf{x} = \text{QBatchNormalization}()(\mathbf{x})$
$\mathbf{x} = \mathbf{Q}$ Activation('quantized_relu(6,0)')(\mathbf{x})
$\mathbf{x} = \mathbf{QDense}(32,$
kernel_quantizer = $quantized_bits(6,0,alpha=1)$,
$bias_quantizer = quantized_bits(6,0,alpha=1))(x)$
$\mathbf{x} = \text{QBatchNormalization}()(\mathbf{x})$
$\mathbf{x} = \mathbf{Q}\mathbf{A}\mathbf{ctivation}(\mathbf{'quantized_relu}(6,0)')(\mathbf{x})$
$\mathbf{x} = \mathbf{QDense}(5,$
kernel_quantizer = $quantized_bits(6,0,alpha=1)$,
$bias_quantizer = quantized_bits(6,0,alpha=1))(x)$
$\mathbf{x} = \text{Activation}(\text{'softmax'})(\mathbf{x})$
· · · · · ·

The second code change is to pass appropriate quantizers, e.g. quantized_bits. In the example above, QKeras is instructed to quantize the kernel and bias to a bitwidth of 6 and 0 integer bits. The parameter alpha can be used to change the absolute scale of the weights while keeping them discretized within the chosen bit width. For example, in a binary network, rather than using the representations ± 1 , one can use \pm alpha. In QKeras, by setting alpha="auto", we also allow for the value of alpha to be computed during training from the absolute scale of the weights in question. Further details can be found in Methods XII B.

QKeras works by tagging all variables, weights and biases created by Keras as well as the output of arithmetic layers, by quantized functions. Quantized functions are specified directly as layer parameters and then passed to QActivation, which acts as a merged quantization and activation function. Quantizers and activation layers are treated interchangeably. To minimize code changes, the quantizers' parameters have carefully crafted and predefined defaults or are computed internally for optimal setup.

The quantized_bits quantizer used above performs

 $^{^2\,}$ The only exception is the <code>QBatchNormalization</code> layer. Here, when

Listing 2. Quantized QKeras model example.

no quantizers are specified, a power-of-2 quantizer is used for γ , σ and β , while μ remains unquantized. This has worked best when attempting to implement quantization efficiently in hardware and software (γ and σ become shift registers and β maintains the dynamic range aspect of the center parameter).

mantissa quantization:

$$2^{int-b+1}clip(round(x*2^{b-int-1}), -2^{b-1}, 2^{b-1}-1),$$
 (1)

where x is the input, b specifies the number of bits for the quantization, and **int** specifies how many bits of bits are to the left of the decimal point.

The quantizer used for the activation functions in Listing. 2, quantized_relu, is a quantized version of ReLU [61]. Two input parameters are passed, namely the precision, in this case 6 bits, and number of integer bits, in this case zero, respectively. The class has further attributes, for instance allowing for stochastic rounding of the activation function, all of which are described in detail in Ref. [60]. Fig. III shows the quantized ReLU function for three different bit widths and two different numbers of integer bits.



FIG. III. The quantized_relu function as implemented in QKeras for a 2- (purple), 3-bit (green and blue) and 6-bit (yellow) precision and for 0 or 1 integer bits. The unquantized ReLU function is shown for comparison (orange).

Through simple code changes like those above, a large variety of quantized models can be created. The full list of quantizers and layers is given in Methods XII A, or in the QKeras code repository [60].

We use QKeras to create a range of deep homogeneously quantized models, trained quantization-aware and based on the same architecture as the baseline model, which will provide a direct comparison between post-training quantization and models trained using QKeras. The model in Listing. 2 is an example of such a homogeneously quantized model. Finally, we want to create an optimally heterogeneously quantized QKeras model with a significantly reduced resource consumption, without compromising the model accuracy. The search space for finding such a configuration is large and exponential in layers. We therefore attempt to automatize the process by allowing users to scan through all the available quantizers in QKeras to find the configuration which fits the available chip area while maintaining high accuracy.

VI. AUTOQKERAS: RESOURCE-AWARE AUTOMATIC QUANTIZATION

As described in Section II, there are several methods for finding the optimal quantization configuration for a given model. These usually proceed by calculating the sensitivity of a given layer to quantization through evaluation of how small disturbances within that layer influence the loss function.

Often, as for instance in Refs [29, 30], only maximization of the model's accuracy and ability to generalize is considered. However, when doing inference on the edge, resources are often limited and shared between multiple applications. This is for instance the case in particle detectors, where a single FPGA is used to perform multiple different tasks. The desired accuracy and size constraints of the model in question are known in advance, and it is desirable to optimize the precision configuration considering both model accuracy and size. Some methods, like HAQ [28], does perform such a hardware-aware optimization. However, only the weight precision per layer is considered. When models are strongly quantized, it is often the case that more or fewer filters in convolutional layers, or neurons in densely connected layers, are necessary. A fine-tuning of the number of units per layer is therefore crucial to achieve the highest possible accuracy at the lowest resource cost.

In this paper, we introduce a method for performing automatic quantization where the user can trade off model area or energy consumption by accuracy in an applicationspecific way. The per-layer weight precision, as well as the number of neurons or filters per layer, are optimized simultaneously. By defining a *forgiving factor* based on the tolerated drop in accuracy for a given reduction in resource-cost, the best quantization configuration and number of units per layer, for a set of given size or energy constraints, can be found. We consider both energy minimization and bit-size minimization as a goal in the optimization.

A. Approximating relative model energy consumption

In order to target a reduction in model energy consumption, a high-level estimate of the model energy is needed. Here, we only concern ourselves with the difference in energy-consumption when comparing models using different quantizations, and not the absolute energy, as this is highly hardware specific. To this end, we assume an energy model where the energy consumption of a given layer is defined as

$$E_{layer} = E_{input} + E_{parameters} + E_{MAC} + E_{output}$$

These correspond to the energy cost of reading inputs E_{input} , parameters $E_{parameters}$ and output E_{output} , and the energy required to perform Multiply and Accumulate (MAC) operations E_{MAC} . For the first three, in a similar

way to compulsory accesses in cache analysis [62], we only consider the first access to the data, as only compulsory accesses are independent of the hardware architecture and memory hierarchy of an accelerator, when comparing models using the same architecture. We also assume a fully unrolled implementation on the hardware (as is the case with hls4ml). For the MAC energy estimation, we only consider the energy needed to compute the MAC. We do not include energy usage of registers, or glue and pipeline logic that will affect the overall energy profile of the device. For a given architecture this energy consumption is known, and here we assume a 45 nm process and follow the energy table given in Ref. [63].

Although this model provides a good initial estimate, it has high-variance concerning the actual energy consumption one finds in practice, especially for different architectural implementations. However, when comparing the energy of two different models, or models of different quantizations, both implemented in the same technology, this simple energy model is sufficient. The reason is that one can assume that the real energy of a layer is some linear combination of the high-level energy model, i.e. $E_{layer}^{Real} = k_1 \times E_{layer} + k_2$, where k_1 and k_2 are constants that depend on the architecture of the accelerator and in the implementation process technology. The slope can be considered as a factor accounting for the additional storage needed to keep the model running, and the offset corresponds to logic that is required to perform the operations. When comparing the energy consumption of two layers with different quantizations, L1 and L2, for the same model architecture, we have that $E_{L1}^{Real} > E_{L2}^{Real}$ if, and only if, the estimated energy $E_{L1} > E_{L2}$.

For these reasons, only relative energy estimates are considered during the automatic quantization, and users cannot target a specific energy value.

To facilitate easy estimation of the relative energy consumption or model bit-size when comparing different QKeras models, we have implemented a tool in the QKeras library, QTools, which performs both data type map generation and energy consumption estimation. A data type map for weights, biases, multipliers, etc., of each layer is generated, where the data type map includes operation types, variable sizes, quantizer types and bits. The output is an estimate of the per-layer energy consumption in pico-Joules, as well as a dictionary of data types per layer. Included in the energy calculation is a set of other tuneable specifications, like whether parameters and activations are stored on static random-access memory (SRAM) or dynamic random-access memory (DRAM), or whether data is loaded from DRAM to SRAM. The precision of the input can also be defined for a better energy estimate. The full list of options can be found in Ref. [60]. The *QTools* library provides an additional metric for model tuning when both accuracy and energy consumption, or model size, needs to be considered.

B. Defining a forgiving factor

With the high-level estimate of a given layers energy consumption provided by QTools, we define a forgiving factor to be targeted during automatic quantization of the model, providing a total loss function which combines energy cost and accuracy. The forgiving factor allows one to tolerate a degradation in a given metric, such as model accuracy, if the model gain in terms of some other metric, like model size, is significantly larger. Here, we allow the forgiving metric to be either minimization of the model bit-size or minimization of the model energy consumption. The forgiving factor is defined by

$$FF = 1 + \Delta_{acc} \times \log_{R}(S \times \frac{C_{ref}}{C_{trial}}),$$
 (2)

where $\Delta_{\rm acc}$ is the tolerated reduction in accuracy in percent, R is the factor stating how much smaller energy the optimized model must have compared to the original model (as a multiplicative factor to the FF metric) and S is a parameter to reduce the reference size, effectively forcing the tuner to choose smaller models. The parameters C_{ref} and C_{trial} refer to the cost (energy or bits) of the reference model and the quantization trial model being tested, respectively. The forgiving factor can be interpreted in the following way: If we have a linear tolerance for model accuracy degradation (or any other performance metric), we should be able to find a multiple of that degradation in terms of the cost reduction of the implementation. It enables an automatic quantization procedure to compensate for the loss in accuracy when comparing two models, by acting as a multiplicative factor.

tAutomatic quantization and re-balancing are then performed by treating quantization and re-balancing of an existing DNN as a hyper parameter search in Keras Tuner [64] using random search, hyperband [65] or Gaussian processes. We design an extension to Keras Tuner called AutoQKeras, which integrates the forgiving factor defined in Eq. 2 and the energy estimation provided by QTools. This allows for simultaneously tuning of the model quantization configuration and the model architecture. For instance, AutoQKeras allows for tuning of the number of filters in convolutional layers and the number of neurons in densely connected layers. This fine-tuning is critical, as when models are strongly quantized, more or fewer filters might be needed. Fewer filters might be necessary in cases where a set of filter coefficients get quantized to the same value.

Consider the example of quantizing two set of filter coefficient [-0.3, 0.2, 0.5, 0.15] and [-0.5, 0.4, 0.1, 0.65]. If we apply a binary quantizer with scale = $\left\lceil \log_2\left(\frac{\sum |w|}{N}\right) \right\rceil$, where w are the filter coefficients and N is the number of coefficients, we will end up with the same filter binary([-0.3, 0.2, 0.5, 0.15]) = binary([-0.5, 0.4, 0.1, 0.65]) = $[-1, 1, 1, 1] \times 0.5$. In this case, we are assuming a scale is a power-of-2 number so that it can be efficiently implemented as a shift operation. On the other hand, more filters might be needed as deep quantization drops information. To recover some of the boundary regions in layers that perform feature extraction, more filters might be needed when the layer is quantized. Lastly, certain layers are undesirable to quantize, often the last layer of a network. In principle, we do not know if by quantizing a layer we need more or less filters or neurons, and as a result, there are advantages to treating these problems as co-dependent problems, as we may be able to achieve a lower number of resources. Note that AutoQKeras does not completely remove model layers.

In AutoQKeras, one can specify which layers to quantize by specifying the index of the corresponding layer in Keras. If attempting to quantize the full model in a single shot, the search space becomes very large. In AutoQKeras, there are two methods to cope with this: grouping layers to use the same choice of quantization, or quantization by blocks. For the former, regular expressions can be provided to specify layer names that should be grouped to use the same quantization. In the latter case, blocks are quantized sequentially, either from inputs to outputs or by quantizing higher energy blocks first. If blocks are quantized one-by-one, assuming each block has N choices and the model consists of B blocks, one only needs to try $N \times B$, rather than N^B options. Although this is an approximation, it is a reasonable trade-off considering the explosion of the search space for individual filter selections, weight and activation quantization. Whether to quantize sequentially from inputs to outputs or starting from the block that has the highest energy impact, depends on the model. For example for a network like ResNet [66]. and if filter tuning is desirable, one needs to group the layers by the ResNet block definition and quantize the model sequentially to preserve the number of channels for the residual block. A few optimizations are performed automatically during model training. First, we dynamically reduce the learning rate for the blocks that have already been quantized so that they are still allowed to train, but at a slower pace. Also, we dynamically adjust



FIG. IV. ROC curves of false positive rate (FPR) versus true positive rate (TPR) for the Baseline Full (BF), quantized 6-bit (Q6), AutoQKeras Energy Optimized (QE) and AutoQKeras Bits Optimized (QB) models.

the learning rate for the layer we are trying to quantize as opposed to the learning rate of the unquantized layers. Finally, we transfer the weights of the model blocks we have already quantized whenever possible (when shapes remain the same).

We then use AutoQKeras to find the optimal quantization configurations for the baseline model for extremely resource-constrained situations, one targeting a minimization of the model's footprint in terms of model energy (QE) and one minimizing the footprint in terms of model bit-size (QB), using the different available targets in AutoQKeras. We want to reduce the resource footprint by at least a factor of 4 while allowing the accuracy to drop by at most 5%. We also allow for tuning of the number of neurons for each dense layer, for the same reason given above for model filter tuning. The model is quantized sequentially per block, where one block consists of a Dense layer and a ReLU layer. The resulting quantization configuration is listed in Table II. A very aggressive quantization configuration is obtained for both optimizations, with both binary and ternary quantizers and a bit-width of 4 at maximum for kernels. Despite the large search space, the obtained configurations are very similar as is to be expected due to the strong correlation between model energy and bit size. Whenever an input or the kernel has one (binary) or two (ternary) bits, we can completely eliminate multiplication operations in an implementation, saving valuable multiplier resources.

The preferred number of neurons per layer is half that of the original (32, 16, 16 rather than 64, 32, 32).

We then compare the relative energy consumption and bit size of the QE and QB models as computed with QTools, with respect to the simple homogeneously quantized model using a 6-bit precision in Listing 2, hereby referred to as Q6.

The QE and QB model energy consumption is reduced by 75% when compared to the Q6 model and, despite the aggressive quantization and reduction in neurons per layer, only a $\sim 3\%$ degradation in accuracy is observed for both . The total bit size is reduced by 80%. The QB model obtains a slightly smaller energy footprint than the QE model, alluding to some degree of randomness when scanning such a large search space. The relative power consumption when implemented on FPGA hardware will be discussed in Section VII.

All models presented above are trained minimizing the categorical crossentropy loss [67] using the Adam optimizer [68]. A learning rate of 0.0001 is set as the starting learning rate. If there is no improvement in the loss for ten epochs, the learning rate is reduced by 50% until a minimum learning rate of 10^{-6} is reached. The batch size is 1,024 and the training proceeds for 100 epochs. The training time for the models trained quantizationaware with QKeras is increased by ×1.5 with respect to the Keras equivalent.

Figure IV compares the classification performance of the BF, Q6, QE and QB models for two different target classes, top (t) and gluon (g). These classes were chosen



FIG. V. The full workflow starting from a baseline TensorFlow Keras Model, which is then converted into an optimally quantized equivalent through QKeras and AutoQKeras. This model is then translated into highly parallel firmware with hls4ml.

TABLE II. Per-layer quantization configuration and the relative model energy consumption for the AutoQKeras Energy Optimized (QE) and AutoQKeras Bits Optimized (QB) models, compared to the simple homogeneously quantized model, Q6.

Model Acc. [%]				Precision					$\frac{E}{E_{O6}}$	$\frac{\text{Bits}}{\text{Bits}_{O6}}$	
		Dense	ReLU	Dense	ReLU	Dense	ReLU	Dense	Softmax		-6
\mathbf{QE}	72.3	$\langle 4, 0 \rangle$	$\langle 4, 2 \rangle$	Ternary	$\langle 3, 1 \rangle$	$\langle 2, 1 \rangle$	$\langle 4, 2 \rangle$	w: Stoc. Bin. b: $\langle 8, 3 \rangle$	$\langle 16, 6 \rangle$	0.27	0.18
\mathbf{QB}	72.8	$\langle 4, 0 \rangle$	$\langle 4, 2 \rangle$	Stoc. Bin.	$\langle 4, 2 \rangle$	Ternary	$\langle 3, 1 \rangle$	Stoc. Bin.	$\langle 16, 6 \rangle$	0.25	0.17
$\mathbf{Q6}$	74.8	$\langle 6, 0 \rangle$	$\langle 6, 0 \rangle$	1.00	1.00						

as the ones where the original network, introduced in Ref. [13], had the highest and lowest AUC scores, respectively. Specifically, the receiver operating characteristic (ROC) curves of false positive rate (FPR) versus true positive rate (TPR), and the corresponding area under the curve (AUC) is shown. For particle detector trigger applications, it is often desirable to operate the algorithm at very low false positive rates, ensuring that only the most interesting events are kept while staying within the available trigger bandwidth. The classification performance of the Q6 model is almost identical to the BF model for FPRs down to 0.1%. The QE and QB models perform slightly worse, with AUC scores 0.02 points lower than for Q6 and BF. For a fixed FPR of 1%, the TPR for BF/Q6 is 60% and 55% for QE/QB. No significant degradation at very low FPR, where typical trigger algorithms would be operated, is observed.

With AutoQKeras, we give the user full flexibility to optimize the quantization configuration for a given use-case. An estimate of the model size and energy consumption can be computed using QTools and the user can then proceed by instructing AutoQKeras how much energy or bits it is desirable to save, given a certain accuracy-drop tolerance. Going from a pre-defined Keras model to an optimally quantized version (based on available resources) that is ready for chip implementation, is made extremely simple through these libraries.

The final, crucial step in this process is to take these quantized models and make it simple to deploy them in the trigger system FPGAs (or any hardware) while making sure the circuit layout is optimal for ultra lowlatency constraint. We will address this in the following section.

VII. ULTRA LOW-LATENCY, QUANTIZED MODEL ON FPGA HARDWARE

To achieve ultra low-latency inference of QKeras models on FPGA firmware, we introduce full integration of QKeras layers in the hls4ml library. The libraries together, provide a streamlined process for bringing quantized Keras models into particle detector triggering systems, while staying within the strict latency and resource constraints and performing high-accuracy inference.

When converting a QKeras model to an HLS project, the model quantization configuration is passed to hls4ml and enforced on the FPGA firmware. This ensures that the use of specific, arbitrary precision in the QKeras model is maintained during inference. For example, when using a quantizer with a given alpha parameter (i.e., scaled weights), hls4ml inserts an operation to re-scale the layer output. For binary and ternary weights and activations, the same strategies as in [43] are used. With binary layers, the arithmetical value of -1 is encoded as 0, allowing the product to be expressed as an XNOR operation. The full workflow starting from a baseline TensorFlow Keras model and up until FPGA firmware generation is shown in Fig. V. This illustrates how, through two simple steps, Keras models can be translated into ultra-compressed, highly parallel FPGA firmware.

We now compare the accuracy, latency and resource consumption of the different models derived above: The BF, BP, and BH models derived without using QKeras, two models optimized using AutoQKeras minimizing the model energy consumption, QE, and model bit consumption, QB, as well as a range of homogeneously quantized QKeras models scanning bit-widths from three to sixteen³. We compare the resource consumption and latency on chip for each model, to the model accuracy. The resources at disposal on the FPGA are digital signal processors (DSPs), lookup tables (LUTs), memory (BRAM) and flip-flops (FF). In this case the BRAM is only used as a LUT read-only memory for calculating the final Softmax function and is the same for all models, namely 1.5 units corresponding to a total of 54 Kb. For larger NNs using a higher reuse factor, and longer latency, BRAM may also be used to store model weights. The estimated resource consumption and latency from logic-synthesis, together with the model accuracy, are listed in Table III. A fully parallel implementation is used, with an *initiation interval*, the number of clock cycles between new data inputs, of 1 in all cases. Resource utilization is quoted in the percentage of total available resources, with absolute numbers quoted in parenthesis.

The most resource-efficient model is the AutoQKeras Energy Optimized (QE) model, reducing the DSP usage by ~ 98%, LUT usage by ~ 80%, and the FF usage by $\sim 90\%$. The accuracy drop is less than 3% despite using half the number of neurons per layer and overall lower precision. The extreme reduction of DSP utilization is especially interesting as, on the FPGA, DSPs are scarce and usually become the critical resource for ML applications. DSPs are used for all MAC operations, however, if the precision of the incoming numbers is much lower than the DSP precision (which, in this case, is 18 bits) MAC operations are moved to LUTs. This is an advantage, as a representative FPGA for the LHC trigger system has $\mathcal{O}(1000)$ DSPs compared to $\mathcal{O}(1)$ million LUTs. If the bulk of multiplication operations is moved to LUTs, this allows for deeper and more complex models to be implemented. In our case, the critical resource reduces from 56% of DSPs for the baseline to 3.4% of LUTs for the 6-bit QKeras trained model with the same accuracy. The latency is $\mathcal{O}(10)$ ns for all models.

In the final two columns of Table III, we compare the relative energy estimation from QTools with the post place-and-route power report from Vivado for the three QKeras models, in both cases relative to the Q6 model. Since the target clock frequency and model initiation interval is identical across these models, the inference rate is the same and taking the ratio of the power is equivalent 11

to taking the ratio of the energy. Very good agreement between the QTools relative energy estimates and the Vivado relative power estimates is observed for the QE and QB models, and the energy ordering is the same for all models.

Further, we compare the results obtained using the QKeras and hls4ml workflow to LogicNets [47]; another work on extreme low-latency, low-resource, fully-unfolded (II=1) FPGA implementations. The metrics are those quoted in Table III. Two LogicNets models have been evaluated: One using the same architecture as in this paper, JSC-M, and another using a larger architecture (32, 64, 192, 192, 16 number of neurons), JSC-L. For JSC-M, an accuracy of 70.6% is quoted, 1.7 points lower than the most resource-efficient model using QKeras and hls4ml, QE. In addition, QE uses $1.2 \times$ fewer LUTs than JSC-M. No DSPs are used in LogicNets, compared to the 66 DSPs in use by the QE model.

The latency has only been evaluated for JSC-L and is quoted to be 13 ns, using a clock frequency of 384 MHz. The final Softmax function has been removed from this estimate. In high energy physics experiments, the final Softmax layer is crucial since trigger thresholds usually are set based on an algorithms false positive rate. The threshold on the FPR is usually set as high as the trigger bandwidth allows, maximizing the true positive rate while straying within the bandwith-budget.

For a clock period of 5 ns, the QE model has a latency of 55 ns, reduced to 45 ns when ignoring the final Softmax layer. The JSC-L model has a latency of 13 ns for a clock period of 2.6 ns.

Finally, we compare the accuracy and resource consumption of a range of homogeneously quantized QKeras models, scanning bit widths from three to sixteen. In Fig. VI (left) the accuracy relative to the baseline model evaluated with floating point precision is shown as a function of bit width. This is shown for the accuracy as evaluated offline using TensorFlow QKeras (green line) and the accuracy as evaluated on the FPGA (orange line). We compare this to the performance achievable using the baseline model and post-training quantization (purple dashed line). The markers represent the accuracy of the baseline, baseline pruned, baseline heterogeneous and AutoQKeras optimized models (again emphasizing that the AutoQKeras models use half as many neurons per layer as the baseline Keras model). Models trained with QKeras retain performance very close to the baseline using as few as 6-bits for all weights, biases, and activations. Accuracy degrades slightly down to 98% of the baseline accuracy at 3-bits precision.

Post-training homogeneous quantization of the baseline model shows a much more significant accuracy loss, with accuracy rapidly falling away below 14-bits. The model resource utilization as a function of bit width for homogeneously quantized QKeras models is shown in the right plot in Fig. VI. The switch from DSPs to LUTs mentioned above is clearly visible: below a bit width of around 10, MAC operations are moved from the DSPs

³ Each model is trained using QKeras version 0.7.4, translated into firmware using hls4ml version 0.2.1, and then synthesized with Vivado HLS (2019.2), targeting a Xilinx Virtex Ultrascale 9+ FPGA with a clock frequency of 200 MHz.



FIG. VI. Relative accuracy (left) and resource utilization (right) as a function of bit width. The right-hand panel shows the metrics for the heterogeneously quantized models. The relative accuracy is evaluated with respect to the floating-point baseline model. Resources are expressed as a percentage of the targeted FPGA: Xilinx VU9P.

TABLE III. Model accuracy, latency, resource utilization and relative energy estimate for six different models. The latency is evaluated for a clock cycle of 200 MHz. Resources are listed as percentage of total, with absolute numbers quoted in parenthesis. The energy is estimated relative to the Q6 model and correspond to the relative energy computed using QTools (second to last column) and the relative power estimate from the post place-and-route report from Vivado (last column).

Model	Accuracy [%	[] Latency [ns]	Latency [clock cycles]	DSP $[\%]$	LUT [%]	FF [%]	$\frac{E_{QK}}{E_{QK}(\mathrm{Q6})}$	$\frac{P_{HLS}}{P_{HLS}(\rm Q6)}$
BF	74.4	45	9	56.0(1,826)	5.2(48,321)	0.8(20,132)	-	-
BP	74.8	70	14	7.7(526)	1.5(17,577)	0.4(10,548)	-	-
BH	73.2	70	14	1.3(88)	1.3(15,802)	0.3(8,108)	-	-
Q6	74.8	55	11	1.8(124)	3.4(39,782)	0.3(8,128)	1.00	1.00
\mathbf{QE}	72.3	55	11	1.0(66)	0.8(9,149)	0.1(1,781)	0.27	0.30
QB	71.9	70	14	1.0(69)	0.9(11,193)	0.1 (1,771)	0.25	0.25
LogicNets JSC-M [47]	70.6	N/A ^a	N/A	0 (0)	1.2(14,428)	0.02~(440)	-	-
LogicNets JSC-L [47]	71.8	13 ^b	5	0 (0)	3.2 (37,931)	0.03(810)	-	-

^a Not evaluated.

^b Using a clock frequency of 384 MHz.

to the LUTs and the critical resource consumption is significantly reduced. For instance, in this case, using a model quantized to 6-bit precision will maintain the same accuracy while reducing resource consumption by $\sim 70\%$. The markers in Fig. VI show the resource consumption of the heterogeneously quantized models. The only model comparable in accuracy and resource consumption to that of the AutoQKeras optimized models. QE and QB, is the baseline heterogeneous. However, in contrast to the QKeras models, BH has been pruned to a weight sparsity of 70% which further reduces the resource consumption (all zero multiplications are removed). In addition, the process of manually quantizing a model post-training is time-consuming and cumbersome, and not guaranteed to always succeed due to its lossy nature. AutoQKeras and hls4ml allows to quantize automatically through quantization-aware training, with specific tolerances in terms of accuracy and area, greatly simplifying the process.

In Ref. [69], the QKeras+hls4ml workflow has also been demonstrated on convolutional architectures benchmarked on the Streetview House Numbers dataset [70]. High accuracy matching the floating point model accuracy can be maintained down to 6-bit precision with QKeras, executed with 5 μ s latency. For larger convolutional architectures like ResNet [66], hls4ml doesn't scale due to the very low latency target. Our main application is the efficient implementation of tiny, custom models targeting $\mathcal{O}(10)$ ns - $\mathcal{O}(1)$ μ s latency.

VIII. CONCLUSION AND FUTURE WORK

We have introduced a novel library, QKeras, providing a simple method for uncovering optimally heterogeneously quantized DNNs for a set of given resource or accuracy constraints. Through simple replacement of Keras layers, models with heterogeneous per-layer, per-parameter type precision, chosen from a wide range of novel quantizers, can be defined and trained quantization-aware. A model optimization algorithm which considers both model area and accuracy is presented, allowing users to maximize the model performance given a set of resource constraints, crucial for high-performance inference on edge. Support for these quantized models has been implemented in hls4ml, providing the necessary chip layout instruction components to enable ultra-fast inference of these tiny-footprint models on a chip. We have demonstrated how on-chip resource consumption can be reduced by a factor of 50 without much loss in model accuracy while performing inference within $\mathcal{O}(10)$ ns. The methods presented here provide crucial tools for inference on the extreme low-area and low-latency edge, like that in particle detectors where a latency of $\mathcal{O}(1)\mu$ s is enforced. Taking a pre-trained model and making it suitable for hardware implementation on the edge, both in terms of latency and size, is one of the bottlenecks for bringing ML applications into extremely constrained computing environments (e.g. a detector at a particle collider), and the workflow presented here will allow for a streamlined and simple process, ultimately resulting in a great improvement in the quality of physics data collected in the future.

The generality and flexibility of the QKeras+hls4ml workflow opens up for a wide array of possible future work. This includes integration with other quantization libraries targeting non-FPGA hardware, like TensorFlow Lite, as well as those targeting FPGA synthesis, like FINN (and the quantization library Brevitas) and HAQ. In addition, while the energy estimator provides a good baseline for relative energy consumption, as demonstrated, we hope to extend the library to provide more device-specific absolute energy estimates. We also plan to explore using a combination of block energy and the curvature of the weight space, as done in HAQ, when quantizing a network one block at a time. Finally, work is ongoing to use the QKeras+hls4ml workflow to deploy ML algorithms for the next data taking period at CERN LHC both on FPGAs and ASICs.

IX. CODE AVAILABILITY

The QKeras library, which also includes AutoQKeras and QTools, is available under github.com/ google/qkeras, where the work presented here is using QKeras version 0.7.4. Examples on how to run the library can be found in the notebook The hls4ml library is available at subdirectory. github.com/fastmachinelearning/hls4ml and all ver-0.2.1 support QKeras models (the work sions \geq presented here is based on version 0.2.1). For examples on how to use QKeras models in hls4ml, the notebook part4_quantization at github.com/ fastmachinelearning/hls4ml-tutorial serves as a general introduction.

The data used in this study are openly available at Zenodo at Ref. [58] under DOI 10.5281/zenodo.3602260.

XI. AUTHOR INFORMATION

A. Corresponding author

Correspondence and material requests can be e-mailed to T. Aarrestad (thea.aarrestad@cern.ch).

B. Contributions

C. N. C., A. K., S. L., and H. Z. conceived and designed the QKeras, AutoQKeras and QTools software libraries. T. A., V. L, M. P., A. A. P., S. S. and J. N. designed and implemented support for QKeras in hls4ml. S. S. conducted the experiments. T. A., A. A. P. and S. S. wrote the paper.

ACKNOWLEDGMENT

M. P. and S. S. are supported by, V. L. and A. A. P. are partially supported by, the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program (grant agreement n^o 772369). V. L. is supported by Zenseact under the CERN Knowledge Transfer Group. A. A. P. is supported by CEVA under the CERN Knowledge Transfer Group. We acknowledge the Fast Machine Learning collective as an open community of multi-domain experts and collaborators. This community was important for the development of this project.

XII. METHODS

A. Additional layers, quantizers and methods in QKeras

In this section, we will give an overview of available layers, quantizers and methods in QKeras.

The summary of available layers in QKeras is listed in Table IV.

For several quantizers (including quantized_bits), a parameter called keep_negative can be set.

If keep_negative is true, negative numbers are not clipped. With a lower number of bits, the rounding adds more bias to the number system. Ref. [71] suggested using stochastic rounding, which uses the fractional part of the number as a probability to round the number up or down.

Stochastic rounding for quantized_bits quantizers can be turned on by setting use_stochastic_rounding =

Layers	Quantizers
QDense,	quantized_bits,
QConv1D,	binary,
QConv2D,	ternary,
QDepthwiseConv2D,	bernoulli,
QSeparableConv2D,	${\tt stochastic_ternary},$
QActivation,	<pre>stochastic_binary,</pre>
QAveragePooling2D,	<pre>smooth_sigmoid,</pre>
$\ttilde{QBatchNormalization},$	hard_sigmoid,
QOctaveConv2D,	binary_sigmoid,
QSimpleRNN,	$\texttt{smooth}_{-}\texttt{tanh},$
QLSTM,	hard_tanh,
QGRU	binary_tanh,
	quantized_relu,
	$quantized_ulaw,$
	$quantized_tanh,$
	quantized_po2,
	quantized_relu_po2

TABLE IV. List of available layers and quantizers in QKeras.

True. However, when an efficient hardware or software implementation is considered, this flag should be avoided in activation functions as it may affect the implementation efficiency.

Activations have been migrated to QActivation, but activation parameters passed directly in in convolutional and dense layers will be recognized as well.

The **bernoulli** and **stochastic** functions rely on stochastic versions of the activation functions, so they are best suited for weights and biases. They draw a random number with uniform distribution from sigmoid of the input x, adding additional regularization. The result is based on the expected value of the activation function. The **temperature** parameter determines the steepness of the sigmoid function.

The quantized guantized_relu and quantized_tanh are quantized versions of ReLU [61] and tanh functions.

The quantized_po2 and quantized_relu_po2 quantizers perform exponent quantization, as defined in [72]. The main advantage of this quantizer is that it provides a representation that is very efficient for multiplication. The parameter max_value defines maximum value.

It should also be noted that the QSeparableConv2D layer is implemented as a depthwise, followed by pointwise quantized expansions, which is an extended form of the SeparableConv2D implementation of MobileNet [73]. The reason we chose to use this version is that MobileNet's SeparableConv2D has an activation between the depthwise convolution and the pointwise convolution, where we need to at least apply some form of quantization.

Besides the drop-in replacement of Keras layers, we have written a few utility functions.

model_quantize function converts a non-quantized model into a quantized version, by applying a specified configuration for layers and activations. The method model_save_quantized_weights saves the quantized weights in the model compatible with an inference or writes the quantized weights in the file filename for production. The method load_qmodel loads and compiles quantized Keras model. Methods print_model_sparsity and print_qstats print sparsity for the pruned layers in the model and statistics of the number of operations per operation type and layer. quantized_model_debug allows for debugging and plotting model weights and activations. Finally, extract_model_operations estimates which operations are required for each layer of the quantized model, e.g. xor, mult, adder etc.

B. Variance shift handling in QKeras

A critical aspect when training quantized versions of tensors and trainable parameters, is the variance shift. During training with very few bits, the variance may shift a lot from its initialization. With popular initialization methods, e.g. glorot_normal, during the initial steps of the training, all of the output tensors will become zero. Consequently, the network will not be trained. For example, in a VGG network [74] the fully connected layers have 4096 elements, and any quantized representation with less than 6 bits will turn the output of these layers to be 0, as $\log_2(\sqrt{4096}) = 6$. For layer *i*, and minimum quantization threshold Δ , the weights w_i are quantized by quantizer (w_i) operation. When the gradient is computed, the quantized weights will appear as a result of the chain rule computation, as depicted in Fig. VII. With the absolute values of all weights below Δ , the gradient will vanish in all layers that transitively generates the inputs to layer *i*. This applies to any large DNN.



FIG. VII. Variance shift and the effect of initialization in gradient descent

QKeras mitigates this challenge by re-scaling the initialized weights appropriately. The parameter alpha is used as a scaling factor. It can be considered as a way to compute a shared exponent when used in weights [75]. It can be set to a given value manually, or overridden by setting it to auto or auto_po2. With alpha = "auto", we compute scale as $\sum q(x)x / \sum q(x)q(x)$ as in [24] for the quantization function q, with a different value for each output channel or output dimension of tensor x. This provides a learned scaling factor that can be used during training. With alpha = "auto_po2" [19], the scaling factor is set to be a power-of-2 number. For the ternary and stochastic_ternary quantizers, we iterate between scale computation and threshold computation, as presented in [76], which searches for threshold and scale tolerant to different input distributions. This is especially important when we need to consider that the threshold shifts depending on the input distribution, affecting the scale as well, as pointed out by [77]. When computing the scale in these quantizers with alpha = "auto", we compute the scale as a floating point number.

- Lin, S.-C. et al. The architectural implications of autonomous driving: Constraints and acceleration. SIG-PLAN Not. 53, 751-766 (2018). URL https://doi.org/ 10.1145/3296957.3173191.
- [2] Ignatov, A. et al. Ai benchmark: Running deep neural networks on android smartphones. In Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2018).
- [3] Leber, C., Geib, B. & Litz, H. High frequency trading acceleration using fpgas. In 2011 21st International Conference on Field Programmable Logic and Applications, 317–322 (2011).
- [4] The LHC Study Group. The Large Hadron Collider, Conceptual Design. Tech. Rep., CERN/AC/95-05 (LHC) Geneva (1995).
- [5] Apollinari, G., Béjar Alonso, I., Brüning, O., Lamont, M. & Rossi, L. High-luminosity large hadron collider (hl-lhc): Preliminary design report. Tech. Rep., Fermi National Accelerator Lab.(FNAL), Batavia, IL (United States) (2015).
- [6] Iandola, F. N. et al. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and; 0.5 mb model size. arXiv preprint arXiv:1602.07360 (2016).
- [7] Howard, A. G. et al. Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861 (2017).
- [8] Sandler, M., Howard, A., Zhu, M., Zhmoginov, A. & Chen, L.-C. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer* vision and pattern recognition, 4510–4520 (2018).
- [9] Ma, N., Zhang, X., Zheng, H.-T. & Sun, J. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In Proceedings of the European conference on computer vision (ECCV), 116–131 (2018).
- [10] Howard, A. et al. Searching for mobilenetv3. In Proceedings of the IEEE International Conference on Computer Vision, 1314–1324 (2019).
- [11] Ding, X. et al. Global sparse momentum sgd for pruning very deep neural networks. In Advances in Neural Information Processing Systems, 6382–6394 (2019).
- [12] He, Y., Zhang, X. & Sun, J. Channel pruning for accelerating very deep neural networks. In *Proceedings of* the *IEEE International Conference on Computer Vision*, 1389–1397 (2017).
- [13] Duarte, J. et al. Fast inference of deep neural networks in FPGAs for particle physics. JINST 13, P07027 (2018). 1804.06913.
- [14] Nagel, M., Baalen, M. v., Blankevoort, T. & Welling, M. Data-free quantization through weight equalization and bias correction. In *Proceedings of the IEEE International*

With alpha = "auto_po2", we enforce the scale to be a power of 2, meaning that an actual hardware or software implementation can be performed by just shifting the result of the convolution or dense layer to the right or left by checking the sign of the scale (positive shifts left, negative shifts right), and taking the log₂ of the scale. This behavior is compatible with shared exponent approaches, as it performs a shift adjustment to the channel.

Conference on Computer Vision, 1325–1334 (2019).

- [15] Meller, E., Finkelstein, A., Almog, U. & Grobman, M. Same, same but different-recovering neural network quantization error through weight factorization. arXiv preprint arXiv:1902.01917 (2019).
- [16] Zhao, R., Hu, Y., Dotzel, J., De Sa, C. & Zhang, Z. Improving neural network quantization without retraining using outlier channel splitting. arXiv preprint arXiv:1901.09504 (2019).
- [17] Banner, R., Nahshan, Y. & Soudry, D. Post training 4-bit quantization of convolutional networks for rapiddeployment. In Advances in Neural Information Processing Systems, 7950–7958 (2019).
- [18] Moons, B., Goetschalckx, K., Berckelaer, N. V. & Verhelst, M. Minimum energy quantized neural networks. arXiv preprint arXiv:1711.00215 abs/1711.00215 (2017). URL http://arxiv.org/abs/1711.00215. 1711.00215.
- [19] Courbariaux, M., Bengio, Y. & David, J.-P. Binaryconnect: Training deep neural networks with binary weights during propagations. In Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M. & Garnett, R. (eds.) Advances in Neural Information Processing Systems 28, 3123–3131 (Curran Associates, Inc., 2015).
- [20] Zhang, D., Yang, J., Ye, D. & Hua, G. Lq-nets: Learned quantization for highly accurate and compact deep neural networks. In *Proceedings of the European conference on computer vision (ECCV)*, 365–382 (2018).
- [21] Li, F. & Liu, B. Ternary weight networks. arXiv preprint arXiv:1605.04711 abs/1605.04711 (2016). URL http: //arxiv.org/abs/1605.04711. 1605.04711.
- [22] Zhou, S. et al. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. arXiv preprint arXiv:1606.06160 (2016).
- [23] Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R. & Bengio, Y. Quantized neural networks: Training neural networks with low precision weights and activations. *The Journal of Machine Learning Research* 18, 6869–6898 (2017).
- [24] Rastegari, M., Ordonez, V., Redmon, J. & Farhadi, A. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European conference on computer vision*, 525–542 (Springer, 2016).
- [25] Micikevicius, P. et al. Mixed precision training. arXiv preprint arXiv:1710.03740 (2017).
- [26] Zhuang, B., Shen, C., Tan, M., Liu, L. & Reid, I. Towards effective low-bitwidth convolutional neural networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018).
- [27] Wang, N., Choi, J., Brand, D., Chen, C.-Y. & Gopalakrishnan, K. Training deep neural networks with 8-bit float-

ing point numbers. In Advances in neural information processing systems, 7675–7684 (2018).

- [28] Wang, K., Liu, Z., Lin, Y., Lin, J. & Han, S. HAQ: hardware-aware automated quantization. arXiv preprint arXiv:1811.08886 abs/1811.08886 (2018). URL http: //arxiv.org/abs/1811.08886. 1811.08886.
- [29] Dong, Z., Yao, Z., Gholami, A., Mahoney, M. & Keutzer, K. Hawq: Hessian aware quantization of neural networks with mixed-precision. In 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 293–302 (2019).
- [30] Dong, Z. et al. Hawq-v2: Hessian aware trace-weighted quantization of neural networks (2019). 1911.03852.
- [31] Wu, B. et al. Mixed precision quantization of convnets via differentiable neural architecture search. arXiv preprint arXiv:1812.00090 (2018).
- [32] Chollet, F. et al. Keras (2015). URL https://github. com/fchollet/keras.
- [33] Abadi, M. et al. TensorFlow: Large-scale machine learning on heterogeneous systems (2015). URL http: //tensorflow.org/.
- [34] Paszke, A. et al. PyTorch: An imperative style, highperformance deep learning library. In Wallach, H. et al. (eds.) Advances in Neural Information Processing Systems 32, 8024 (Curran Associates, Inc., 2019). 1912.01703.
- [35] Open Neural Network Exchange Collaboration. ONNX. https://onnx.ai/ (2017). URL https://onnx.ai/.
- [36] Venieris, S. I., Kouris, A. & Bouganis, C.-S. Toolflows for mapping convolutional neural networks on FPGAs: A survey and future directions. ACM Comput. Surv. 51 (2018). 1803.05900.
- [37] Guo, K., Zeng, S., Yu, J., Wang, Y. & Yang, H. A survey of FPGA-based neural network inference accelerators. ACM Trans. Reconfigurable Technol. Syst. 12 (2019). 1712.08934.
- [38] Shawahna, A., Sait, S. M. & El-Maleh, A. FPGA-based accelerators of deep learning networks for learning and classification: A review. *IEEE Access* 7, 7823–7859 (2019). 1901.00121.
- [39] Abdelouahab, K., Pelcat, M., Serot, J. & Berry, F. Accelerating cnn inference on FPGAs: A survey (2018). 1806.01683.
- [40] Intel. Intel High Level Synthesis Compiler. https: //www.intel.com/content/www/us/en/software/ programmable/quartus-prime/hls-compiler.html (2020).
- [41] Mentor/Siemens. Catapult High-Level Synthesis. https://www.mentor.com/hls-lp/ catapult-high-level-synthesis (2020).
- [42] Iiyama, Y. et al. Distance-weighted graph neural networks on fpgas for real-time particle reconstruction in high energy physics. Frontiers in Big Data 3, 44 (2021). URL https://www.frontiersin.org/article/ 10.3389/fdata.2020.598927.
- [43] Ngadiuba, J. et al. Compressing deep neural networks on fpgas to binary and ternary precision with hls4ml (2020). 2003.06308.
- [44] Umuroglu, Y. et al. FINN: A framework for fast, scalable binarized neural network inference. In Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (ACM Press, 2017). 1612.07119.
- [45] Blott, M. et al. FINN-R: An end-to-end deep-learning framework for fast exploration of quantized neural networks. ACM Trans. Reconfigurable Technol. Syst. 11

(2018). 1809.04570.

- [46] Alessandro, Franco, G., nickfraser, Umuroglu, Y. & vfdev. Xilinx/brevitas: Release version 0.2.1 (2021). URL https: //doi.org/10.5281/zenodo.4507794.
- [47] Umuroglu, Y., Akhauri, Y., Fraser, N. J. & Blott, M. Logicnets: Co-designed neural networks and circuits for extreme-throughput applications (2020). 2004.03021.
- [48] Guan, Y. et al. FP-DNN: An automated framework for mapping deep neural networks onto FPGAs with RTL-HLS hybrid templates. In 2017 IEEE 25th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM), 152 (2017).
- [49] Sharma, H. et al. From high-level deep neural models to fpgas. In Microarchitecture (MICRO), 2016 49th Annual IEEE/ACM International Symposium on, 1 (IEEE, 2016).
- [50] Gokhale, V., Zaidy, A., Chang, A. X. M. & Culurciello, E. Snowflake: An efficient hardware accelerator for convolutional neural networks. In 2017 IEEE International Symposium on Circuits and Systems (ISCAS), 1–4 (2017).
- [51] Venieris, S. I. & Bouganis, C.-S. fpgaConvNet: A toolflow for mapping diverse convolutional neural networks on embedded FPGAs. In NIPS 2017 Workshop on Machine Learning on the Phone and other Consumer Devices (2017). 1711.08740.
- [52] Venieris, S. I. & Bouganis, C.-S. fpgaConvNet: Automated Mapping of Convolutional Neural Networks on FPGAs. In Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, 291 (ACM, 2017).
- [53] Venieris, S. I. & Bouganis, C.-S. fpgaConvNet: A Framework for Mapping Convolutional Neural Networks on FPGAs. In 2016 IEEE 24th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM), 40 (IEEE, 2016).
- [54] Huimin Li et al. A high performance fpga-based accelerator for large-scale convolutional neural networks. In 2016 26th International Conference on Field Programmable Logic and Applications (FPL), 1–9 (2016).
- [55] Zhao, R. et al. Hardware compilation of deep neural networks: An overview. In 2018 IEEE 29th International Conference on Application-specific Systems, Architectures and Processors (ASAP), 1–8 (2018).
- [56] Google. TensorFlow Lite (2020). URL https://www. tensorflow.org/lite.
- [57] Moreno, E. A. et al. JEDI-net: a jet identification algorithm based on interaction networks. Eur. Phys. J. C 80, 58. 15 p (2019). URL https://cds.cern.ch/record/2688535. 9 figures, 7 tables.
- [58] Pierini, M., Duarte, J. M., Tran, N. & Freytsis, M. Hls4ml lhc jet dataset (150 particles) (2020). URL https://doi. org/10.5281/zenodo.3602260.
- [59] Zhu, M. & Gupta, S. To prune, or not to prune: exploring the efficacy of pruning for model compression (2017). 1710.01878.
- [60] Coelho, C. Qkeras (2019). URL https://github.com/ google/qkeras.
- [61] Nair, V. & Hinton, G. E. Rectified linear units improve restricted boltzmann machines. In *ICML* (2010).
- [62] Hennessy, J. L. & et al. Computer architecture: A quantitative approach – sixth edition.
- [63] Horowitz, M. Computing's energy problem (and what we can do about it). In 2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC), 10-14 (2014).

- [64] O'Malley, T. et al. Keras Tuner. https://github.com/ keras-team/keras-tuner (2019).
- [65] Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A. & Talwalkar, A. Hyperband: A novel bandit-based approach to hyperparameter optimization. *The Journal of Machine Learning Research* 18, 6765–6816 (2017).
- [66] He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition (2015). 1512.03385.
- [67] Goodfellow, I., Bengio, Y. & Courville, A. Deep Learning (MIT Press, 2016). http://www.deeplearningbook.org.
- [68] Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations (ICLR) 2015, Conference Track Proceedings (2015). 1412.6980.
- [69] Aarrestad, T. *et al.* Fast convolutional neural networks on fpgas with hls4ml (2021). 2101.05108.
- [70] Netzer, Y. et al. Reading digits in natural images with unsupervised feature learning. In NIPS 2011 Workshop on Deep Learning and Unsupervised Feature Learning (2011). URL https://deeplearningworkshopnips2011. files.wordpress.com/2011/12/12.pdf.
- [71] Gupta, S., Agrawal, A., Gopalakrishnan, K. & Narayanan,P. Deep learning with limited numerical precision. In

International Conference on Machine Learning, 1737–1746 (2015).

- [72] Kwan, H. K. & Tang, C. Z. A design method for multilayer feedforward neural networks for simple hardware implementation. In 1993 IEEE International Symposium on Circuits and Systems, 2363–2366 vol.4 (1993).
- [73] Howard, A. G. et al. Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861 abs/1704.04861 (2017). URL http://arxiv.org/abs/1704.04861. 1704.04861.
- [74] Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image recognition. In *Interna*tional Conference on Learning Representations (2015).
- [75] Das, D. et al. Mixed precision training of convolutional neural networks using integer operations. arXiv preprint arXiv:1802.00930 abs/1802.00930 (2018). URL http: //arxiv.org/abs/1802.00930. 1802.00930.
- [76] Hwang, K. & Sung, W. Fixed-point feedforward deep neural network design using weights+ 1, 0, and- 1. In 2014 IEEE Workshop on Signal Processing Systems (SiPS), 1–6 (IEEE, 2014).
- [77] Li, F., Zhang, B. & Liu, B. Ternary weight networks. arXiv preprint arXiv:1605.04711 (2016).