





1st Mesoamerican Workshop on Reconfigurable X-ray Scientific Instrumentation for Cultural Heritage

Machine learning and model compression for FPGA/SoC

Antigua Guatemala, June 2025

Romina Soledad Molina, Ph.D.









The Abdus Salam International Centre for Theoretical Physics



Outline

- Introduction.
- Machine Learning.
- Machine Learning and System-On-Chips based on FPGA.
- ML and Model Compression Techniques.
- An End-to-End Workflow to Compress and Deploy DNN on FPGA.
 - DNN Training and Compression.
 - Integration with a Hardware Synthesis tool for ML.
 - Hardware Assessment Framework.
- ML and SoC-based FPGA for real-case applications.
- Final remarks.









A machine/system capable of imitating human behavior













Classification



Classification





Classification





Classification



















Generalization



Image from Togootogtokh, E., & Amartuvshin, A. (2018). Deep Learning Approach for Very Similar Objects Recognition Application on Chihuahua and Muffin Problem. ArXiv, abs/1801.09573.



Machine learning for classification





Machine learning for classification



- In a classifier, an input is mapped to a specific class.
- Supervised training phase: the network compares its current output with the desired output. The difference between these two values is corrected using backpropagation.











































Low latency



Low latency

Energy Efficiency



Low latency

Energy Efficiency

High parallelism



Low latency

Energy Efficiency

High parallelism

Scalability



Low latency

Energy Efficiency

High parallelism

Scalability

Customizable Al Acceleration











• Rise of real-time ML

- ML is now used in latency-critical domains: HEP, robotics, autonomous systems, IoT.
- The need for fast and efficient inference
 - Low-latency, low-power inference.
- Why FPGAs?
 - Highly parallel and reconfigurable, tuned for latency and energy efficiency.


Machine Learning and System-On-Chips based on FPGA

• Rise of real-time ML

- ML is now used in latency-critical domains: HEP, robotics, autonomous systems, IoT.
- The need for fast and efficient inference
 - Low-latency, low-power inference.

• Why FPGAs?

• Highly parallel and reconfigurable, tuned for latency and energy efficiency.

• The challenge

- ML models are software-native.
- Hardware mapping is complex and manual.
- Needs automation and abstraction.



- **Compression,** in the context of Machine Learning, involves techniques aiming to reduce the size of models or datasets while preserving performance.
 - Model Compression: Reducing the size of machine learning models (such as neural networks) without significantly affecting their accuracy.
 - Data compression: Reducing the size of the data used for training, validation, and testing. Example: Autoencoders.





• Trade-off Between Compression and Accuracy

The Abdus Salam International Centre for Theoretical Physics

> Balance between reducing the size of the model or data and maintaining its performance.





- Problem:
 - \circ Very large models \rightarrow high memory consumption and inference time.

The most accurate models are large and expensive in terms of memory and processing time.



The most accurate models (such as deep neural networks) are large and expensive in terms of memory and processing time.

Model	Parameters (millions)	Disk size (MB)
ResNet-50	~25.6M	~98 MB
BERT-base	~110M	~440 MB
BERT-large	~340M	~1.3 GB
VGG-16	~138M	~528 MB
VGG-19	~144M	~548 MB
YOLOv3	~62M	~236 MB
SpineNet-49S (small)	~11M	~45 MB
MobileNetV3-Large	~5.4M	~20 MB



- Problem:
 - \circ Very large models \rightarrow high memory consumption and inference time.

- Solution:
 - Compression



Distilled versions

Model	Parameters (millions)	Disk size (MB)	Accuracy
DistilBERT	66M (↓40%)	~250 MB	97% of BERT
SqueezeNet	1.2M (↓99%)	~4.8 MB (↓99%)	Similar to VGG-16
YOLOv8-Nano	3.2M (↓92%)	~8 MB (↓90%)	Good trade-off with YOLOv8-L
YOLO-Fastest	0.2M (↓99.5%)	~1 MB (↓99%)	Lower accuracy



















Pruning	Quantization	Knowledge distillation
Ļ	\downarrow	\downarrow
Remove neurons and connections.	Selection of the number of bits to represent the weights and bias.	Transfers the knowledge from a teacher network to a smaller and faster target network.
	Fully on-chip deployment	





Pruning













An end-to-end workflow to efficiently compress and deploy DNN on SoC/FPGA



End-to-end workflow

A- DNN training and compression



Publication:

Molina, R. S., Morales, I. R., Crespo, M. L., Costa, V. G., Carrato, S., & Ramponi, G. (2023). "An end-to-end workflow to efficiently compress and deploy DNN classifiers on SoC/FPGA". *IEEE Embedded Systems Letters*, *16*(3), 255-258.

Available at https://github.com/RomiSolMolina/workflowCompressionML



End-to-end workflow

A- DNN training and compression



Publication:

Molina, R. S., Morales, I. R., Crespo, M. L., Costa, V. G., Carrato, S., & Ramponi, G. (2023). "An end-to-end workflow to efficiently compress and deploy DNN classifiers on SoC/FPGA". *IEEE Embedded Systems Letters*, *16*(3), 255-258.

Available at https://github.com/RomiSolMolina/workflowCompressionML



End-to-end workflow

A- DNN training and compression



Publication:

Molina, R. S., Morales, I. R., Crespo, M. L., Costa, V. G., Carrato, S., & Ramponi, G. (2023). "An end-to-end workflow to efficiently compress and deploy DNN classifiers on SoC/FPGA". *IEEE Embedded Systems Letters*, *16*(3), 255-258.



A. DNN training and compression



DNN training and compression Stage 1 - Teacher training





DNN training and compression Stage 2 - Student training











🐼 ONNX





















https://github.com/fastmachinelearning/











Integration with a hardware synthesis tool for ML High-level synthesis

- High-Level Synthesis
 - It provides the facility to create RTL from a high level of abstraction.
 - Several implementations are possible from the same source description.
 - Implements the design based on defaults and user applied directives.
 - It allows the optimization of the input code using directives to:
 - Reduce latency.
 - Improve performance and throughput.
 - Reduce resource utilization.

Without optimization, HLS tool will look to minimize latency and improve concurrency.



HLS Component Development Flow



https://docs.amd.com/r/en-US/ug1399-vitis-hls/Introduction-to-Vitis-HLS-Components



HLS Component Development Flow



https://docs.amd.com/r/en-US/ug1399-vitis-hls/Introduction-to-Vitis-HLS-Components


HLS Component Development Flow



https://docs.amd.com/r/en-US/ug1399-vitis-hls/Introduction-to-Vitis-HLS-Components



HLS Component Development Flow



https://docs.amd.com/r/en-US/ug1399-vitis-hls/Introduction-to-Vitis-HLS-Components



HLS Component Development Flow



https://docs.amd.com/r/en-US/ug1399-vitis-hls/Introduction-to-Vitis-HLS-Components





ML framework support:

• (Q)Keras

- PyTorch
- (Q)ONNX

[hls4ml] https://fastmachinelearning.org/hls4ml/





ML framework support:

• (Q)Keras

- PyTorch
- (Q)ONNX

Neural networks architectures:

- Fully Connected NN
- Convolutional NN
- Recurrent NN
- Graph NN

[hls4ml] https://fastmachinelearning.org/hls4ml/





ML framework support:

• (Q)Keras

- PyTorch
- (Q)ONNX

Neural networks architectures:

- Fully Connected NN
- Convolutional NN
- Recurrent NN
- Graph NN

HLS backends:

- Vivado HLS
- Intel HLS
- Vitis HLS
- Catapult HLS
- oneAPI (experimental)



High-Level Synthesis for Machine Learning



hls4ml is tested on the following platforms:

Vivado HLS versions 2018.2 to 2020.1 Intel HLS versions 20.1 to 21.4. Versions > 21.4 have not been tested. Vitis HLS versions 2022.2 to 2024.1. Versions <= 2022.1 are known not to work. Catapult HLS versions 2024.1_1 to 2024.2 oneAPI versions 2024.1 to 2025.0





How does it work?

With hls4ml, each layer of output values is calculated independently in sequence, using pipelining to speed up the process by accepting new inputs after an initiation interval. The activations, if nontrivial, are precomputed.







How does it work?

With hls4ml, each layer of output values is calculated independently in sequence, using pipelining to speed up the process by accepting new inputs after an initiation interval. The activations, if nontrivial, are precomputed.

Simplifying the input network must be done before using hls4ml to generate HLS code, for optimal compression to provide a sizable speedup.







Trade-off between latency and FPGA resource usage determined by the parallelization of the calculations in each layer.





Trade-off between latency and FPGA resource usage determined by the parallelization of the calculations in each layer.

Reuse factor: number of times a multiplier is used to do a computation.







Trade-off between latency and FPGA resource usage determined by the parallelization of the calculations in each layer.

Reuse factor: number of times a multiplier is used to do a computation.



Fewer resources, Lower throughput, Higher latency

More resources, Higher throughput, Lower latency

[hls4ml] https://fastmachinelearning.org/hls4ml/









io_parallel

io_stream

[hls4ml] https://fastmachinelearning.org/hls4ml/





io_parallel

Data is passed in **parallel** between the layers.

This style allows for **maximum parallelism** and is well suited for MLP networks and small CNNs which aim for lowest latency.





io_stream

Data is passed **one "pixel" at a time**.

Each pixel is an **array of channels**, which are always sent in parallel. This method for sending data between layers is recommended for larger CNN and RNN networks.





Strategy: the implementation of core matrix-vector multiplication routine, which can be **latency-oriented**, **resource-saving oriented**, **or specialized**.

Different strategies will have an **impact on overall latency and resource consumption** of each layer, and users are advised to choose based on their design goals.





Activations - Implementation parameter

- **latency:** Good latency, high resource usage. **It does not work well if there are many output classes**.
- **stable:** Slower but with better accuracy, useful in scenarios where **higher accuracy is needed**.
- **legacy:** An older implementation with poor accuracy, but good performance. Usually the latency implementation is preferred.
- **argmax:** If you don't care about normalized outputs and only care about which one has the highest value, using argmax saves a lot of resources. This sets the highest value to 1, the others to 0.





Profiling

- The tools in **hls4ml.model.profiling** can help to choose the right precision for the model.
- hls4ml.model.profiling.numerical method with three objects: a Keras model object, test data, and an HLSModel.





Python integration

hls4ml°

1 from tensorflow.keras.models import load_model 2 from sklearn.metrics import accuracy_score 3 model = load_model('model_keras_MLP.h5')

4 model.summary()

Model: "sequential"

Layer (type)	Output	Shape	Param #
fcl (Dense)	(None,	60)	3900
relu1 (Activation)	(None,	60)	0
fc0 (Dense)	(None,	40)	2440
relu0 (Activation)	(None,	40)	0
fc2 (Dense)	(None,	30)	1230
relu2 (Activation)	(None,	30)	0





Python integration





Python integration



```
import hls4ml
import plotting
hls4ml.model.optimizer.OutputRoundingSaturationMode.layers = ['Activation']
hls4ml.model.optimizer.OutputRoundingSaturationMode.rounding mode = 'AP RND'
hls4ml.model.optimizer.OutputRoundingSaturationMode.saturation mode = 'AP SAT'
config = hls4ml.utils.config from keras model(model, granularity='name')
config['Model'] = {'Precision' : 'ap fixed<17,16>', 'ReuseFactor' : 1, 'Strategy' : 'Latency'}
config['LayerName']['fc1']['Precision']['weight'] = 'ap fixed<9, 1>'
config['LayerName']['softmax']['Precision'] = 'ap fixed<32,15>'
print("-----
hls model = hls4ml.converters.convert from keras model(model.
                                                    hls config=config.
                                                    output dir='model 3/MLP student smr3765'
hls model.compile()
```



QKeras for quantization-aware training (QAT)



```
# MLP architecture
# Create the student OKERAS
studentQ MLP = keras.Sequential(
       Input(shape=(30,)),
        ODense(20. name='fcl'.
                 kernel quantizer=quantized bits(9,1,alpha=1), bias quantizer=quantized bits(23,15,alpha=1)),
        QActivation(activation=quantized relu(16,15), name='relu1'),
        QDense(10, name='fc2',
                 kernel guantizer=guantized bits(9,1,alpha=1), bias guantizer=guantized bits(23,15,alpha=1)),
        OActivation(activation=guantized relu(16,15), name='relu2'),
        ODense(10, name='fc6'.
                 kernel quantizer=quantized bits(9,1,alpha=1), bias quantizer=quantized bits(23,15,alpha=1)),
        QActivation(activation=quantized relu(16,15), name='relu3'),
        QDense(4, name='output',
                 kernel quantizer=quantized bits(32,15,alpha=1), bias quantizer=quantized bits(32,15,alpha=1)),
        Activation(activation='softmax', name='softmax')
    1,
   name="student".
print qstats(studentQ MLP)
```



High-level synthesis project generated with hls4ml



layer2_t layer2_out[N_LAYER_2];
#pragma HLS ARRAY_PARTITION variable=layer2_out complete dim=0
nnet::dense<input_t, layer2_t, config2>(input_4, layer2_out, w2, b2); // fc1

layer3_t layer3_out[N_LAYER_2];
#pragma HLS ARRAY_PARTITION variable=layer3_out complete dim=0
nnet::linear<layer2_t, layer3_t, linear_config3>(layer2_out, layer3_out); // fc1_linear

layer5_t layer5_out[N_LAYER_5];
#pragma HLS ARRAY PARTITION variable=layer5_out complete dim=0
nnet::dense<layer3_t, layer5_t, config5>(layer3 out, layer5 out, w5, b5); // fc2

layer6_t layer6_out[N_LAYER_5];
#pragma HLS ARRAY_PARTITION variable=layer6_out complete dim=0
nnet::linear<layer5_t, layer6_t, linear_config6>(layer5_out, layer6_out); // fc2_linear

layer8_t layer8_out[N_LAYER_8];
#pragma HLS ARRAY PARTITION variable=layer8_out complete dim=0
nnet::dense<layer6_t, layer8_t, config8>(layer6_out, layer8_out, w8, b8); // fc3



























- Gamma/Neutron discrimination.
- Gamma/Neutron discrimination for diamond detector.
- Pest classification in fruit crops.
- Pulse shape discriminator for cosmic rays studies.
- Volcanic seismic event detection.
- Water quality monitoring applied to Dunav river.
- Object detection for adverse weather conditions, particularly haze and fog.



Gamma/Neutron discrimination

Gamma/neutron discrimination

International Centre



- Tagged dataset of **gamma and neutron events** from Deuterium-Deuterium (DD) and Deuterium-Tritium (DT) generators.
- The dataset was recorded at the Neutron Science Facility (NSF) of the Nuclear Science and Instrumentation Laboratory (NSIL), IAEA.
- The total gamma and neutron events in this dataset are 10,913 and 27,696, respectively.



Gamma/neutron discrimination



Morales, I. R., Crespo, M. L., Bogovac, M., Cicuttin, A., Kanaki, K., & Carrato, S. (2023). Gamma/neutron classification with SiPM CLYC detectors using frequency-domain analysis for embedded real-time applications. Nuclear Engineering and Technology.

Dataset from https://doi.org/10.5281/zenodo.8037059



Gamma/neutron discrimination







Publication:

Morales, I. R., Molina, R. S., Bogovac, M., Jovalekic, et al. (2024). *Gamma/neutron online discrimination based on machine learning with CLYC detectors*. IEEE Transactions on Nuclear Science.




Publication:

Morales, I. R., Molina, R. S., Bogovac, M., Jovalekic, et al. (2024). *Gamma/neutron online discrimination based on machine learning with CLYC detectors*. IEEE Transactions on Nuclear Science.





Teacher architecture with **2,623** parameters distributed in 6 hidden layers (MLP).

Compressed architecture with 217 parameters, distributed in 6 hidden layers (MLP).

Input size reduction: 35 samples of the leading edge of the pulse.

Morales, I. R., Molina, R. S., Bogovac, M., Jovalekic, et al. (2024). Gamma/neutron online discrimination based on machine learning with CLYC detectors. IEEE Transactions on Nuclear Science.



Overall accuracy

Publication:

- Teacher architecture (original): **99.00%**
- Student architecture (compressed): **98.20%**

Gamma/neutron discrimination



Morales, I. R., Molina, R. S., Bogovac, M., Jovalekic, et al. (2024). *Gamma/neutron online discrimination based on machine learning with CLYC detectors*. IEEE Transactions on Nuclear Science.



Overall accuracy

Publication:

- Teacher architecture (original): **99.00%**
- Student architecture (compressed): **98.20%**
- SoC memory footprint in terms of resource utilization @200MHz
 Artix-7 platform: below 35%

Gamma/neutron discrimination



Morales, I. R., Molina, R. S., Bogovac, M., Jovalekic, et al. (2024). *Gamma/neutron online discrimination based on machine learning with CLYC detectors*. IEEE Transactions on Nuclear Science.



Overall accuracy

- Teacher architecture (original): **99.00%**
- Student architecture (compressed): **98.20%**
- SoC memory footprint in terms of resource utilization @200MHz
 - Artix-7 platform: **below 35%**
- SoC latency
 - Zedboard platform: 45 clk cycles (@200MHz)

Publication:

Morales, I. R., Molina, R. S., Bogovac, M., Jovalekic, et al. (2024). *Gamma/neutron online discrimination based on machine learning with CLYC detectors*. IEEE Transactions on Nuclear Science.

Gamma/neutron discrimination





Volcanic seismic event detection



Copahue volcano seismic event detection

Copahue Volcano - Geological setting and data





Image from https://news.sky.com/story/chile-volcano-red-alert-issued-for-copahue-10444642

Publication:



Copahue volcano seismic event detection

Copahue Volcano - Geological setting and data



Between December 2017 and March 2018, the National University of Río Negro deployed a temporary network (called CP) of six broad-band and two-short period seismic stations.

CP covered an area of 12 km in the East-West (E-W) direction and 14 km in the North-South (N-S) direction in Caldera del Agrio.

Publication:



Copahue volcano seismic event detection

Copahue Volcano - Geological setting and data



2018-01-24T00:00:00 - 2018-01-24T23:59:59.99

Raw data* acquired from sensors installed at different stations that formed the CP network. Each station generates signals in **three channels: vertical, eastern, and northern**. These files have MSEED format and each is **24h long**.

Publication:



Copahue volcano seismic event detection

Copahue Volcano - Geological setting and data



2018-01-24T00:00:00 - 2018-01-24T23:59:59.99

Seismic traces recorded at the HIGI station.

Publication:



Copahue volcano seismic event detection

Copahue Volcano - Geological setting and data



Publication:







Copahue volcano seismic event detection

Model compression

ML-model with features as input

- Quantization 8 bits
- Pruning (30 % sparsity)



Publication:



Copahue volcano seismic event detection

Model deployment



Publication:



Copahue volcano seismic event detection

Model deployment



Publication:



Final remarks

- The proposed workflow successfully generates compressed models, leading to a **fully on-chip memory-mapped implementation** on the FPGA.
- The **integration of KD** into the ensemble of compression techniques contributes to achieving a balanced student model in terms of size, computational efficiency, and accuracy.
- The workflow addresses the entire development cycle: from the ML-based architecture training to the hardware deployment, overcoming the limitations outlined in previous works.







1st Mesoamerican Workshop on Reconfigurable X-ray Scientific Instrumentation for Cultural Heritage

Machine learning and model compression for FPGA/SoC

Antigua Guatemala, June 2025

Romina Soledad Molina, Ph.D.









The Abdus Salam International Centre for Theoretical Physics