



### **FPGAs in Detector Instrumentation: concepts and trends**

Michele Caselle

19-21 February 2020 Best Western Hotel

Helmholtz Research Field Matter



# **Challenges for the future**



The next generation of detectors are extremely challenging: HEP, astrophysics, photon science, etc.

#### **Particle physics**

Unprecedent luminosity operations,

4D detectors: excellent spatial and time resolution

*HL-CMS*: Unprecedented data rate of up to 50 Tb/s to be processed in < 4  $\mu$ s with high efficiency





#### **Photon Science**

Terapixel per second imaging:

- 100 million pixel,
- MHz- frame rates
- High dynamic range

#### **Astroparticle physics**

Cryogenic detectors of unique energy resolution for dark matter searches and neutrino physics

CTA, IceCube-Gen2, etc.





# Accelerators and beam physics

Complex dynamics on short time scale,

Multi-spectral THz detectors for beam diagnostics, for plasma accelerators

### **FPGAs in Detector Instrumentation - Outlook**



**Motivations** - How to cope with the data deluge from the next generation of detectors?

**Novel heterogenous programmable devices** – What is the evolution of Field Programmable Gate Array (FPGA) devices?

**Machine learning (ML) and artificial intelligence** - *Will enable us to process the data deluge in real-time?* 

Al on FPGA - How develop a fast ML inference on FPGA?

**Applications:** fast track-reconstruction based on ML, control of the beam dynamics in complex synchrotron machines e.g. autonomous accelerator?



The first commercially viable Field Programmable Gate Array (FPGA), named XC2064, has been invented in 1985 by Ross Freeman and Bernard Vonderschmitt, the Xilinx co-founders.

• A FPGA is a semiconductor *Integrated Circuit* (*IC*) device on which the function can be defined after manufacturing ("in the field") using software-like languages (ex: VHDL, Verilog). FPGAs can be reconfigured at any times.



Michele Caselle

• A Look-Up Table (LUT) can implement any arbitrary Boolean function of its inputs and can be cascaded to other LUTs to perform more complex functions







6





Zynq ZU11EG Ultrascale+ (Xilinx)





#### Block RAM 21 Mbit + 22 Mbit (Ultra RAM)

DSP -Digital Signal Processing (2928 units)

7





Zynq ZU11EG Ultrascale+ (Xilinx)



CLBs - Configurable Logic Blocks (> 1 Million of CLBs)





Zynq ZU11EG Ultrascale+ (Xilinx)



CLBs - Configurable Logic Blocks (> 1 Million of CLBs)

Sea of Programmable Logic and Routing resourses

Latency & bandwidth comparable to ASICs

9

Michele Caselle



### From FPGA ...

### ... to Adaptive Compute Acceleration Platform



Michele Caselle

### A New Class of Devices for Today's Challenges



**Device Category** 

# Heterogeneous: ZYNQ MPSoC technology



Heterogeneous platform of the Zynq System-on-Chip (SoC) integrates, in a monolithic device, FPGA resources with a back-end software running on a hard-core ARM-based processor.





System-on-Chip (SoC) workshop – CERN , 12-14 June 2019 https://indico.cern.ch/event/799275/

Ref. https://www.xilinx.com/support/documentation/selection-guides/zynq-ultrascale-plus-product-selection-guide.pdf

# Heterogeneous: ZYNQ RFSoC technology



The bandwidth bottleneck of previous ZYNQ MPSoC was the bandwidth and complexity of the

JESD204B standard communication for fast ADC/DAC



 $Ref.\ https://www.xilinx.com/support/documentation/selection-guides/zynq-usp-rfsoc-product-selection-guide.pdf$ 

# Versal – ACAP (Overview)



New Xilinx architecture Versal - ACAP (Adaptive Compute Acceleration Platform) develop in TSMC 7nm







### Next generation of data processing

# **Machine Learning**





The generation of detectors will face a peak luminosity of  $30 \times 10^{34}$  cm<sup>-2</sup>/s with a pile-up 1000 and radiation levels that are 1-2 orders of magnitude larger than those at the HL-LHC. Timing precision of 5 – 10 ps is required for pile-up mitigation. *The resulting data rates lie in the hundreds of TB/s.* 



# What's the different between Artificial lintelligence, Machine Learning, and Deep Learning?

Al involves machines that can perform tasks that are characteristic of human intelligence While this is rather general, it includes things like planning, understanding language, recognizing objects and sounds, learning, and problem solving



#### At its core, machine learning is simply a way of achieving AI

*"the ability to learn without being explicitly programmed."* You see, you can get a traditional program without using machine learning, but this would require building millions of lines of codes with complex rules and decision-trees.



# Deep learning is one of many approaches to machine learning

Deep learning was inspired by the structure and function of the brain, namely the interconnecting of many neurons. Artificial Neural Networks (ANNs) are algorithms that mimic the biological structure of the brain.

# **Machine learning**





**Supervised**: All data is *labeled* and the algorithms learn to predict the output from the input data.

**Unsupervised**: All data is *unlabeled* and the algorithms learn to inherent structure from the input data

**Reinforcement:** The learning algorithm is trained not on present data but rather based on a *feedback* system.



# **Supervised Machine Learning**





# Reinforcement ML – Inverted pendulum application



Inverted pendulum

- GOAL: to keep the pendulum in equilibrium
- **Action:** to apply a "torque" or "rotational force"  $(\tau)$
- Reinforcement ML: which learns from the rewards / mistakes
- ML running on FPGA



Reinforcement





# **Machine learning**





**Supervised**: All data is *labelled* and the algorithms learn to predict the output from the input data.

**Unsupervised**: All data is *unlabelled* and the algorithms learn to inherent structure from the input data

**Reinforcement:** The learning algorithm is trained not on present data but rather based on a *feedback* system.

The idea behind a deep neural network is to mimic the biological structure of the brain with a similar structure with layers of artificial neurons

**Deep learning** is part of a broader family of machine learning methods based on artificial neural networks. Learning can be supervised, unsupervised and reinforcements

# **Deep Neural Network**















Filter 1







Filter 3



Filter 4













A deep neural network consists of a hierarchy of layers, whereby each layer transforms the input data into more abstract representations (e.g. edge -> nose -> face). The output layer combines those features to make predictions





#### Forward propagation



# How to implement a fast ML inference on FPGA?

#### DNN produces many parameters for the model, which increases the compute cost and requires high memory bandwidth. There are two main ways to optimize a DNN application:

**Pruning:** This is a form of DNN compression. It reduces the number of "synaptic" connections and "neurons" and so that the overall amount of data is reduced. Typically, weights close to zero are removed.

#### What we lose in accuracy?

Comparison GPU (FP 32) vs FPGA (INT8)

Top-5 Accuracy	FP-32	FIXED-16 (INT16)	FIXED-8 (INT8)	Difference vs FP32
VGG-16	86.6%	86.6%	86.4%	(0.2%)
GoogLeNet	88.6%	88.5%	85.7%	(2.9%)
SqueezeNet	81.4%	81.4%	80.3%	(1.1%)

**Quantization:** To bring the neural network to a reasonable size while also achieving high-performance accuracy. In this method, the process of approximating a neural network that uses floating-point numbers (FTP32) by a neural network of low-bit width numbers (INT8) is performed.

### High-Performance DNN on FPGA





### Xilinx Deep Learning Processing Unit (DPU) which is a configurable computation engine dedicated to

convolutional neural networks

fpgaConvNet: http://cas.ee.ic.ac.uk/people/sv1310/ fpgaConvNet.html, CNNECST, many others ...

Custom frameworks: Snowflake: arXiv:1708.02579, DNNWeaver: http://act-lab.org/artifacts/dnnweaver/,

- 2 \_\_\_\_\_
- 33 7<sup>th</sup> KSETA Plenary Workshop 2020

# Machine learning on FPGA

Nowadays, several High Level tools are available to generate "fast" ML inferences running on FPGA

- High Level Synthesis four Machine Learning, developed at CERN
  - What is hls4ml → framework removes major barrier on hardware development of ML algorithms allowing developers with little or no FPGA expertise to program the FPGA
- Machine learning for Data Quality Monitoring (on-line), developed at CERN (HL-CMS)
  - Based on supervised learning  $\rightarrow$  binary classification problem: good plot vs bad plot









# hls4ml – Workflow



Very success framework, developed at CERN for CMS jet classification is now integrated and distributed in the Xilinx Software Development Kit (XSDK)

Webpage: https://fastmachinelearning.org/hls4ml/



# hls4ml – reconstruction chain for jet



**GOAL:** Machine Learning for low-level trigger system based on FPGA  $\rightarrow$  fast jet substructure





FPGA	vs	GPU
------	----	-----

Res_V	FPGA	Python Keras		
	prediction	calculation (GPU)		
Gluon (g)	0.118164	0.12993355		
Quark (q)	0.639648	0.6487177		
Boson (W)	0.118164	0.10633943		
Boson (Z)	0.118164	0.10616959		
top quark (t)	0.015625	0.00883975		
y time:	Cou	Courtesy: Weijia Wang		

GPUs: of ~ tens µs

FPGA: 16 clock cycles @ 200 MHz = 80 ns

classification

# Xilinx Deep Learning Processing Unit (DPU)

- Karlsruher Institut für Technologie
- The Xilinx Deep Learning Processor Unit (DPU) is a configurable computation engine dedicated for convolutional neural networks.
- The DNNDK (Deep Neural Network Development Kit) is designed as an integrated framework, which aims to simplify and accelerate deep learning application development and deployment on the Deep Learning Processor Unit (DPU).





more parallels DPUs could be instantiated which will work in parallel

# Examples: Objects recognition/classification

Karlsruher Institut für Technologie

- Neural network implemented resnet50
- Xilinx Deep Learning Processor Unit (DPU)



#### Courtesy: Weijia Wang



High-Flex 2 multi-purpose PCIe card



cat.jpg



xian.jpg



cab.jpg



HF-2.jpg



Michele Caselle

### Examples: Advanced driver-assistance systems (ADAS)

Karlsruher Institut für Technologie

**GOAL:** Real-time recognition of cars, trucks, people, bikes, motorcycles, etc.

Video stored in the on-board ETH DDR4 memory and processed by yolov3detection on FPGA UFO Camera Courtesy: Weijia Wang tioning of the states in the Video Analysis @Xilinx DPU High-Flex 2 multi-purpose PCIe card nvolution laye **IPE – ADAS Demo** Yolov3 architecture

# Advanced beam control (KARA)



- **Goal:** to keep the coherent synchrotron radiation (CSR) THz intensity stable by fast feedback to RF System
- Complex and nonlinear dynamics in longitudinal / transverse bunch profiles → described by a nonlocal nonlinear partial differential vlasov fokker planck equation describing the time evolution of the probability distribution of a particle in synchrotron machine.



0.2

 $^{-4}$ 

energy deviation ( $\sigma_{E,0}$ 

- We would like to control the phase space of the bunches by ML on FPGA
- Action: to control the nonlinear dynamics of the phase space of the beam by a modulation of the RF system
- **RF System parameters are**: RF Amplitude, phase and frequency



#### Institute for Data Processing and Electronics (IPE)

10

10 time (T<sub>s</sub>)

RF amplitude modulation

15

15

0.136

0.13/

1.006

1.004

1.002

0.998

0.996

0.994





42 7<sup>th</sup> KSETA Plenary Workshop 2020

#### Courtesy: Weijia Wang and Tobias Boltz

# **PANDA** – fast track-reconstruction



- Trigger-less readout system (no hardware trigger)
- **GOAL**: to explore ML for fast-track reconstruction of PANDA
  - Forward Tracking System detector
  - Barrel Tracking detector

Courtesy: Waleed Esmail (Forschungszentrum Jülich) and Weijia Wang

### **Two-steps fast-track reconstruction**



**STEP 1**: Find "local" *Track Segments* by DNN (under investigation by unsupervised ML)



### **Two-steps tracking reconstruction**



STEP 1: Find "local" *Track Segments* by CNN (under investigation by unsupervised ML)

To find local track segments



### **Two-steps tracking reconstruction**



- By Unsupervised (*K*-means) and supervised (recurrent neural network) learning:
- STEP 1: Create *Track Segments* by using K-means clustering algorithm

To find local track segments

STEP 2: interpolate the track segments from the different parts of the FTS to form a full track candidate, it is based on a Recurrent Neural Network (RNN)





Recurrent neural networks have a memory that enables them to remember important events that happened many time steps in the past.

For: handwriting recognition or speech recognition.

# **PANDA – A possible ML implementation**

- Karlsruher Institut für Technologie
- Distributed AI architecture running on heterogenous FPGA-GPUs infrastructure, local track reconstruction on FPGA and RNN on GPUs and integrated within the on-line event selection framework



Courtesy: Waleed Esmail (Forschungszentrum Jülich) and Weijia Wang

	FT1,FT2	FT3,FT4	FT5,FT6			
Purity	99%	99%	99%			
Efficiency	96%	95%	96%			
Local track segment "ONLY" Preliminar						

The *purity* specifies which fraction of hits in one track come from the correct particle.

The *efficiency* is defined as the ratio of the number of correctly reconstructed tracks to all generated tracks.

### Impact

... we are just at the beginning of a

new disruptive technology ....



Motivations - Next generation of detectors will generate huge and complex data volume (petabyte/sec), serious technological challenges → to push the technology envelope (especially for trigger and data acquisition)



- **Novel programmable** devices families with dedicated *AI engines* opens new prospective for future detectors.
- New generation of development tools for users without knowledge of deep-learning or FPGA.

What's about FPGA also for data analysis ?

Thank you for your attention



### **Backup slides**

# **FPGA** in detectors





etc. (Wideband Readout System for Microwave-Resonator Multiplexed Sensors - Nick Karcher)

• Low-level trigger process for real-time event filtering (*FPGA-based real-time track reconstruction for the CMS phase-2 tracker upgrade -Luis Eduardo Ardila Perez*)

VME –CERN cards with # 5 FPGAs: Real-time data processing, merging and formatting and low-level trigger

# Three Ages of FPGA



# **Charting an Aggressive Course Forward**



# Two FPGA technology solutions





Virtex Ultrascale+



HBM (High Bandwidth Memory) integrated into the FPGA

- DRAM: up to 4 GByte
- Bandwidth: 1.84 Tb/s

Ref. Xilinx white paper WP380

Institute for Data Processing and Electronics (IPE)



#### Institute for Data Processing and Electronics (IPE)



# High-performance heterogeneous FPGA-GPU DAQ

- Modern photon science detectors generate huge raw data volumes (~120 Gb/s)
- Observed slow changes in synchrotron machine (e.g. current)  $\rightarrow$  sec hrs



PCIe readout card of UFO DAQ Platform

- DMA working close to theoretical limit of data link
- Data latency 5 times better than other DMA architectures
- M. Caselle et al., JINST 12 C03015 (2017)

- Heterogeneous FPGA/GPU-based readout system, the UFO DAQ platform, has been developed <a href="http://ufo.kit.edu/ufo">http://ufo.kit.edu/ufo</a>
- The core component is a "novel" Direct Memory Access (DMA) architecture
- Direct FPGA ↔ GPU communication enables real-time data processing



# Versal – Adaptable Engines

For traditionalists: This is the FPGA part

#### Some known facts

- 6 Input LUTs
- Each CLB has 32 LUTs and 64 FF (4x density compared to US+)
- 16 LUTs in a slice can be
  - a 64 bit RAM
  - 32-bit shift registers (SRL32) or two SRL16
- Internal connection of LUTs possible
- 4x clock, 4x set/reset, 16 clock enable
- 3 step voltage-scaling supported





# Versal – Al tile architecture

Karlsruher Institut für Technologi

1.3 GHz VLIW / SIMD vector processors

#### Parallelity

- VLIW: 7+ operations / clock cycle
- SIMD: 512 bit vector datapath (8 / 16 / 32 bit & SPFP operands)
- Up to 128 INT8 MACs / clock cycle / core

#### Memory

- 16 KB Internal program memory
- 32 KB data memory (parallel)
- Integrated DMA logic



# CPU & GPU vs Versal-ACAP

GPU & CPU - Mismatched Throughput



In heterogenous CPU / GPU systems the GPU, the other performance-critical functions of the application must still run in software (CPU)

Data copy from system memory to GPU global memory is still a limitation

High flexibility degree

Xilinx – Matched Throughput



In Versal devices the other performance-critical functions will be implemented in the FPGA, therefore, deep pipeline implementation on massive parallel operations

Data copy/moving by a dedicated high-bandwidth "Network on Chip" bus

Less flexible, but ...



Vitis simplifies the use of deep-learning neural networks, even for users without knowledge of deep-learning or FPGAs. The Vitis AI Library allows users to focus more on the development of their applications, rather than the underlying hardware.

# **RF-ADC/DAC Implementation steps**



# **Unsupervised Learning**



Unsupervised learning finds hidden patterns or intrinsic structures in data. **Clustering** is the most common unsupervised learning technique. *It is used for exploratory* 



In cluster analysis, data is partitioned into groups based on some measure of similarity or shared characteristic. Clusters are formed so that objects in the same cluster are very similar and objects in different clusters are very distinct.

### **Unsupervised Learning**





Result: Dendrogram

relationship

between clusters

showing the hierarchical

#### k-Means

data into k number of mutually exclusive clusters

#### Best Used...

- When the number of clusters is known
- For fast clustering of large data sets



How to work: Produces nested sets of clusters by analyzing similarities between pairs of points

#### Best Used...

When you don't know in advance how many clusters are in your data

#### k-Medoids

How to work: Similar to k-means, but with the requirement that the cluster centers coincide with points in the data.

**Result:** Cluster centers that coincide with data points

#### **Self-Organizing Map**

How It Works: Neural-network based clustering that transforms a dataset into a topology-preserving 2D map.

#### Best Used...

To visualize high-dimensional data in 2D or 3D



**Result:** Lower-dimensional (typically 2D) representation



<sup>63 7&</sup>lt;sup>th</sup> KSETA Plenary Workshop 2020



### **Convolution calculation**

### Convolution



#### 



Neural network for physics:

- **Supervised:** Neural network trained by Monte Carlo simulations
- Unsupervised (no trained): to find "New Physics" without being explicitly programmed





### **CMOS technology**



#### 

# **Machine learning**





**Supervised**: All data is *labelled* and the algorithms learn to predict the output from the input data.

**Unsupervised**: All data is *unlabelled* and the algorithms learn to inherent structure from the input data

**Reinforcement:** The learning algorithm is trained not on present data but rather based on a *feedback* system.

**Classification techniques** predict discrete responses, classify input data into categories. Example: whether an email is genuine or spam, or whether a tumor is cancerous or benign.

- **Regression techniques** predict continuous responses. Example: changes in temperature or fluctuations in power demand.
- **Clustering** is employed for exploratory data analysis to find hidden patterns or groupings in data.
- Reinforcement is employed for fast feedback control systems (i.e. beam stability in accelerators).

### From data to understanding ...





The generation of detectors will face a peak luminosity of  $30 \times 10^{34} \text{ cm}^{-2}/\text{s}$  with a pile-up 1000 and radiation levels that are 1-2 orders of magnitude larger than those at the HL-LHC. Timing precision of 5 – 10 ps is required for pile-up mitigation. *The resulting data rates lie in the hundreds of TB/s.* 



#### Data Analysis

Traditional data analysis techniques in HEP use a sequence of Boolean decisions followed by statistical analysis on the selected data.

The last big success of the traditional science is the Higgs boson



Find "a new physics" without being explicitly written in the analysis

Mathematical models learnt from data

code.

#### From Detectors and DAQ systems

- The data deluge MUST BE readout and real-time processed without any data/information loss
- The radiation load on the detectors  $\rightarrow$  will dramatically impact on the noise/efficiency of each single sensitive cell



# **Examples: Objects recognition/classification**





Michele Caselle