



ALMA MATER STUDIORUM  
UNIVERSITÀ DI BOLOGNA



Istituto Nazionale di Fisica Nucleare

# Accelerating Machine Learning inference using FPGAs: the PYNQ framework tested on an AWS EC2 F1 Instance

*Dr. Marco Lorusso*<sup>1,2</sup>

Prof. Daniele Bonacorsi<sup>1,2</sup> Dr. Riccardo Travaglini<sup>2</sup>

<sup>1</sup>University of Bologna - Department of Physics and Astronomy

<sup>2</sup>National Institute for Nuclear Physics - Bologna Division

26<sup>th</sup> October 2022



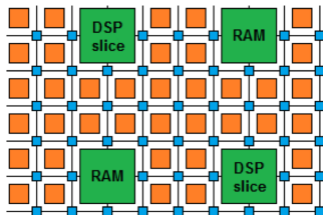
# Field Programmable Gate Array

## Field Programmable

**Gate Arrays (FPGAs)** → Middle ground between ASICs and multipurpose CPUs:

- ▶ **Programmables** to perform a wide range of tasks;
- ▶ **Low-level/Near-metal** implementation of algorithms → low latency;
- ▶ Blend the **benefits of both hardware and software**;
- ▶ Internal layout made up of *logic blocks* (**LUTs, flipflops, Digital Signal Processor slices**), embedded in a **general routing structure**.

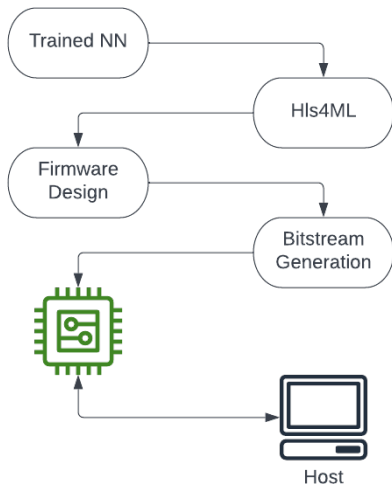
## FPGA diagram





# Implementing a Neural Network on an FPGA

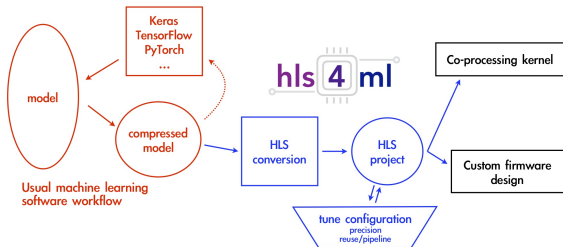
- ▶ **NN Translation into HLS** (C++) using *hls4ml* (see next slide);
- ▶ **Firmware design** (I/O interfaces);
- ▶ **Synthesis and implementation** of the design;
- ▶ Production of the **bitstream and programming** of the FPGA;
- ▶ **Running** of the inference using an application on the **host** machine.





# The hls4ml package

<https://fastmachinelearning.org/hls4ml>

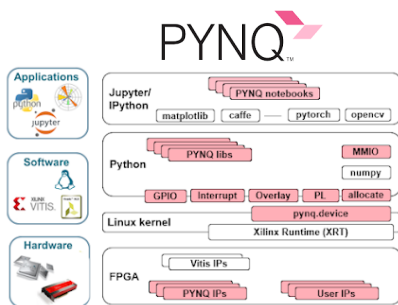


- ▶ Developed by members of the HEP community to translate **ML algorithms** written in **Python** into **High Level Synthesis** code;
- ▶ HLS allows the **generation** of **hardware descriptive code** (HDL) from *behavioral descriptions* contained in C++ program;
- ▶ The translated Python objects can be injected in the automatic workflow of proprietary software like Vivado from Xilinx Inc.



# The PYNQ project

- ▶ *PYNQ* is an **open-source** project from Xilinx®;
- ▶ It provides a **Jupyter-based framework** with Python APIs for using Xilinx platforms;
- ▶ The **Python language** opens up the **benefits of programmable logic (PL)** to people **without** in-depth knowledge of **low-level programming languages**.



<https://pynq.readthedocs.io>



## An introduction to PYNQ

- ▶ The **overlay** class is the **core** of the library;
- ▶ An overlay object is built providing the **FPGA design** to run on the PL;
- ▶ FPGA is **programmed** and relevant **interface** is available through **PYNQ API function** calls;
- ▶ It is possible to **accelerate** a software **application**, or to customize the hardware platform for a particular application.

```
1 from pynq import Overlay
2
3 overlay = Overlay("designbitstream.xclbin") # or .awsxclbin
4 result = overlay.<function described in FPGA design>
```



## The testing ground: AWS F1 Instances

Cloud computing is used to test the capabilities of these tools in preparation for deployment of FPGA accelerator cards in a local server.

- ▶ Part of the **AWS Cloud Computing** catalogue;
- ▶ EC2 F1 instances use **FPGAs** to enable **delivery** of **custom hardware accelerations**;
- ▶ Packaged with **tools** to **develop**, simulate, debug, and **compile** a design.

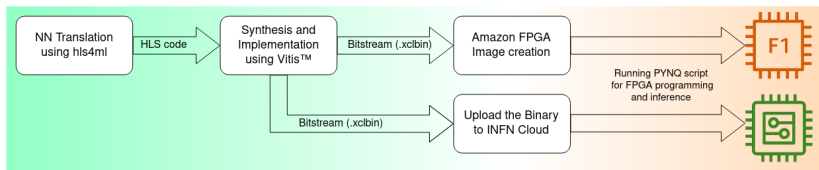




## Deploying on F1

- ▶ Follow the *Application Acceleration development flow*, offered by Vitis™, targeting data center acceleration cards;
- ▶ **Upload the bitstream** to a S3 bucket and request the **creation** of an *Amazon FPGA Image* (AFI) accessible from all F1 instances;
- ▶ Write a **Python script** using PYNQ APIs.

A "more traditional" approach is to use **OpenCL** to write the host application: both ways follow the **same** list of **basic instructions**.







## From Python to HLS Code

```
1 import tensorflow as tf
2 from qkeras.qlayers import QDense, QActivation
3
4 netinputs = tf.keras.layers.Input(shape=(4,), dtype=X_train.dtype, name="input_1")
5 x = QActivation(activation=quantized_relu(16,6,relu_upper_bound=6.0),
6               name='qrelu1')(inputs)
7 x = QDense(16, kernel_quantizer=quantized_bits(16,5,alpha=1),
8           bias_quantizer=quantized_bits(16,5,alpha=1),
9           kernel_initializer='random_normal', name='qdense_1')(x)
10 x = QActivation(activation=quantized_relu(16,6), name='qrelu2')(x)
11 ...# List of layers and activation functions
12 output = tf.keras.layers.Activation('softmax', name='soft1')(x)
13 model = tf.keras.Model(inputs=netinputs, outputs=output, name="model")
14 model.compile(optimizer='adam', loss='sparse_categorical_crossentropy')
15 history = model.fit(X_train, Y_train, epochs=num_epochs, validation_data=(X_test, Y_test))
```



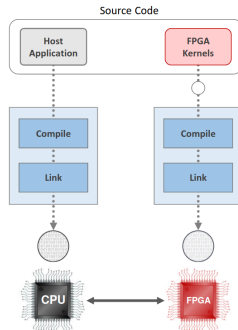
```
1 import hls4ml
2
3 config = hls4ml.utils.config_from_keras_model(model, granularity='model')
4 hls_model = hls4ml.converters.convert_from_keras_model(model,
5               hls_config=config, part='<id of FPGA model>')
6 hls_model.compile()
7 hls_model.build(csim=False, synth=False)
```



## Producing the Bitstream with Vitis

The *build* function creates the HLS code to import in the Vitis Software Platform developed by Xilinx.

- ▶ An **application project** with the target platform is created;
- ▶ The HLS code from *hls4ml* is imported as **source** for the *kernel* of the application;
- ▶ A *Hardware function* is associated to the **main C++ function** in the code;
- ▶ The **host application** is usually written in OpenCL (see next slide);
- ▶ The whole application is build for **hardware deployment** → **Bitstream**.





## OpenCL vs PYNQ

The first thing to do in both cases, is to **program the device and initialize** the software context.

```
1 auto devices = xcl::get_xil_devices();
2 auto fileBuf = xcl::read_binary_file(binaryFile);
3 cl::Program::Binaries bins{{fileBuf.data(),
  ↪ fileBuf.size()}};
4 OCL_CHECK(err, context = cl::Context({device}, NULL, 1
  ↪ NULL, NULL, &err));
5 OCL_CHECK(err, q = cl::CommandQueue(context, {device}, 2
  ↪ CL_QUEUE_PROFILING_ENABLE, &err));
6 OCL_CHECK(err, cl::Program program(context, {device} 3
  ↪ bins, NULL, &err));
7 OCL_CHECK(err, krnl_vector_add = cl::Kernel(program,
  ↪ "vadd", &err));
8
```

```
import pynq
ov =
  ↪ pynq.Overlay("model_binary.awsxcclbin")
nn = ov.myproject
```

In OpenCL host and FPGA **buffers** need to be handled separately and linked after creation; with PYNQ, the user is only presented with a single interface for both:

```
1 std::vector<int, aligned_allocator<int>>
  ↪ source_in1(DATA_SIZE);
2 OCL_CHECK(err, 1::Buffer buffer_in1(context, 1 inp = pynq.allocate(27, 'u2')
3 CL_MEM_USE_HOST_PTR | CL_MEM_READ_ONLY, 2 out = pynq.allocate(1, 'u2')
  ↪ vector_size_bytes,
4 source_in1.data(), &err))
```



## OpenCL vs PYNQ (cont'd)

To **initiate data transfers** the direction as a function parameter must be specified in OpenCL, while in PYNQ the same is done with a specific function:

```
1 OCL_CHECK(err, err =
↳ q.enqueueMigrateMemObjects({buffer_input}, 0 /*01   inp.sync_to_device()
↳ means from host*/,NULL,&eventinp));
```

To **run the kernel** in OpenCL each kernel argument need to be set explicitly using the `setArgs()` function, before starting the execution with `enqueueTask()`; in PYNQ, the `.call()` function does everything in a single line.

```
1 OCL_CHECK(err, err = myproject.setArg(0, buffer_input));
2 OCL_CHECK(err, err = myproject.setArg(1, buffer_output));
3 //[...]
4 OCL_CHECK(err, err =                               1   nn.call(inp,out)
↳ q.enqueueTask(myproject,NULL,&eventker));
5 // wait for all kernels to finish their operations
6 OCL_CHECK(err, err = q.finish());
```

Finally, the **output is retrieved** in both cases similarly to the input transfer:

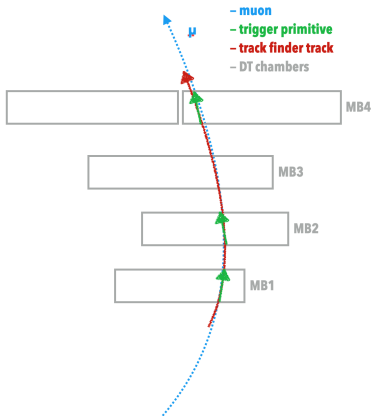
```
1 OCL_CHECK(err, err =
↳ q.enqueueMigrateMemObjects({buffer_output},   1   out.sync_from_device()
2 CL_MIGRATE_MEM_OBJECT_HOST));
```



## The tested model

To **test** the **workflow** and the **performance**, a **Neural Network** has been considered:

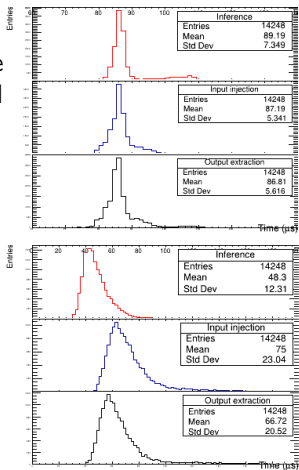
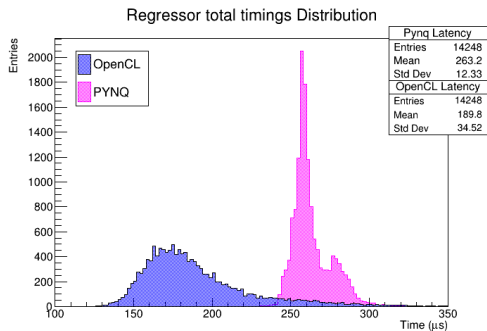
- ▶ **Regressor** in the context of Level-1 **triggering** at the CMS experiment at CERN:
  - ▶ NN predicts **transverse momentum** of **muons** using their position and direction in the detector.





# Timing Comparison

A difference in **computation times** can be seen between the same algorithm deployed with PYNQ and OpenCL:

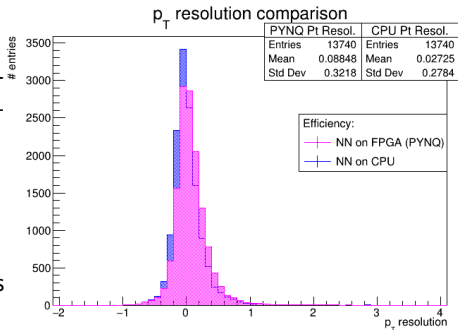




## Inference comparison

The **FPGA's output** has been **validated** against the NN run on a consumer **CPU**:

- ▶ **small difference** traceable to **quantization** of floating point to **fixed** point;
- ▶ small bias towards higher values of  $\Delta p_T/p_T$ .





## Summary and conclusions

- ▶ This work is still in **progress** (i.e. kernel optimization);
- ▶ The possibility of **deploying** a **Neural Network** on a **FPGA** inside an **AWS instance** has been explored;
- ▶ A **fast** and **easy-to-use** alternative to host applications written in **OpenCL** has been found in **PYNQ** using the **Python** programming language;
- ▶ There seems to be **no important drawbacks** from using this new approach.





Thank you!



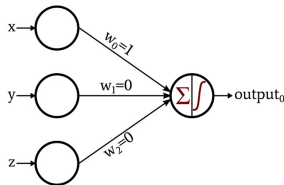
# Backup



# Artificial Neural Networks

The  $p_T$  assignment is currently carried out using **precompiled LUTs**. An alternative was explored using **Artificial Neural Network (ANN)**:

- ▶ An ANN is a network designed to tackle **non-linear learning problems**;
- ▶ The **Fully Connected Multilayer Perceptrons (MLPs)** are made up of single units called *Perceptrons*;
- ▶ Perceptrons can be stacked together to **build arbitrarily deep custom networks**;
- ▶ The NN *learns* during the **training** process by receiving **input patterns** together with the corresponding true target variable and finding the **best set of weights**;
- ▶ The weights are used to **predict the output with unseen data**.



Graphical representation of a Perceptron.



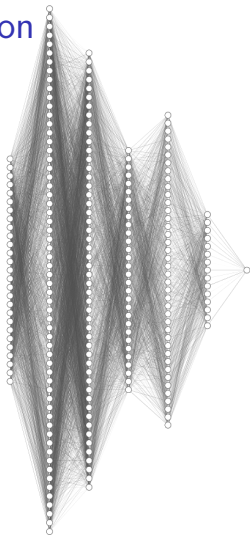
## Neural Network for regression

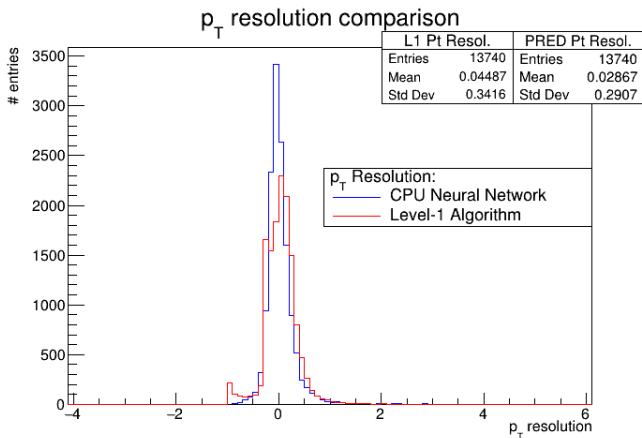
### A Fully

**Connected MLP** was built using QKeras with:

- ▶ **Input layer:** 27 features;
- ▶ **6 hidden layers:** 35, 20, 25, 40, 20, 15 nodes;
- ▶ **Output layer:** returns the  $p_T$  value.
- ▶ **Activation function:** Rectified Linear Unit;
- ▶ **Weight pruned.**

The model was **tested using a consumer CPU** before the hardware implementation.





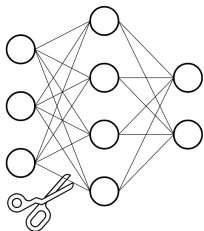
**Figure:** Transverse momentum resolution histograms computed for the machine learning model (blue) and Level-1 trigger (red) based momentum assignment.



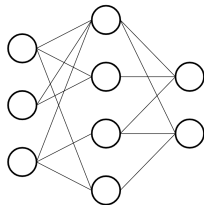
## Optimization techniques

To produce an **optimized NN** for **implementation** on an FPGA:

- ▶ **Quantization:**  
the parameters were converted **from double precision floating-points to fixed points** to exploit the efficiency of DSPs;
- ▶ **Pruning: connections**  
between nodes with low influence were **cut** to **minimize** the number of **parameters** and operations during inference and **reduce the resources** needed for implementation.



Before pruning



After pruning

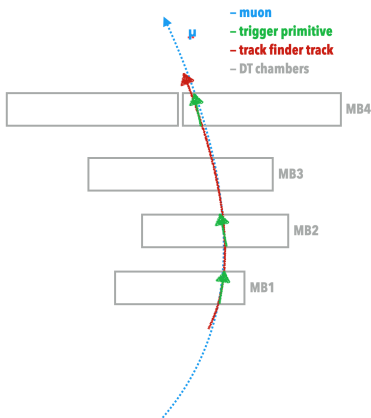




## Dataset to train and test the NN

The entire **dataset** contains about **300000 simulated muons** with a range in  $p_T$  from **3 to 200 GeV/c**. A set of **information** is included in order to **predict** the muon  $p_T$ :

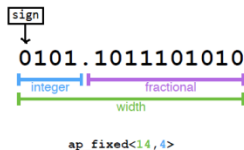
- ▶ **Trigger segments' position** (wheel, sector,  $\phi$ ) for each station crossed by the particle;
- ▶ Their **direction** in CMS global coordinates ( $\phi_b$ ).
- ▶ Trigger primitives' **quality** (i.e. number of hits used to build a segment).





## Quantization

In order to produce an **optimized NN** for **implementation** on an FPGA, the models were *quantized*:



- ▶ *Quantization* is the conversion from **high-precision floating-points to normalized low-precision integers** (*fixed-point*) parameters;
- ▶ *QKeras* is a Python package developed as a collaboration between Google and HEP researchers to **build NN with quantized parameters**;
- ▶ It has an easy-to-use API: there are **drop-in replacements** for the most common layers used with Keras (e.g. Dense → QDense).

```
1 QDense(64, kernel_quantizer = quantized_bits(6,0),  
2       bias_quantizer = quantized_bits(6,0)(x))  
3 QActivation('quantized_relu(6,0)')(x)
```

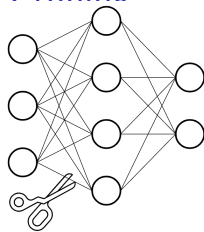




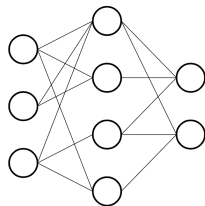
## Slimming techniques - Weight Pruning

When building a NN model, the final hardware platform where the inference computation will be run, has to be considered.

- ▶ *Weight Pruning* is the elimination of unnecessary values in the weight tensor;
- ▶ Connections between nodes with low influence are "cut" during the synthesis of the HLS design;
- ▶ This is aimed at minimizing the number of parameters and operations involved in the inference computation.



Before pruning



After pruning



## Regressor total timings Distribution

