

Accelerating Machine Learning inference using FPGAs: a crash course

Dr. Marco Lorusso^{1,2}

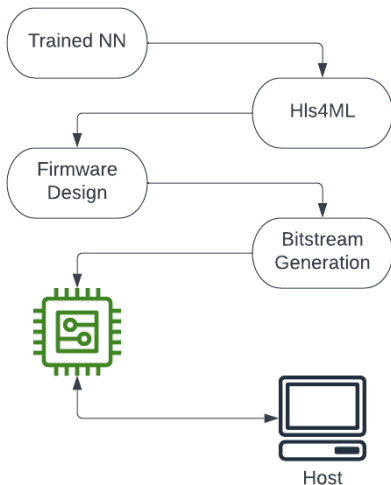
¹University of Bologna - Department of Physics and Astronomy

²National Institute for Nuclear Physics - Bologna Division

4th November 2022

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.



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=netoutput, 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))
  
```



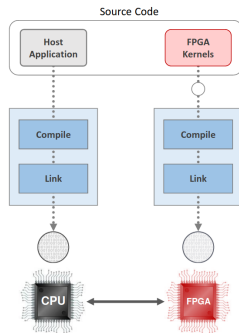
```

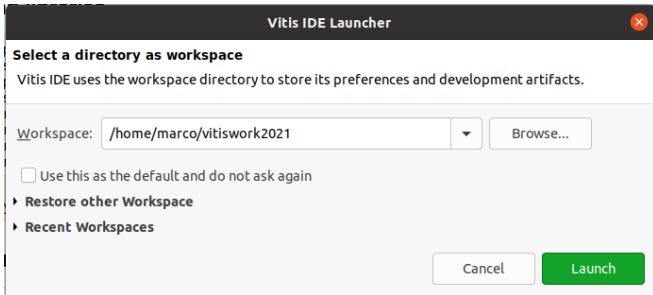
1  import hls4ml
2
3  config = hls4ml.utils.config_from_keras_model(model, granularity='model')
4  hls_model = hls4ml.converters.convert_from_keras_model(model,
5                hls_config=config, part='<id of FPGA model>')
6  hls_model.compile()
7  hls_model.build(csim=False, synth=False)
  
```

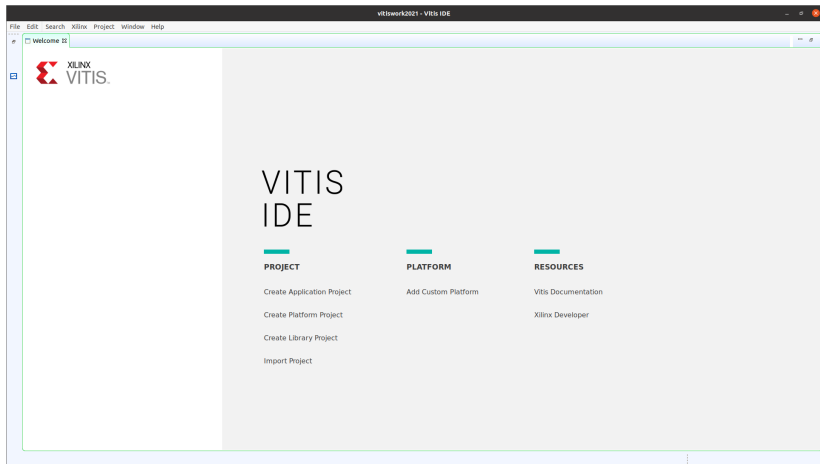
Producing the Bitstream with Vitis

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;
- ▶ The whole application is build for **hardware deployment** → **Bitstream**.





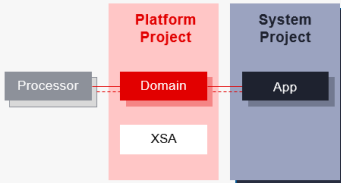


□ ×
New Application Project

Create a New Application Project

This wizard will guide you through the 4 steps of creating new application projects.

1. Choose a **platform** or create a **platform project** from Vivado exported XSA
2. Put application project in a **system project**, associate it with a processor
3. Prepare the application runtime – **domain**
4. Choose a template for application to quick start development



The diagram illustrates the project structure. On the left is a grey box labeled 'Processor'. To its right is a red box labeled 'Platform Project' which contains a red box labeled 'Domain' and a white box labeled 'XSA'. To the right of the 'Platform Project' is a blue box labeled 'System Project' which contains a dark grey box labeled 'App'. Dashed lines connect the 'Processor' to the 'Domain' box, and the 'Domain' box to the 'App' box.

- A platform provides hardware information and software environment settings.
- A system project contains one or more applications that run at the same time.
- A domain provides runtime for applications, such as operating system or BSP.
- A workspace can contain unlimited platforms and unlimited system projects.

Skip welcome page next time. (Can be reached with Back button)

?

< Back
Next >
Cancel
Finish




New Application Project

Platform

Choose a platform for your project. You can also create an application from XSA through the 'Create a new platform from hardware (XSA)' tab.

Select a platform from repository
 Create a new platform from hardware (XSA)

Find: + Add ⚙ Manage

Name	Board	Flow	Vendor	Path
 xilinx_aws-vu9p-f1_shell-v042618	aws-vu9p-f	DataCenter Accel	xilinx	/home/marco/aws-fpga/Vitis/aws_platform/xilinx_
 xilinx_u50_gen3x16_xdma_20192	u50	DataCenter Accel	xilinx	/opt/xilinx/platforms/xilinx_u50_gen3x16_xdma_20
 xilinx_u250_gen3x16_xdma_3_1_	u250	DataCenter Accel	xilinx	/opt/xilinx/platforms/xilinx_u250_gen3x16_xdma_3

Platform Info

General Info

Name: xilinx_aws-vu9p-f1_shell-v042618

Part: xcvu9p-flgb2104-2-i

Family: virtexuplus

Description: {No description given}

FPGA part: xcvu9p-flgb2104-2-i

Number of DDRs: four

Acceleration Resources

Clock Frequencies	
Clock	Frequency (MHz)
CPU	1
PL 0	250.000000
PL 1	500.000000
PL 2	250.000000
PL 3	125.000000

Domain Details

Domains	
Domain name	Details
x86_0	CPU: x86_0 OS: Linux OS

New Application Project
□ ×

Application Project Details
☰

✖ Invalid application project name. Project name must be specified

Application project name:

System Project

Create a new system project for the application or select an existing one from the workspace i

Select a system project

- awsiris_system
- awsirislite_system
- awsnn_system
- + Create new...

System project details

System project name:

Target processor

Select target processor for the Application project.

Processor	Associated applications
x86 SMP	

Show all processors in the hardware specification i

?

< Back
Next >
Cancel
Finish

New Application Project
□ ×

Templates

Select a template to create your project.

Available Templates:

Find: ☰ ⊕

▼ Acceleration templates with PL and AIE accelerators

- Empty Application
- Empty Application (XRT Native API's)
- Vector Addition

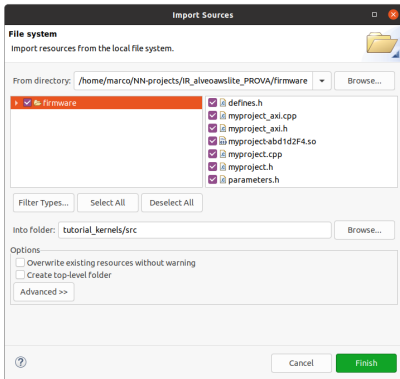
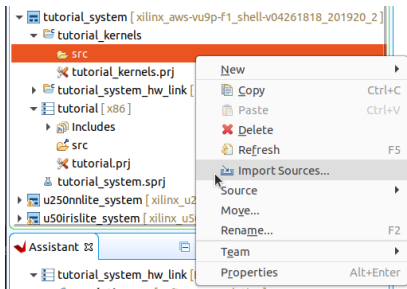
Empty Application

Creates a new Empty application

Show only certified examples

?

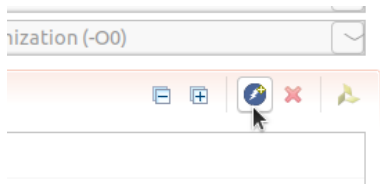
< Back
Next >
Cancel
Finish

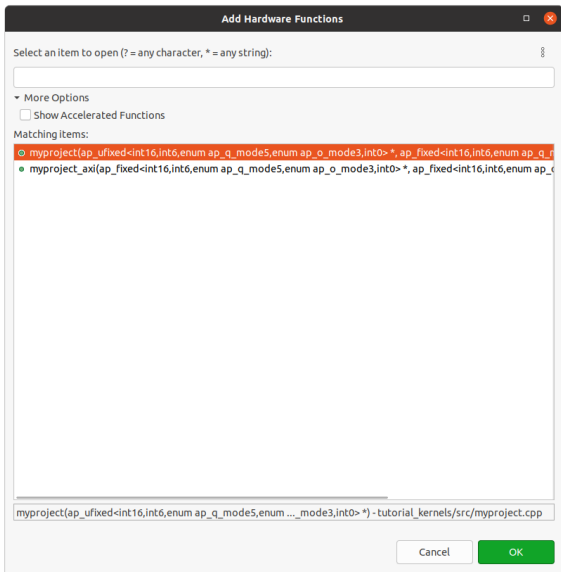


```

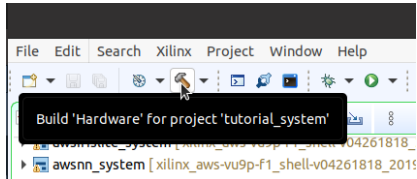
void myproject(
    input t input_1[N_INPUT_1],
    result_t layer8_out[N_LAYER_6],
) {
    //hls-fpga-machine-learning insert IO
    //#pragma HLS ARRAY_RESHAPE variable=input_1 complete dim=0
    //#pragma HLS ARRAY_PARTITION variable=layer8_out complete dim=0
    //#pragma HLS INTERFACE ap_vld port=input_1,layer8_out
    #pragma HLS PIPELINE

    unsigned short const_size_in_1 = N_INPUT_1;
    unsigned short const_size_out_1 = N_LAYER_6;
  
```





Active build configuration: Hardware



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.

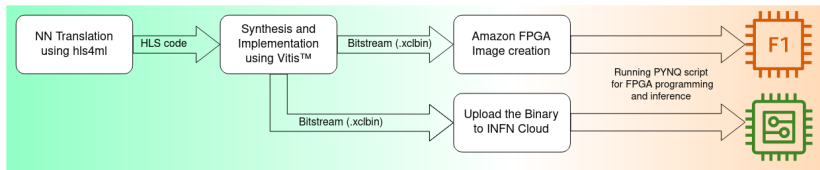
- ▶ Part of the **AWS Cloud Computing** catalogue;
- ▶ EC2 F1 instances use **FPGAs** to enable **delivery** of **custom hardware accelerations**;
- ▶ Packaged with **tools** to **develop**, simulate, debug, and **compile** a design.



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**.



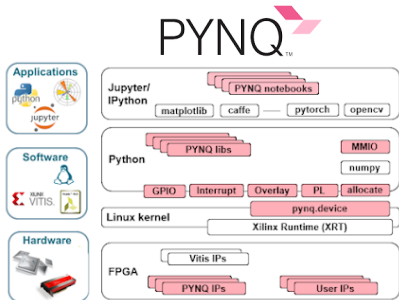
Amazon FPGA Image

- ▶ The *aws-fpga* repository contains all the tools needed for deploying (and developing) on a F1 instance;
- ▶ The *awsxclbin* (AFI) can be created by running the `create_vitis_afi.sh` script which is included in the `Vitis/tools/` directory;
- ▶ Before running the command, make sure that `aws-fpga/vitis_setup.sh` has been sourced;
- ▶ Remember to configure the AWS CLI and set up the bucket region, e.g. `aws configure set region us-east-1`;
- ▶ Create an AFI by running:

```
1 aws-fpga/Vitis/tools/create_vitis_afi.sh -xclbin=<filename>.xclbin  
↔ -s3_bucket=<bucket-name> -s3_dcp_key=<dcp-folder-name>  
↔ -s3_logs_key=<logs-folder-name>
```

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>
```

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  import pynq
   ↪ NULL, NULL, &err));
5  OCL_CHECK(err, q = cl::CommandQueue(context, {device}, 2  ov =
   ↪ CL_QUEUE_PROFILING_ENABLE, &err));
6  OCL_CHECK(err, cl::Program program(context, {device}) 3  nn = ov.myproject
   ↪ bins, NULL, &err));
7  OCL_CHECK(err, krnl_vector_add = cl::Kernel(program,
   ↪ "vadd", &err));
8

```

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, l1::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());
  
```

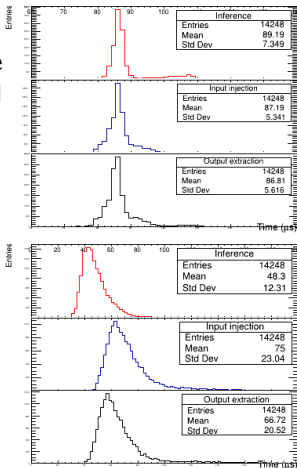
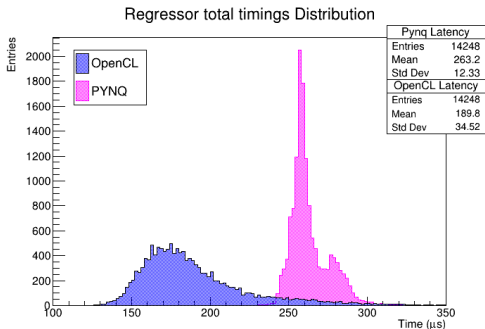
Finally, the **output is retrieved** in both cases similarly to the input transfer:

```

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

Timing Comparison

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



Thank you!

Backup

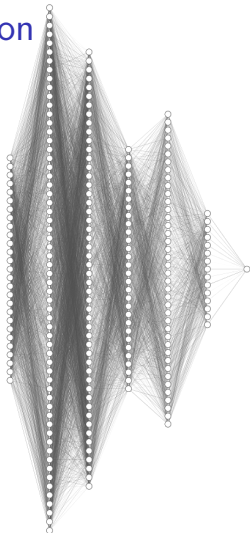
Neural Network for regression

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.



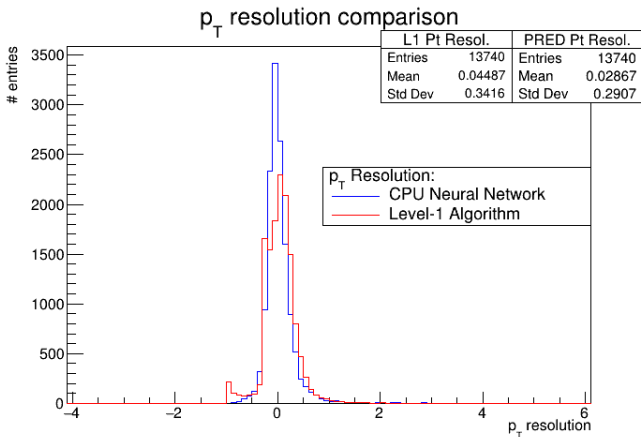
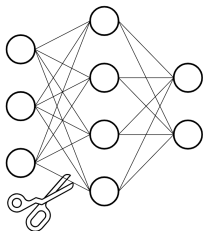


Figure: Transverse momentum resolution histograms computed for the machine learning model (blue) and Level-1 trigger (red) based momentum assignment.

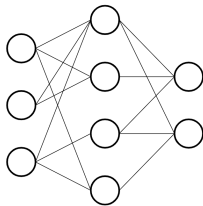
Optimization techniques

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

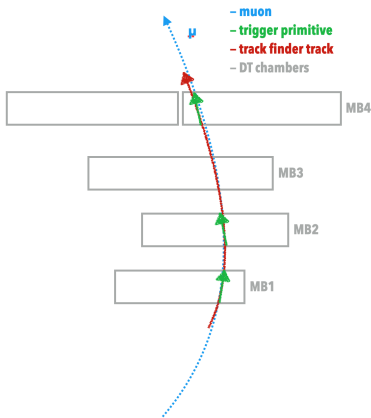


After pruning

Dataset to train and test the NN

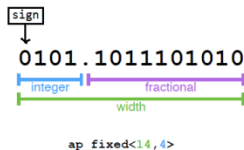
The entire **dataset** contains about **300000 simulated muons** with a range in p_T from **3 to 200 GeV/c**. A set of **information** is included in order to **predict** the muon p_T :

- ▶ **Trigger segments' position** (wheel, sector, ϕ) for each station crossed by the particle;
- ▶ Their **direction** in CMS global coordinates (ϕ_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).

```

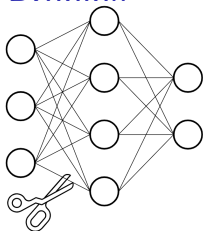
1 QDense(64, kernel_quantizer = quantized_bits(6,0),
2         bias_quantizer = quantized_bits(6,0)(x))
3 QActivation('quantized_relu(6,0)')(x)

```

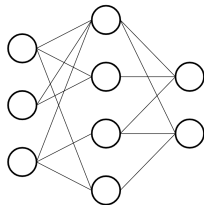
Slimming techniques - Weight Pruning

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

