



Finanziato  
dall'Unione europea  
NextGenerationEU



Ministero  
dell'Università  
e della Ricerca



Italiadomani  
PIANO NAZIONALE  
DI RIPRESA E RESILIENZA

# Report su WP4

S. Gennai (INFN-MiB) & A. Pompili (UniBA)



## WP4: **missione**, **obiettivi**, **“milestone”** (in a nutshell)

### **Boosting the computational performance of Theoretical & Experimental Physics algorithms:**

- **4.1** Tools & guidelines for developing and porting heterogeneous codes and algorithms on modern architectures
- **4.2** Competence & Training Center for heterogeneous computing

O4.1 : document best practises & SW tools for codes' development/porting to heterogeneous platforms (GPUs, FPGAs)

O4.2 : prepare and support the R&D testbed to offer multiple architectures; **optimize single-node performance**

O4.3 : organize training opportunities open to external users; trained personnel will help to boost the activities.

M9-15 : report on best practises for heterogeneous computing

M22-26 : first training opportunity; testbeds ready for users; user support in place

M25-36 : results from testbed & benchmarking activities

M36 : final report on technologies, training & support system. White Paper for use cases external to the CN.



## Meeting nel 2023

- Per ora non si discutono argomenti tecnici nei nostri meeting
  - Organizziamo la preparazione degli eventi formativi  
... quindi i meeting si addensano a ridosso di questi
- Siamo comunque in contatto con alcuni WP e seguiamo le discussioni nei loro meeting
  - e.g. Flagship di nostro interesse in WP2
  - WP1 e WP3 sono - almeno per ora - "piu' distanti"
- Non appena le varie applicazioni cominciano "a macinare" vorremmo avere presentazioni dedicate che affrontino in dettaglio la parte più tecnologica coinvolta (rispetto magari alle presentazioni in WP1, 2, 3).  
In alternativa potremmo pensare a meeting congiunti?

### December 2023

14 Dec [Bi-weekly meeting WP4](#)

### November 2023

27 Nov - 30 Nov [Introductory course to HLS FPGA programming.](#)

16 Nov [Bi-weekly meeting WP4](#)

### October 2023

26 Oct [Bi-weekly meeting WP4](#)

12 Oct [Bi-weekly meeting WP4](#)

### September 2023

28 Sept [Bi-weekly meeting WP4](#)

### July 2023

11 Jul [Bi-weekly meeting WP4](#)

### June 2023

19 Jun - 21 Jun [First course about the porting on GPUs of code and algorithms](#)

15 Jun [Bi-weekly meeting WP4](#)

### May 2023

18 May [Bi-weekly meeting WP4](#)

04 May [Bi-weekly meeting WP4](#)

### April 2023

21 Apr [Bi-weekly meeting WP4](#)

06 Apr [Bi-weekly meeting WP4](#)

### March 2023

23 Mar [Bi-weekly meeting WP4](#)

### September 2022

30 Sept [kick-off meeting](#)



## Landscape document



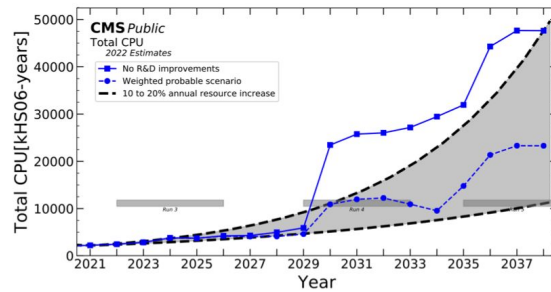
### Best practices for heterogeneous computing

#### Index:

- **Introduction**
- **Estimates for computing needs for the current/future generation of HEP experiments**
- **The need for heterogeneous computing**
  - Performance/cost considerations
  - Power consumption considerations
  - The utilization of HPC systems
- **Solutions available/under R&D**
  1. GPU
    - Portability solutions under performance studies
  2. ARM
  3. FPGA
- **Link to the proposed flagships**
- **Bibliography**



# Landscape document



## Best practices for heterogeneous computing

### Index:

- **Introduction**
- **Estimates for computing needs for the current/future generation of HEP experiments**
- **The need for heterogeneous computing**
  - Performance/cost considerations
  - Power consumption considerations
  - The utilization of HPC systems
- **Solutions available/under R&D**
  1. GPU
    - Portability solutions under performance studies
  2. ARM
  3. FPGA
- **Link to the proposed flagships**
- **Bibliography**

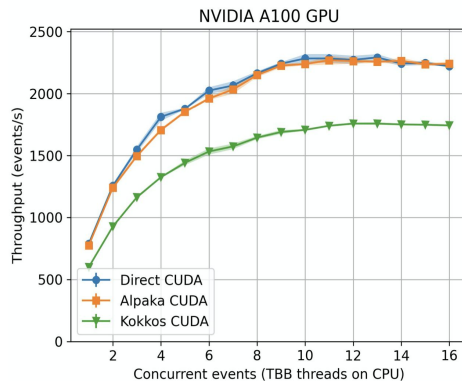
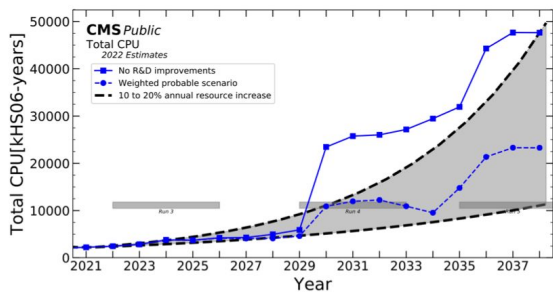
## Introduction

Heterogeneous computing refers to the use of different types of hardware platforms and accelerators (e.g., in practical cases today GPUs, FPGAs, but in principle on whatever technology) alongside traditional CPUs to improve computational performance and efficiency. Its use has become increasingly prevalent quite recently in many scientific fields especially when dealing with large datasets and complex computations in order to reduce the power consumption and/or the cost of the computing infrastructure.

Utilization of heterogeneous computing is for example gaining traction in domains like genomics [1], weather forecasts [2], medical diagnostic [3], just to cite a few clear examples. Also the industrial system has embraced utilization of heterogeneous computing, as for example in Ref. [4], often in association with Machine Learning techniques.



# Landscape document



## Best practices for heterogeneous computing

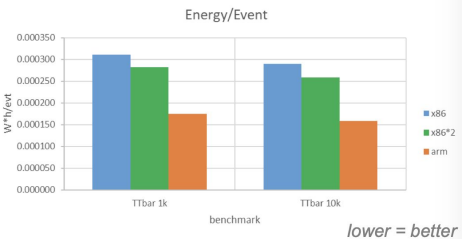
### Index:

- Introduction
- Estimates for computing needs for the current/future generation of HEP experiments
- The need for heterogeneous computing
  - Performance/cost considerations
  - Power consumption considerations
  - The utilization of HPC systems
- Solutions available/under R&D
  1. GPU
    - Portability solutions under performance studies
  2. ARM
  3. FPGA
- Link to the proposed flagships
- Bibliography

## Introduction

eous computing refers to the use of different types of hardware platforms and rs (e.g., in practical cases today GPUs, FPGAs, but in principle on whatever y) alongside traditional CPUs to improve computational performance and efficiency. s become increasingly prevalent quite recently in many scientific fields especially ling with large datasets and complex computations in order to reduce the power ion and/or the cost of the computing infrastructure.

of heterogeneous computing is for example gaining traction in domains like [1], weather forecasts [2], medical diagnostic [3], just to cite a few clear examples. ndustrial system has embraced utilization of heterogeneous computing, as for Ref. [4], often in association with Machine Learning techniques.



Algorithm	Platform	Number of Devices	Batch Size	Inf./s [Hz]	Bandwidth [Gbps]
FACILE	AWS EC2 F1	1	16,000	36 M	23
FACILE	Alveo U250	1	16,000	86 M	55
FACILE	T4 GPU	1	16,000	8 M	5.1
ResNet-50	AWS EC2 F1	8	10	1400	6.7
ResNet-50	V100 GPU	8	10	1,700	8.1
ResNet-50	ASE	1	1	460	2.2
ResNet-50	T4 GPU	1	10	250	1.2



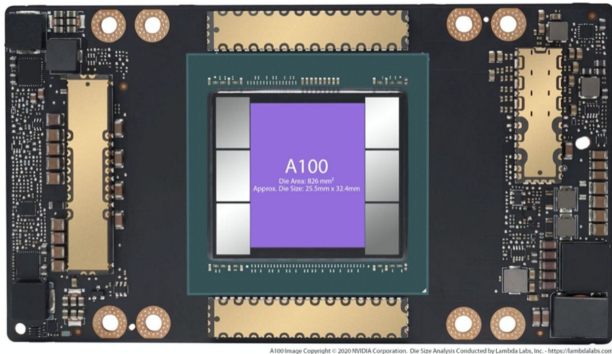
## Introductory course on GPU programming and portability tools

First course about the porting on GPUs of code and algorithms

19 Jun 2023, 11:00 → 21 Jun 2023, 20:00 Europe/Rome

Alexis Pompili (Istituto Nazionale di Fisica Nucleare), Simone Gennai (MI8)

Description



- Lecturers:
  - Andrea Bocci (CERN),
  - Felice Pantaleo (CERN),
  - Francesco Visconti (INAF).
- Tutors:
  - Lorenzo Capriotti,
  - Adriano Di Florio,
  - Tommaso Diotallevi,
  - Aurora Perego,
  - Giorgio Pizzati.
- Technical support:
  - Gioacchino Vino and Adriano Di Florio.
  - ReCas HPC for granting access to the virtual nodes and GPUs



Finanziato  
dall'Unione europea  
NextGenerationEU



Ministero  
dell'Università  
e della Ricerca



Italiadomani  
PIANO NAZIONALE  
DI RIPRESA E RESILIENZA



Centro Nazionale di Ricerca in HPC,  
Big Data and Quantum Computing



Istituto Nazionale  
di Fisica Nucleare

RECAS BARI







Finanziato  
dall'Unione europea  
NextGenerationEU



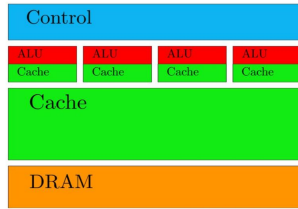
Ministero  
dell'Università  
e della Ricerca



Italiadomani  
PIANO NAZIONALE  
DI RIPRESA E RESILIENZA



Centro Nazionale di Ricerca in HPC,  
Big Data and Quantum Computing



**CPU**

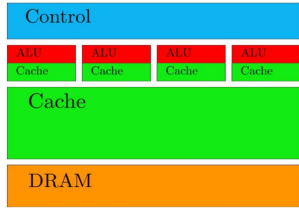


**GPU**



**RECAS BARI**





CPU



GPU

```
import cupy as cp
from cupyx.profiler import benchmark

# Define input arrays
b = cp.arange(10, 20)
c = cp.arange(20, 30)

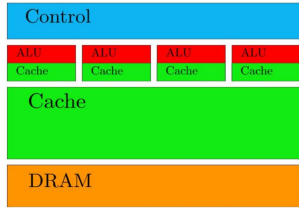
@cp.fuse(kernel_name='squared_diff')
def squared_diff(x, y):
    return (x - y) * (x - y)

print(benchmark(squared_diff, (b, c), n_repeat=10))
```



RECAS BARI





CPU



GPU

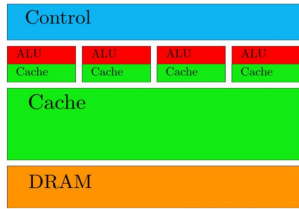
```
import cupy as cp
from cupyx.profiler import benchmark

# Define input arrays
b = cp.arange(10, 20)
c = cp.arange(20, 30)

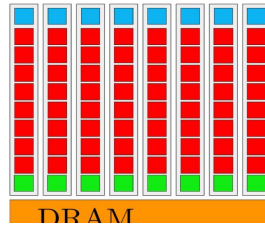
@cp.fuse(kernel_name='squared_diff')
def squared_diff(x, y):
    return (x - y) * (x - y)

print(benchmark(squared_diff, (b, c), n_repeat=10))
```





CPU



DRAM

```
import cupy as cp
from cupyx.profiler import benchmark

# Define input arrays
b = cp.arange(10, 20)
c = cp.arange(20, 30)

, n_repeat=10))
```

## Partecipazione nei vari giorni

- 75 (45) day one
- 50 (30) day two
- 35 (10) day three

Abbiamo anche lanciato una survey per capire il livello di apprezzamento del corso e suggerimenti per dove migliorare (see back up)





## Introductory course on HLS, Conifer and the BondMachine

Introductory course to HLS FPGA programming.

27–30 Nov 2023  
Europe/Rome timezone

Enter your search term



Overview

Timetable

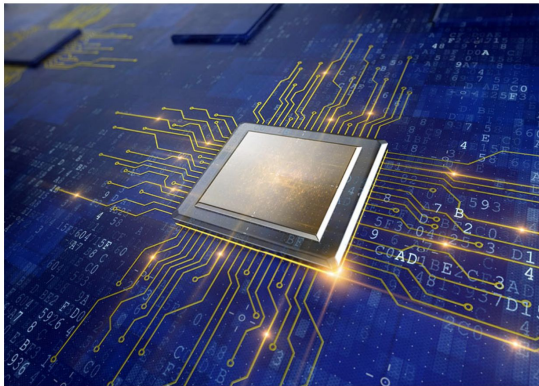
Contribution List

Registration

Participant List



Introduction to FPGA programming.



- Lecturers
  - Giovanni Petrucciani (CERN)
  - Sioni Summers (CERN)
  - Mirko Mariotti (Perugia)
- Tutors
  - Marco Lorusso
  - Giulio Bianchini
- Technical support:
  - Carmelo Pellegrino
  - Diego Michelotto
  - **ML\_INFN for providing the HW**



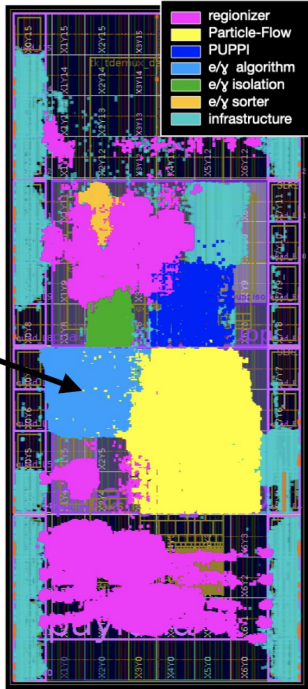
Finanziato  
dall'Unione europea  
NextGenerationEU

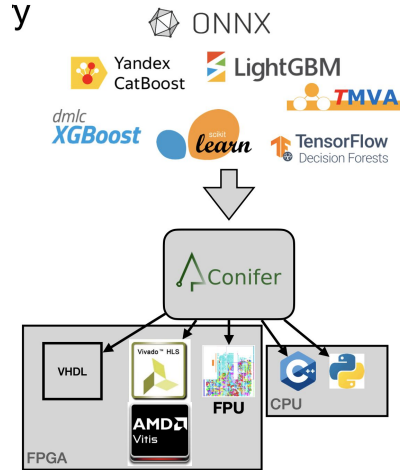
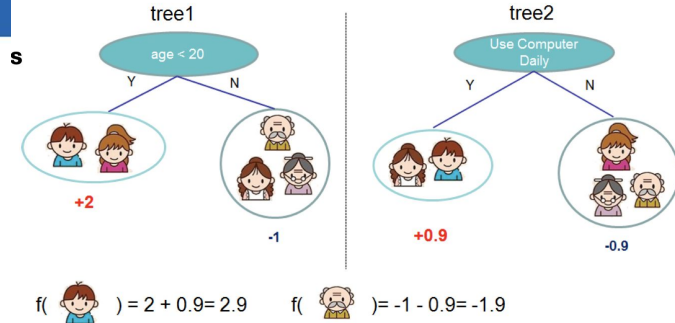
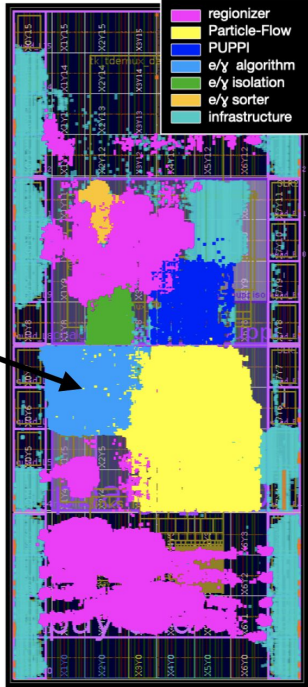


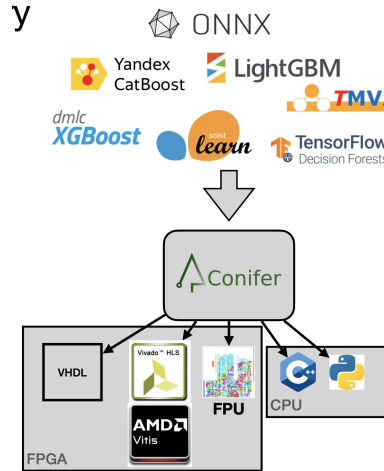
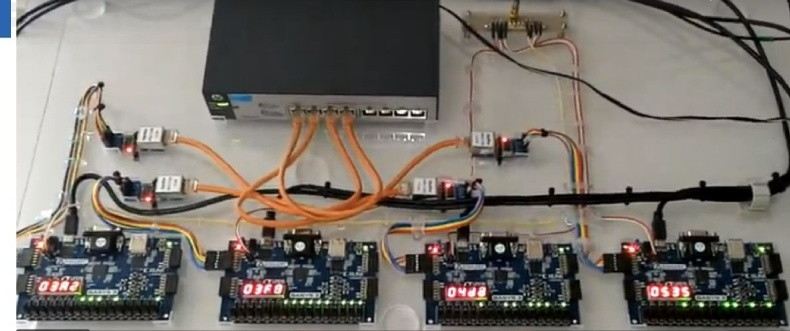
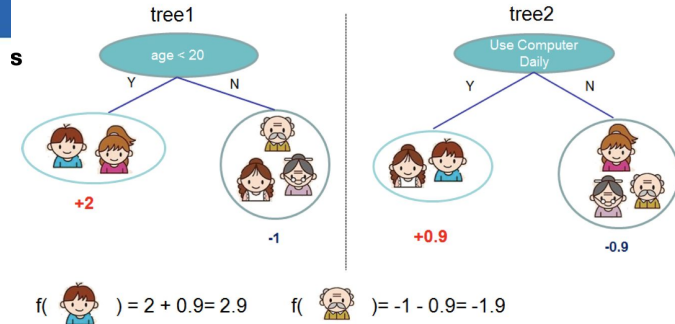
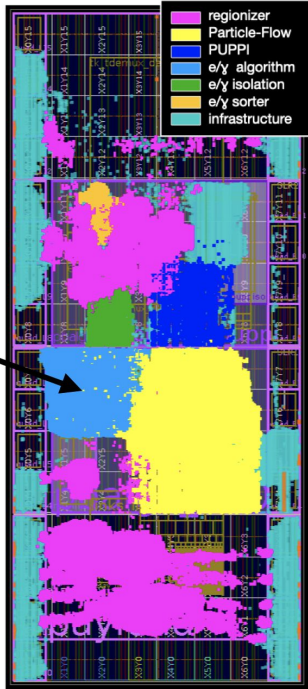
Ministero  
dell'Università  
e della Ricerca



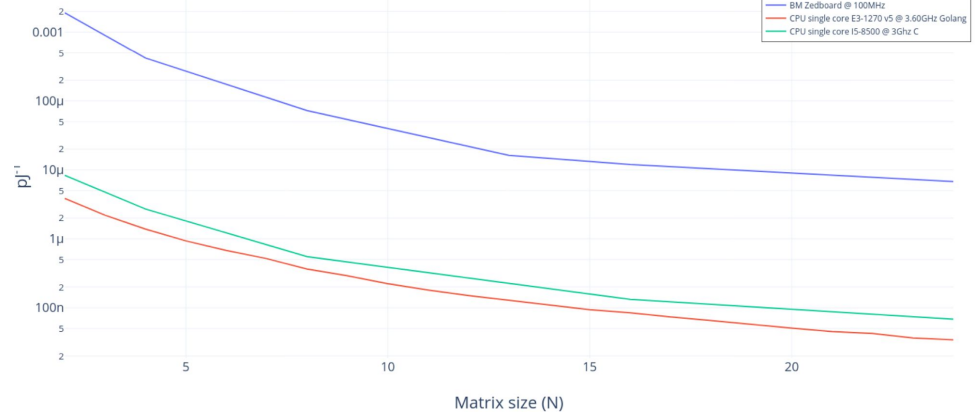
Italiadomani  
PIANO NAZIONALE  
DI RIPRESA E RESILIENZA



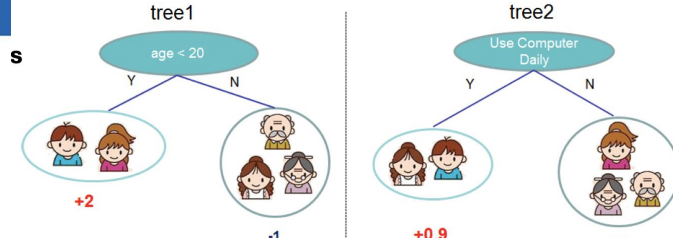
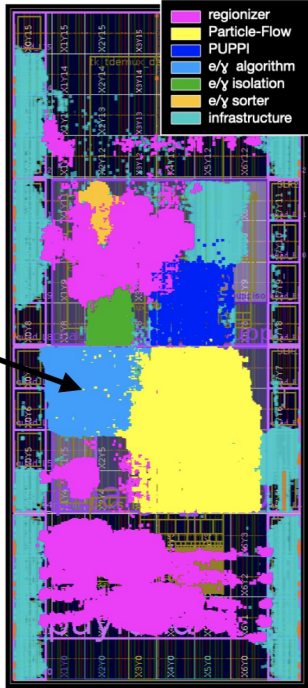




Energy efficiency



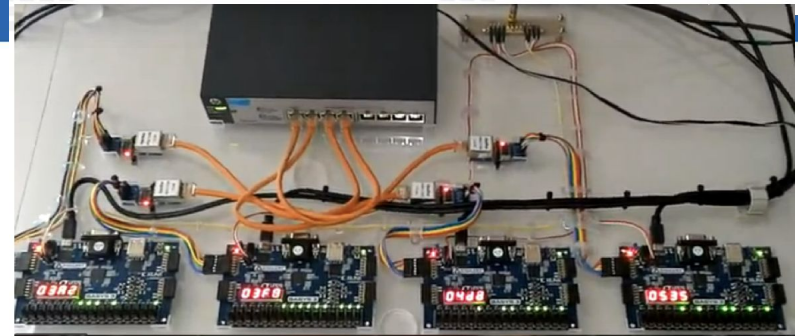
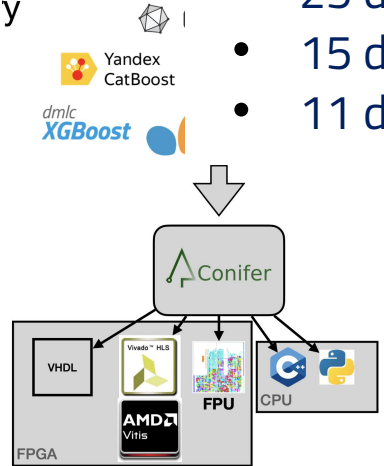




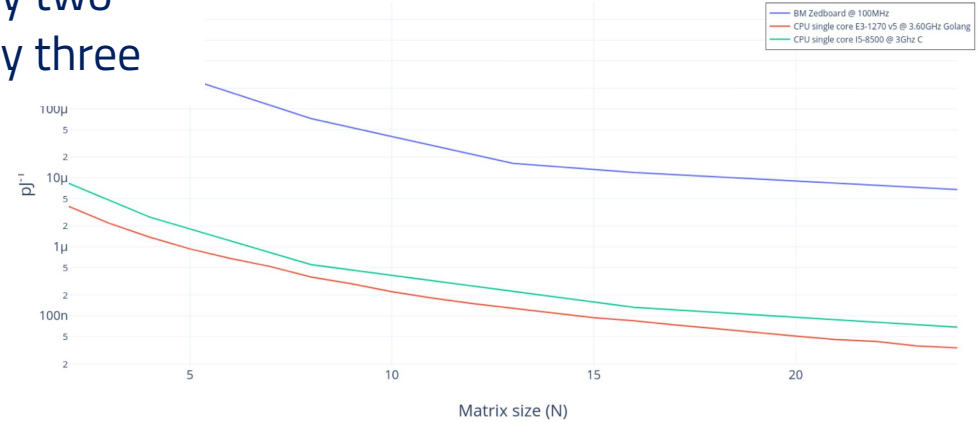
$f(\text{child}) = 2 + 0.9 = 2.9$

## Attendance

- y
- 25 day one
  - 15 day two
  - 11 day three



Energy efficiency





## Raccolta di contributi

Evento	Data	Titolo	Partecipanti
Corso / Tutorial	06/2023	<a href="#">First course about the porting on GPUs of code and algorithms (19-21 June 2023)</a>	79
Workshop	11/2023	<a href="#">From Physics to Medicine: XAI workshop</a>	101
Corso / Tutorial	11/2023	<a href="#">Introductory course on HLS and FPGA</a>	25

Titolo	Autori	Link a paper/conference
Fast Neural Network Inference on FPGAs for Triggering on Long-Lived Particles at Colliders	Andrea Coccaro Francesco Armandò Di Bello Stefano Giagu Lucrezia Rambelli and Nicola Stocchetti	<a href="https://arxiv.org/pdf/2307.05152.pdf">https://arxiv.org/pdf/2307.05152.pdf</a>
Sviluppo di acceleratori per il Machine Learning e sistemi di Inference as a Service su FPGA	Daniele Spiga Diego Ciangottini Giacomo Surace Giulio Bianchini Loriano Storchi Mirko Mariotti	<a href="#">Workshop Loano</a>
KServe inference extension for a FPGA vendor-free ecosystem	Daniele Spiga Diego Ciangottini Giacomo Surace Giulio Bianchini Loriano Storchi Mirko Mariotti	<a href="#">CHEP 2023</a>
Deep Learning techniques for reconstruction on ASTRI Mini-Array Monte Carlo data	Saverio Lombardi, Francesco Visconti, Michele Mastropietro	<a href="https://pos.sissa.it/444/713/pdf">https://pos.sissa.it/444/713/pdf</a>
A novel explainable approach in radiomics pipeline for local recurrence prediction of lung cancer: a feasibility study exploiting high energy physics potential to evaluate the model	Mariagrazia Monteleone, Simone Gennai, Pietro Govoni, Chiara Paganelli	ACM ISBN 979-8-4007-0815-2/23/09. <a href="https://doi.org/10.1145/3632047.3632074">https://doi.org/10.1145/3632047.3632074</a>
Triggerless data acquisition pipeline for Machine Learning based statistical anomaly detection	Gaia Grosso, Nicolò Lai, Matteo Migliorini, Jacopo Pazzini, A	<a href="#">CHEP 2023</a>
40MHz Triggerless Readout of the CMS Drift Tube Muon Detector	Matteo Migliorini, Jacopo Pazzini, Andrea Triossi, Marco Zan	<a href="#">TWEPP 2023</a>
Front-End RDMA Over Converged Ethernet, real-time firmware simulation	Gabriele Bortolato, Antonio Bergnoli, Damiano Bortolato, Dar	<a href="#">TWEPP 2023</a>
Front-End Rdma Over Converged Ethernet, real-time firmware simulation	Gabriele Bortolato, Antonio Bergnoli, Damiano Bortolato, Dar	<a href="#">TIPP 2023</a>



## Formazione: what's next ?

- Hackaton per gpu:
  - Con quale forma e con quali tematiche (p.es. GPU per python?) ?
  - Come possiamo coinvolgere gli altri WP ? (p.es. WP3?)
- Ottimizzazione codice gpu: profiling e best practices (in un corso avanzato)?
  - scheduler
  - profiling
  - tensorflow/pytorch
- Low level programming for FPGA
  - corso previsto per primo trimestre 2024
- Siamo in contatto con WP5 per pensare a eventi congiunti tra i 2 WP
  - Daniele ne parlera' nel suo talk
    - partendo anche dalla survey che avevamo lanciato all'inizio del progetto per vedere le competenze presenti tra le persone interessate:
    - [https://docs.google.com/forms/d/1b0iyq77\\_vnXa3HbEYthjDh4I4WNIhAuT6NuJrjLPS40/edit#responses](https://docs.google.com/forms/d/1b0iyq77_vnXa3HbEYthjDh4I4WNIhAuT6NuJrjLPS40/edit#responses)



## Food for thoughts/possible issues

- l'idea originale che i WP tecnologici funzionassero in modalita' push non ha funzionato molto bene
  - I contatti che si sono formati in realta' sono tra gruppi che si conoscevano gia'
    - gli usecase in cui abbiamo sinergie sono principalmente con WP2
- poi c'e' la mancanza di expertise per porting su gpu di codice non C++
  - Qui c'e' da intensificare i contatti con WP3 per capire l'expertise che sia già presente nel WP (a parte Francesco Visconti ...)
- Di conseguenza l'attendance ai nostri meeting e' abbastanza limitata
  - Forse dovremmo organizzare meeting condivisi con WP1 e WP3 in cui si discute la parte tecnologica dei flagship?
- ultimo punto: come si accedono le FPGA al CNAF?
  - a chi si chiede?



## Additional considerations

### GPU

- Ci sono già in corso diverse attività sia teoriche che sperimentali, al momento l'interesse maggiore e' lo sviluppo di framework platform agnostic
  - qui la maggiore esperienza viene dalla parte sperimentale, dove l'effort in questa direzione e' iniziato da tempo
  - serve rafforzare l'expertise italiana e fare formazione diffusa sul porting di algoritmi
  - al di la' del porting degli algoritmi c'e' poi da capire il fine tuning di diversi parametri per ottimizzare la memoria e l'esecuzione ovviamente questo puo' essere molto project dependent

### FPGA

- Area con interesse crescente negli ultimi tempi, e **molto varia per quanto riguarda progetti e anche tipologia di HW**, attività principalmente lato sperimentale con forte liason a tematiche legate a trigger selection e reconstruction
- Probabilmente l'area dove c'e' meno expertise e quindi anche dove serve maggiormente la formazione di nuovo person-powe
  - c'e' una conoscenza diffusa dei vari codici per la scrittura di firmware, manca forse la conoscenza per il setup di **cluster di FPGA** su media-grande scala (forse più indicato per WP5)
  - esistono vari testbed già funzionanti per alcuni progetti ma su scala ridotta, in altri casi si usano risorse messe a disposizione dai vari esperimenti (principalmente al CERN).



## Conclusions

- organizzati primi eventi formativi con un buon riscontro (con impostazione piu' vicino all'ambito WP2 e un po' WP1 e WP3)
- bisogna insistere in questa direzione cercando di coinvolgere maggiormente WP1 e WP3
- adesso che il lavoro nell'ambito delle varie flagship e' iniziato vorremmo focalizzare nelle discussioni la parte piu' squisitamente tecnologica
  - (capire esigenze, individuare bottleneck, condividere strategie,...)
  - sia in meeting WP4 dedicati sia/o in meeting co-organizzati con gli spoke WP1-3



Finanziato  
dall'Unione europea  
NextGenerationEU



Ministero  
dell'Università  
e della Ricerca



Italiadomani  
PIANO NAZIONALE  
DI RIPRESA E RESILIENZA



## Back-up



## Overall satisfaction (19 answers)

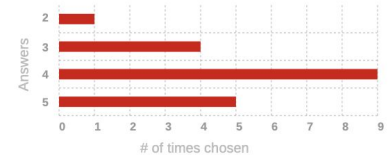
All the scores are between 3.5 and 4.0 out of 5

**How much would you rate the C/C++ to CUDA part of the course?** **Answered:** 19 **Please choose a number from 1 to 5, 1 means you really did not find the course useful and/or interesting while 5 means you find it very useful and interesting**

Average: 3.95

Min: 2

Max: 5

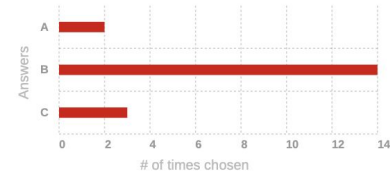


**Do you think the C/C++ to CUDA part was easy to follow , considering it being an introductory one?** **Answered:** 19

A. Not easy at all: 2 (10.53%)

B. Enough easy for not experienced people: 14 (73.68%)

C. Not easy for me but i made it; likely easy for people with more experience than me: 3 (15.79%)







## Overall satisfaction (19 answers)

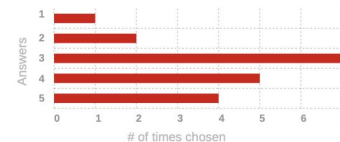
All the scores are between 3.5 and 4.0 out of 5

How much would you rate the CUDA-Python part of the course? **Answered: 19** Please choose a number from 1 to 5, 1 means you really did not find the course useful and/or interesting while 5 means you find it very useful and interesting

Average: 3.47

Min: 1

Max: 5

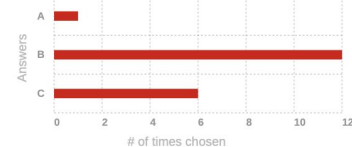


Do you think the python-cuda part was easy to follow, considering it being an introductory one? **Answered: 19**

A. Not easy at all: 1 (5.26%)

B. Enough easy for not experienced people: 12 (63.16%)

C. Not easy for me but i made it; likely easy for people with more experience than me: 6 (31.58%)





## Overall satisfaction (19 answers)

All the scores are between 3.5 and 4.0 out of 5

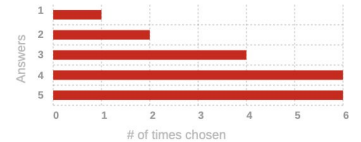
How much would you rate the introduction to Alpaka part of the course?

**Answered:** Please choose a number from 1 to 5, 1 means you really did not find the course useful and/or interesting while 5 means you find it very useful and interesting  
19

Average: 3.74

Min: 1

Max: 5



Do you think the Alpaka part was easy to follow, considering it being an introductory one?

Answered : 19

A. Not easy at all: 5 (26.32%)

B. Enough easy for not experienced people: 7 (36.84%)

C. Not easy for me but i made it; likely easy for people with more experience than me: 7 (36.84%)

