

Towards EXtreme scale Technologies and Accelerators for euROhpc hw/Sw Supercomputing Applications for exascale



textarossa

WP6 Applications and Use Cases

D6.3 Final Assessment and Guidelines



This project has received funding from the European Union's Horizon 2020 research and innovation programme, EuroHPC JU, grant agreement No 056921



textarossa

TEXTAROSSA

Towards EXtreme scale Technologies and Accelerators for euROhpc
hw/Sw Supercomputing Applications for exascale
Grant Agreement No.: 956831

Deliverable: D6.3 Final Assessment and Guidelines

Project Start Date: 01/04/2021

Duration: 36 months

Coordinator: AGENZIA NAZIONALE PER LE NUOVE TECNOLOGIE, L'ENERGIA E LO SVILUPPO ECONOMICO SOSTENIBILE - ENEA, Italy

Deliverable No	D6.3
WP No:	WP6
WP Leader:	PSNC
Due date:	M36 (March 31, 2024)
Delivery date:	5/06/2024

Dissemination

Level:

PU	Public	X
PP	Restricted to other programme participants (including the Commission Services)	
RE	Restricted to a group specified by the consortium (including the Commission Services)	
CO	Confidential, only for members of the consortium (including the Commission Services)	



This project has received funding from the European Union's Horizon 2020 research and innovation programme, EuroHPC JU, grant agreement No 956831



DOCUMENT SUMMARY INFORMATION

Project title:	Towards EXtreme scale Technologies and Accelerators for euROhpc hw/Sw Supercomputing Applications for exascale
Short project name:	TEXTAROSSA
Project No:	956831
Call Identifier:	H2020-JTI-EuroHPC-2019-1
Unit:	EuroHPC
Type of Action:	EuroHPC - Research and Innovation Action (RIA)
Start date of the project:	01/04/2021
Duration of the project:	36 months
Project website:	TEXTAROSSA.eu

WP6 Applications and Use Cases

Deliverable number:	D6.3					
Deliverable title:	Final Assessment and Guidelines					
Due date:	M36					
Actual submission date:	05/06/2024					
Editors:	Pasqua D’Ambra, Michal Kulczewski, Alessandro Lonardo					
Authors:	Massimo Bernaschi, Mauro Carrozzo, Alessandro Celestini, Pasqua D’Ambra, Daniel Jaschke, Paolo Palazzari, Michal Kulczewski, Ariel Oleksiak, Miłosz Ciżnicki, Wojciech Szeliga, Daniel Jaschke, Alessandro Lonardo, Luca Pontisso, Cristian Rossi, Francesco Simula, Martin Kuehn					
Work package:	WP6					
Dissemination Level:	Public					
No. pages:	150					
Authorized (date):	05/06/2024					
Responsible person:	Pasqua D’Ambra					
Status:	Plan	Draft	Working	Final	Submitted	Approved

Revision history:

Version	Date	Author	Comment
0.1	2023-10-27	Pasqua D’Ambra	Draft structure
0.2	2024-02-15	Pasqua D’Ambra	CNR section
0.3	2024-03-04	Sergio Saponara	UNIPI (CINI)
0.4	2024-03-15	Carlos Álvarez	RAIDER OmpSs@FPGA results
0.5	2024-03-15	Antonio Filgueras	OmpSs@FPGA IDV-E results
0.6	2024-03-15	All	Remaining contributions
0.7	2024-03-18	Michał Kulczewski	Ready for first internal review

0.8	2024-03-28	All	Final remaining contributions, addressing internal reviewers comments
0.9	2024-03-29	Pasqua D'Ambra	Ready for submission
1.0	2024-05-01	Massimo Celino	Submitted
2.0	2024-04-26	All	Enhanced version ready for final submission
3.0	2024-06-05	All	Revised version

Quality Control:

Checking process	Who	Date
Checked by internal reviewer 1	Alessandro Lonardo	2024-06-04
Checked by internal reviewer 2	Bérenger Bramas	2024-06-04
Checked by Task Leader	Pasqua D'Ambra	2024-06-05
Checked by WP Leader	Michał Kulczewski	2024-06-05
Checked by Project Coordinator	Massimo Celino	2024-06-05

COPYRIGHT

© Copyright by the **TEXTAROSSA** consortium, 2021-2024

This document contains material, which is the copyright of TEXTAROSSA consortium members and the European Commission, and may not be reproduced or copied without permission, except as mandated by the European Commission Grant Agreement No. 956831 for reviewing and dissemination purposes.

ACKNOWLEDGEMENTS

This project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement no 956831. The JU receives support from the European Union's Horizon 2020 research and innovation programme and Italy, Germany, France, Spain, Poland.

Please see <http://TEXTAROSSA.eu> for more information on the TEXTAROSSA project.

The partners in the project are AGENZIA NAZIONALE PER LE NUOVE TECNOLOGIE, L'ENERGIA E LO SVILUPPO ECONOMICO SOSTENIBILE (ENEA), FRAUNHOFER GESELLSCHAFT ZUR FOERDERUNG DER ANGEWANDTEN FORSCHUNG E.V. (FHG), CONSORZIO INTERUNIVERSITARIO NAZIONALE PER L'INFORMATICA (CINI), INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET AUTOMATIQUE (INRIA), BULL SAS (BULL), E4 COMPUTER ENGINEERING SPA (E4), BARCELONA SUPERCOMPUTING CENTER-CENTRO NACIONAL DE SUPERCOMPUTACION (BSC), INSTYTUT CHEMII BIOORGANICZNEJ POLSKIEJ AKADEMII NAUK (PSNC), ISTITUTO NAZIONALE DI FISICA NUCLEARE (INFN), CONSIGLIO NAZIONALE DELLE RICERCHE (CNR), IN QUATTRO SRL (in4). Linked third parties of CINI are POLITECNICO DI MILANO (CINI-POLIMI), Università di Torino (CINI-UNITO) and Università di Pisa (CINI-UNIP); linked third party of INRIA is Université de Bordeaux; in-kind third party of ENEA is Consorzio CINECA (CINECA); in-kind third party of BSC is Universitat Politècnica de Catalunya (UPC).

The content of this document is the result of extensive discussions within the TEXTAROSSA © Consortium as a whole.

DISCLAIMER

The content of the publication herein is the sole responsibility of the publishers and it does not necessarily represent the views expressed by the European Commission or its services.

The information contained in this document is provided by the copyright holders "as is" and any express or implied warranties, including, but not limited to, the implied warranties of merchantability and fitness for a particular purpose are disclaimed. In no event shall the members of the TEXTAROSSA collaboration, including the copyright holders, or the European Commission be liable for any direct, indirect, incidental, special, exemplary, or consequential damages (including, but not limited to, procurement of substitute goods or services; loss of use, data, or profits; or business interruption) however caused and on any theory of liability, whether in contract, strict liability, or tort (including negligence or otherwise) arising in any way out of the use of the information contained in this document, even if advised of the possibility of such damage.

Table of contents

Executive Summary 12

1 Introduction..... 15

2 Methodology 17

 2.1 Hardware 17

 2.2 Application tuning, performance and energy measurements 18

3 Results..... 19

 3.1 Smart cities – CINI-UNIFI..... 19

 3.2 MathLib – CNR 23

 3.3 RTM – FRAUNHOFER..... 33

 3.4 HEP - INFN..... 41

 3.5 NEST-GPU – INFN 56

 3.6 RAIDER – INFN..... 66

 3.7 TNM – INFN..... 83

 3.8 ScalFMM (Mathlibs-INRIA) 90

 3.9 Chameleon (Mathlibs-INRIA) 91

 3.10 StarONNX (INRIA)..... 94

 3.11 UrbanAir - PSNC 96

 3.12 FIPLib: FPGA Image Processing Library – ENEA/INFN 109

 3.13 NBody - BSC..... 119

4 Concluding Remarks..... 134

 4.1 Smart Cities - CINI-UNIFI 134

 4.2 MathLib - CNR..... 134

 4.3 RTM - Fraunhofer 135

 4.4 HEP - INFN..... 135

 4.5 NEST-GPU - INFN 136

 4.6 RAIDER - INFN 136

 4.7 TNM - INFN..... 137

 4.8 MathLib - INRIA 137

 4.9 UrbanAir - PSNC..... 138

 4.10 FIPLib - ENEA 139

 4.11 NBody - BSC..... 139

5 Achieved project objectives, lesson learned and guidelines 140

 5.1 Project objectives..... 140

 5.2 Lesson learned and guidelines 147

6 References..... 150

List of Figures

Figure 1: TEXTAROSSA Integrated Platform.....	16
Figure 2: [SmartCities] Computation flow for the YOLOv5 step	20
Figure 3: [SmartCities] Power efficiency measured as FPS/Watt	22
Figure 4: [MathLib CNR] Power behavior of BCMG on IDV-A before the 2-phase cooling system installation.....	28
Figure 5: [MathLib CNR] Power behavior of BCMG on IDV-A after the 2-phase cooling system installation.....	28
Figure 6: [MathLib CNR] Weak scalability on Leonardo up to 64 billion dofs, number of iterations.....	30
Figure 7: [MathLib CNR] Weak scalability on Leonardo up to 64 billion dofs, solve time ...	30
Figure 8: [MathLib CNR] Weak scalability on Leonardo up to 64 billion dofs, execution time	31
Figure 9: [MathLib CNR] Weak scalability on Leonardo up to 64 billion dofs, speedup	31
Figure 10: [MathLib CNR] Weak scalability on Leonardo. BCMGX vs AMGX, number of iterations.....	32
Figure 11: [MathLib CNR] Weak scalability on Leonardo. BCMGX vs AMGX, solve time	33
Figure 12: [MathLib CNR] Weak scalability on Leonardo. BCMGX vs AMGX, time per iteration	33
Figure 13: [RTM] Marmousi results with different combinations of Floating Point formats .	35
Figure 14: [RTM] Marmousi results with different combinations of Floating Point formats continued.....	36
Figure 15: [RTM] Distribution of absolute values in velocity model (left) and all the receiver traces (right).....	36
Figure 16: [RTM] Distribution of absolute values in velocity model (left) and all the receiver traces (right).....	37
Figure 17: [RTM] Distribution of absolute values in velocity model (left) and all the receiver traces (right).....	38
Figure 18: [RTM] Marmousi results with different combinations of Floating Point formats .	40
Figure 19: [RTM] Infinity norm of difference between Float16 and Float32 as percentage of the infinity norm of the Float16 example	41
Figure 20: [HEP] Throughput vs cores for CLUE on Dibona CPU	43
Figure 21: [HEP] Throughput vs cores for Pixeltrack on Dibona CPU	44
Figure 22: [HEP] Energy efficiency vs cores for CLUE on Dibona CPU.....	44
Figure 23: [HEP] Energy efficiency vs cores for Pixeltrack on Dibona CPU.....	45
Figure 24: [HEP] Throughput (a.u.) vs number of used GPUs for CLUE on Dibona node....	46
Figure 25: [HEP] Throughput (a.u.) vs number of used GPUs for Pixeltrack on Dibona node	47
Figure 26: [HEP] Energy efficiency (normalized to one GPU) vs number of used GPUs for CLUE on Dibona node.....	47
Figure 27: [HEP] Energy efficiency (normalized to one GPU) vs number of used GPUs for Pixeltrack on Dibona node.....	48
Figure 28: [HEP] Throughput vs number of used cores for CLUE on IDV-A CPU.....	49
Figure 29: [HEP] Throughput vs number of used cores for Pixeltrack on IDV-A CPU.....	50
Figure 30: [HEP] Energy efficiency vs number of used cores for CLUE on IDV-A CPU	50
Figure 31: [HEP] Energy efficiency vs number of used cores for Pixeltrack on IDV-A CPU	51
Figure 32: [HEP] Throughput vs number of used GPUs for CLUE on IDV-A node.....	52
Figure 33: [HEP] Throughput vs number of used GPUs for Pixeltrack on IDV-A node.....	52

Figure 34: [HEP] Energy efficiency vs number of used GPUs for CLUE on IDV-A node53

Figure 35: [HEP] Energy efficiency vs number of used GPUs for Pixeltrack on IDV-A node53

Figure 36: [HEP] Throughput vs number of used cores for Pixeltrack on IDV-E CPU55

Figure 37: [HEP] Throughput vs number of used cores for CLUE on IDV-E CPU55

Figure 38: [HEP] Energy efficiency vs cores for Pixeltrack on IDV-E CPU.....55

Figure 39: [HEP] Energy efficiency vs cores for CLUE on IDV-E CPU56

Figure 40: [NEST-GPU] Example of idle GPU power envelope57

Figure 41: [NEST-GPU] Single active GPU power reading on IDV-A.....57

Figure 42: [NEST-GPU] 2 active GPUs power reading on IDV-A.....58

Figure 43: [NEST-GPU] 3 active GPUs power reading on IDV-A.....59

Figure 44: [NEST-GPU] 4 active GPUs power reading on IDV-A.....59

Figure 45: [NEST-GPU] Single active GPU power reading on Dibona60

Figure 46: [NEST-GPU] 2 active GPUs power reading on Dibona61

Figure 47: [NEST-GPU] 3 active GPUs power reading on Dibona61

Figure 48: [NEST-GPU] 4 active GPUs power reading on Dibona62

Figure 49: [NEST-GPU] IDV-E power reading62

Figure 50: [NEST-GPU] IDV-E PDU reading63

Figure 51: [NEST-GPU] CPU-only run power reading on IDV-A – 1 GPU size.....64

Figure 52: [NEST-GPU] CPU-only run power reading on IDV-A – 2 GPUs size65

Figure 53: [NEST-GPU] CPU-only run power reading on IDV-A – 3 GPUs size65

Figure 54: [NEST-GPU] CPU-only run power reading on IDV-A – 4 GPUs size66

Figure 55: [RAIDER] The workflow for the generation of CNN kernels in RAIDER.....67

Figure 56: [RAIDER] Example of input images for the CNN67

Figure 57: [RAIDER] Test setup on the Xilinx® Alveo U200 installed in the INFN Roma 1 APE Lab.....69

Figure 58: [RAIDER] Vitis synthesis report of krnl_sender HLS kernel.....69

Figure 59: [RAIDER] Application CPU hosts power profiles (4 Intel Sapphire Rapid).....70

Figure 60: [RAIDER] FPGA power profiles (2 Alveo Xilinx U200 setup).....71

Figure 61: [RAIDER] FPGA power profiles (3 Alveo Xilinx U200 setup).....71

Figure 62: [RAIDER] FPGA power profiles (4 Alveo Xilinx U200 setup).....71

Figure 63: [RAIDER] Throughput scaling with the number of deployed CNN HLS kernels 72

Figure 64: [RAIDER] Vitis synthesis report of hardware project to be deployed on Computing node.....73

Figure 65: [RAIDER] Test setup on the Xilinx® Alveo U280 installed in the IDV-E testbed (single node).....73

Figure 66: [RAIDER] Test setup on the Xilinx® Alveo U280 installed in the IDV-E testbed (double node).....74

Figure 67: [RAIDER] CPU hosts power profiles (2 Ampere Altra Max processor).....75

Figure 68: [RAIDER] FPGA power profiles (2 Alveo Xilinx U280 setup).....75

Figure 69: [RAIDER] FPGA power profiles (3 Alveo Xilinx U280 setup).....76

Figure 70: [RAIDER] FPGA power profiles (4 Alveo Xilinx U280 setup).....76

Figure 71:[RAIDER] Throughput scaling trends80

Figure 72: [RAIDER] Energy Efficiency scaling trends80

Figure 73: [RAIDER] Throughput with different OmpSs@FPGA implementations81

Figure 74: [RAIDER] Different OmpSs@FPGA implementations organization83

Figure 75: [TNM] Execution time vs error varying the number of final sweeps in double precision.....84

Figure 76: [TNM] Energy consumption vs error varying the number of final sweeps in double precision.....85

Figure 77: [TNM] Execution time vs system size87

Figure 78: [TNM] Execution time vs number of sweeps in single precision87

Figure 79: [TNM] Ground energy state (lower is better) vs number of single precision sweeps87

Figure 80: [TNM] Speedup after algorithmic improvements88

Figure 81: [TNM] Improvement ratio in energy efficiency after algorithmic improvements .88

Figure 82: [TNM] Relative CPU time w.r.t icTPO89

Figure 83: [FMM] Performance results for TB-FMM on two hardware configurations.....91

Figure 84: [CHAMELEON] Performance results for Chameleon on two hardware configurations93

Figure 85: [CHAMELEON] Execution trace (emulation) for a Cholesky factorization without using user’s priorities.....93

Figure 86: [CHAMELEON] Execution trace (emulation) for a Cholesky with user’s priorities94

Figure 87: [StarONNX] Performance plots for StarONNX with the GoogLeNet network95

Figure 88: [StarONNX] Performance plots for StarONNX with the GoogLeNet network showing latency as a function of throughput: for each throughput, the batch size that minimizes maximal inference time for a query is chosen.....96

Figure 89: [UrbanAir] Iterations per second (problem size: 6.5M).....98

Figure 90: [UrbanAir] Iterations per second (5M per GPU)98

Figure 91: [UrbanAir] Speedup comparison for different domain sizes99

Figure 92: [UrbanAir] Efficiency (weak speedup) comparison for different domain sizes ..100

Figure 93: [UrbanAir] Time-to-solution on GPUs: V100 vs. H100101

Figure 94: [UrbanAir] Strong speedup comparison for 59M grid points102

Figure 95: [UrbanAir] Iterations per second for 59M grid points102

Figure 96: [UrbanAir] Energy usage characteristic for 59M grid points.....103

Figure 97: [UrbanAir] Energy consumption for 59M grid points, only used GPUs accounted104

Figure 98: [UrbanAir] Energy consumption for 59M grid points, unused GPUs accounted 104

Figure 99: [UrbanAir] Time-to-solution and energy consumption for 59M grid points per GPU.....105

Figure 100: [UrbanAir] Iterations / kW105

Figure 101:[FIBLib] Codelet and graphics showing a simple image processing kernel110

Figure 102: [FIPLib] Basic image processing pipeline used in ALGO1110

Figure 103: [FIPLib] Kernels in ALGO1111

Figure 104: [FIPLib] Power absorption of ALGO1113

Figure 105: [FIPLib] APEIRON project for single FPGA bistream115

Figure 106: [FIPLib] multi-FPGA deployment and execution scheme.....116

Figure 107: [FIPLib] Kernel displacement in multiple FPGA setup. Internal datapath set at 32 bytes116

Figure 108: [FPLIB] Multi FPGA APEIRON Image Processing on 512x512 images execution power profile117

Figure 109: [FPLIB] Multi FPGA APEIRON Image Processing on 4096x4096 images execution power profile117

Figure 110: [FIPLib] CPU, Single- and Multi-FPGA implementation performances comparison.....118

Figure 111: [NBody] Single node n-body FPGA design.....119

Figure 112: [NBody] Diagram sequence of the n-body computation for single FPGA121

Figure 113: [NBody] Temperature and power for an actively air-cooled Alveo U200.....122

Figure 114: [NBody] Temperature and power for a passively air-cooled Alveo U55c.....122

Figure 115: [NBody] Temperature and power air cooled IDV-E's Alveo U280..... 123
 Figure 116: [NBody] Temperature and power 2-phase cooled IDV-E's Alveo U280..... 123
 Figure 117: [NBody] Sequence diagram of multi-fpga n-body execution 125
 Figure 118: [NBody] Sequence diagram of multi-fpga n-body execution 126
 Figure 119: [NBody] Sequence diagram of multi-fpga n-body execution 127
 Figure 120: [NBody] Performance comparison between using 1 and 2 FPGA..... 128
 Figure 121: [NBody] body design diagram showing ompss accelerators and apeiron communication modules 129
 Figure 122: [NBody] Ompss + apeiron sequence diagram..... 130
 Figure 123: [NBody] Performance for different n-body FPGA implementations..... 130
 Figure 124: [NBody] Energy efficiency of n-body implementations..... 131
 Figure 125: [NBody] Measured performance and ideal across different number of devices 132
 Figure 126: [NBody] MareNostrum 4 n-body scalability..... 132

List of Tables

Table 1: [SmartCities] Implementation results and benchmarks for different platforms 21
 Table 2: [MathLib CNR] KPIs 26
 Table 3: [MathLib CNR] Performance and Energy Consumption of the SpMV kernel 26
 Table 4: [MathLib CNR] Performance and Energy Consumption of the SpMM kernel..... 27
 Table 5: [MathLib CNR] Performance and Energy Consumption of the MWM kernel 27
 Table 6: [MathLib CNR] Performance and Energy Consumption of the BCMG solver before installation of the 2-phase cooling system 28
 Table 7: [MathLib CNR] Performance and Energy Consumption of the BCMG solver after installation of the 2-phase cooling system 29
 Table 8: [HEP] KPIs for CLUE on Dibona CPU 42
 Table 9: [HEP] KPIs for Pixeltrack on Dibona CPU..... 43
 Table 10: [HEP] KPIs for CLUE on GPU on Dibona node 46
 Table 11: [HEP] KPIs for Pixeltrack on GPU on Dibona node..... 46
 Table 12 [HEP] KPIs for CLUE on CPU for IDV-A node..... 49
 Table 13: [HEP] KPIs for Pixeltrack on CPU for IDV-A node..... 49
 Table 14: [HEP] KPIs for CLUE on GPU on IDV-A node 51
 Table 15: [HEP] KPIs for Pixeltrack on GPU on IDV-A node 51
 Table 16: [HEP] KPIs for Pixeltrack on CPU on IDV-E node..... 54
 Table 17: [HEP] Throughput vs number of used cores for CLUE on IDV-E CPU..... 54
 Table 18: NEST-GPU comparison between A100 and H100 GPUs..... 60
 Table 19: [NEST-GPU] neural simulation on IDV-A CPU-only – runtimes and power 64
 Table 20: [NEST-GPU] KPIs for the neural network simulation application 66
 Table 21: [RAIDER] Processing Throughput with an increasing number of Computing nodes (and CNN HLS kernels)..... 72
 Table 22: [RAIDER] processing time per event with an increasing number of Computing FPGAs (and CNN HLS kernels)..... 74
 Table 23: [RAIDER] Baseline KPIs evaluated on the execution of the RAIDER application on a single Xilinx Alveo U200 FPGA, processing 2.7M events with a global clock of 100 MHz 77
 Table 24: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a two Xilinx Alveo U200 FPGA system, processing 2.7M events with a global clock of 200 MHz 77
 Table 25: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a three Xilinx Alveo U200 FPGA system, processing 2.7M events with a global clock of 200 MHz 78

Table 26: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a four Xilinx Alveo U200 FPGA system, processing 2.7M events with a global clock of 200 MHz 78

Table 27: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a two Xilinx Alveo U280 FPGA IDV-E node, processing 2.7M events with a global clock of 200 MHz 78

Table 28: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a three Xilinx Alveo U280 FPGA configuration distributed over 2 IDV-E nodes, processing 2.7M events with a global clock of 200 MHz 79

Table 29: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a four Xilinx Alveo U280 FPGA configuration distributed over 2 IDV-E nodes, processing 2.7M events with a global clock of 200 MHz 79

Table 30: [RAIDER] Improvement factors of the KPIs for the tested designs over the baseline.81

Table 31: [TNM] Hybrid simulation with MPI and threading via MKL library. We use simulation of the quantum green tea library out of the Quantum TEA suite.....86

Table 32: [TNM] Measurement of key performance indicators for a single-core simulation with Quantum matcha TEA on IDV-A. The simulations repeat the analysis executed for D6.2 on the most recent machine, i.e., IDV-A.90

Table 33: [UrbanAir] KPIs overview 100

Table 34: [UrbanAir] KPI Time-to-solution (strong) 106

Table 35: [UrbanAir] KPI Iterations / s (strong) 106

Table 36: [UrbanAir] KPI Iterations / kW (strong, used GPUs) 106

Table 37: [UrbanAir] KPI iterations /kW (all GPUs, strong) 107

Table 38: [UrbanAir] KPI time-to-solution (weak) 107

Table 39: [UrbanAir] KPI Iterations / s (weak) 107

Table 40: [UrbanAir] KPI Iterations / kW (used GPUs, weak) 107

Table 41: [UrbanAir] KPI Iterations / kW (all GPUs, weak) 107

Table 42: [FIPLib] Resources used by ALGO1 112

Table 43: [FIPLib] Energy and time used for ALGO1 113

Table 44: [FIPLib] Multi-FPGA APEIRON Implementation performances 117

Table 45: [NBody] Resource usage of the single-node n-body for the IDV-E alveo U280.. 120

Table 46: [NBody] Comparison of performance, power and power efficiency across different devices..... 124

Table 47: [NBody] Comparison of power consumption on air cooled and 2-phase cooled IDV-E..... 125

Table 48: [NBody] Summary of performance, power consumption and power efficiency across different clusters..... 133

Executive Summary

This deliverable summarizes possible adaptation of applications and/or TEXTAROSSA toolchain features which was needed to obtain benefits in terms of performance and energy efficiency, as evaluated by the KPIs defined for each use case in the previous D6.1 deliverable. Therefore, as follow up of all the WP6 activities on the applications, it provides guidelines, possible recommendations and final evaluation of the project efforts. All use cases applied common methodology for performance and energy measurements, discussed in detail in D1.4. The tests were conducted among others on TEXTAROSSA available testbeds, including IDV-A (node with NVIDIA GPU accelerators from Atos) and IDV-E (node with FPGA accelerators from E4). The applications represent a wide range of application types (HPC, AI, HDPA), scientific domains, approaches to parallelization (task-based, streaming-based, distributed/shared memory). In particular, the use cases benefit from developments around one of the three main pillars defined: heterogeneous resources, mixed precision and dynamic runtime systems.

The work carried out corresponds directly to the following overall project objectives:

- Energy efficiency, by the application developments (also as benchmarks to assess the node energy efficiency)
- Sustained application performance, by the application developments;
- Seamless integration of reconfigurable accelerators, by using the APEIRON framework and the BSC Task Scheduler for OmpSS@FPGA;
- Development of new IPs, as the INFN intra/inter-FPGA communication IP behind the APEIRON framework BSC Task Scheduler;
- Integrated Development Platform, by using existing IDV-A and IDV-E.

CNR has developed computational kernels for a TEXTAROSSA MathLib specifically required in sparse matrix computations and iterative linear solvers, which are widely exploited in Scientific Computing and Data Analysis. The library was thought both as main component of the TEXTAROSSA platform and as benchmark tool for IDV-A. Focus was on GPU-kernels efficiency and on scalability when multiple GPUs, also on different computing nodes, are needed for computations because dimensions largely exceed the memory resources of a single GPU. Therefore, on IDV-A, the library stresses the GPU operation capabilities and memory/communication channels bandwidth at the node level. Performance analysis and GPU energy consumption KPIs on IDV-A are reported. Furthermore, scalability results and comparisons with the state-of-the-art Nvidia library on the Leonardo Italian supercomputer are discussed. Results demonstrate benefits of the new algorithms, specifically thought for GPU-accelerated architectures, and of the implementation design patterns, which efficiently exploit high throughput of GPUs and reduce/hide data communication among memories and computing nodes, with respect to Nvidia library implementing same computations. Power and energy consumption results of the kernels on IDV-A are discussed before and after the installation of the new 2-phase cooling system on the IDV-A platform.

INFN provided one AI/High Performance Data Analytics (RAIDER) and three HPC (NEST-GPU, HEP, TNM) applications as benchmarks to drive the co-design and characterization activities of project IDVs.

The RAIDER (Real-time AI-based Data analytics on hEteRogeneous distributed systems) application, described in section 3.6, is driven by the use case of a real-time particle identification system for the CERN NA62 High Energy Physics experiment. Being entirely

FPGA-based, it represents the TEXTAROSSA IDV-E reference application for its characterization in terms of processing throughput and energy efficiency KPIs. RAIDER was co-designed and co-developed with the APEIRON HLS streaming programming framework aimed at the seamless scalable integration of reconfigurable accelerators. The INFN intra/inter-FPGA communication IP, developed according to the project objective “Development of new IPs”, is the key technology behind the APEIRON framework, enabling its scalability feature. NEST-GPU is a GPU-accelerated neural network simulator engine for in-silico experiments which aims for easy reconfigurability and usage by the neurophysiology practitioner while striving for high efficiency and performance. While being self-standing production-ready code, the very significant gains in power consumption and reduced runtimes that it has demonstrated against its sibling application NEST (which is CPU-only) have motivated the current effort for its integration into the larger environment managed by the NEST Initiative and are fostering its employment in a larger number of hybrid CPU+GPU HPC platforms. Simulating a neuronal system distributed over the available GPUs, the NEST-GPU application stresses both the GPU computing and the inter-GPU communication capabilities of the IDV-A node. Furthermore, the CPU-only NEST simulator engine has been used to characterize and compare the processing and energy performance of the IDV-A and IDV-E nodes, as reported in section 3.5.

Regarding High-Energy Physics (HEP) codes on heterogeneous architectures we have selected two representative applications: Pixeltrack, a track reconstruction algorithm, and CLUE, a cluster algorithm for high-granularity calorimeters, both developed for the for the CERN CMS experiment (see section 3.4). Both applications are natively heterogeneous, being implemented with the Alpaka [1] library, and have been run both in CPU and in CPU+GPU configurations in the IDV-E and IDV-A nodes. Finally, the Tensor Network Methods (TNM) application combines multiple frontends for the simulation of quantum systems. For TEXTAROSSA, we considered the following sub-applications: i) Quantum matcha TEA: gate-based quantum circuit emulator for digitized quantum circuits, and ii) Quantum green TEA: solver for the Schrödinger equation or Lindblad equation; in the TEXTAROSSA context, we restricted ourselves to finding the ground state of a system. Both have been used to characterize the processing and power efficiency of hybrid CPU+GPU systems like the IDV-A, as described in section 3.7.

INRIA has developed a new scheduler for task-based applications that utilize heterogeneous architectures, such as FPGA+CPU or GPU+CPU. This scheduler was validated using two classical numerical simulations: Chameleon, a dense linear algebra solver, and ScalFMM, a fast multipole method library. The results indicate that this scheduler complements existing ones by enhancing the performance of irregular test cases. Additionally, Inria has improved StarONNX, a framework designed for performing inference requests on heterogeneous computing platforms, demonstrating that the new version operates efficiently on IDV-A.

PSNC provided UrbanAir, a HPC application for the purpose of benchmarking and hardware characterization, IDV-A in particular. The scope of UrbanAir is to predict air quality in urban environment, by detailed modelling of the wind flow over street canyons and complex building structure. In TEXTAROSSA, work was focused on GPU version, UrbanAir-gcrk, to run the kernels exploiting GPUs accelerators. The primary purpose was to increase energy efficiency and computation performance (lower time-to-solution). The secondary purpose was to exploit multi-GPUs available on multiple nodes to be able to run over domains which dimensions exceed the memory constraints of a single GPU. Results demonstrate that by adapting codes to accelerators it is possible to achieve primary and secondary goals.

BSC has used NBody application to test all its framework developments along the project like the fast task scheduler IP, the OmpSs@FPGA task-based programming model, the task+stream models integration, the FPGA power measurement tools or the multinode programming model Implicit Message Passing. Also, the power and cooling efficiency of the OmpSS@FPGA model in the IDV-E node has been evaluated. Beside the results of these tasks reported in this deliverable, also the integration of the RAIDER application stream kernel with the OmpSs@FPGA task-based model was done in coordination with INFN and reported in the corresponding section.

ENEA used the Vitis High-Level Synthesis (HLS) flow to implement an image processing library (FPGA Image Processing Library - FIPLib) on the Alveo U280 board. This to showcase the benefits in terms of both speed and power usage achievable using FPGA accelerators. The FIPLib has also been used in conjunction with the APEIRON framework, developed by INFN, to demonstrate how a design could be scaled to several FPGA accelerators to sustain a higher throughput, not achievable through a single FPGA for its limited resources.

CINI-UNIPI provides example results of HPC services for smart cities. Particularly, the proposed application refers to processing of video-surveillance applications for smart cities, where data acquired by telecameras are first processed with YOLO (You Only Look Once) to detect scenes containing people, and then people are automatically counted, and their position tracked using DeepSORT. The latter is an evolution of the tracking algorithm SORT: Simple Online and Realtime Tracking. Profiling of the algorithm on several platforms (those with Arm N1 processors are representative of the IDV-E) allows assessing the usefulness of different technologies. Results on RISC-V based platforms are also reported.

FHG has evaluated the usability of reduced precision Floating Point format for seismic Reverse Time Migration calculations (RTM). A small reference implementation has been created to calculate the Marmousi reference example based on 16 bit Posit and 16 bit Float format. The quality of the results has been evaluated based on the images calculated and the numerical stability of the total energy of the source wavefields over time. Storing the wavefields in Posit 16 bit delivered acceptable images. Performing the calculations in Posit 16 bit too required some optimizations but finally delivered also acceptable results. Calculating the image in Posit 16 bit too revealed higher deviations from the reference result but the image quality was still acceptable. In Float 16 bit format both storage of the wavefields and the calculation of the kernel delivered acceptable images. Calculating the image in Float 16 bit too revealed higher deviations from the reference result but the image quality was acceptable.

1 Introduction

Work performed in WP6 is essential to demonstrate the TEXTAROSSA outcomes in both hardware and software perspective. The applications need to use these for the final evaluation of the project, but far more important is to come up with conclusions if and how new hardware and software toolchain can produce benefits in terms of computation and energy efficiency of applications coming from different domains as well as can make available new and easy to use development and evaluation tools.

In the TEXTAROSSA project we proposed some applications related to AI (Artificial Intelligence), HPDA (High Performance Data Analytics) and HPC (High Performance Computing). Provided software represents quite a comprehensive set of different hardware used (CPU, GPU, FPGA), programming models and problems to be solved. Use cases are developed based on three distinctive approaches: i) adaptation to heterogeneous resources, ii) applying posit and mixed precision, and iii) using high-level programming models and their dynamic runtime systems. Therefore, there is a different set of computational, energy efficiency and accuracy metrics defined (KPI – key performance indicator) for each of the applications, though some naturally overlaps. The KPIs were discussed in the previous D6.1 deliverable. In this document we focus on applications development, and on reporting benchmarks and results. The work carried out is related to the following project objectives:

- energy efficiency: by applications developments to adapt to heterogeneous resources, energy efficient accelerators or using mixed precision;
- sustained application performance: by applications development to adapt to more computational efficient accelerators, in some cases applying a re-design of basic algorithms in order to better exploit fine-grained parallelism of accelerators, or/and in using scheduler of task/streaming-based frameworks;
- seamless integration of reconfigurable accelerators: by using the APEIRON framework;
- development of new IPs: by testing INFN intra/inter-FPGA communication IP which works behind the APEIRON framework;
- Integrated Development Platform: by using the available IDV-A and IDV-E platforms for final benchmark results.

In order to give a global picture of the TEXTAROSSA project achievements, in Figure 1 we summarize all the components developed or improved in this project, as an integrated HW/SW layered architecture. At the bottom layer we have the two types of accelerated prototype heterogeneous nodes (IDV-A and IDV-E), equipped with last generation NVIDIA GPUs and Xilinx Alveo FPGAs and the two-phase cooling system developed in this project. On the top of the above infrastructure, we have a basic toolchain including standard compilers and libraries for parallel programming as well as run-time supports and tools by vendors which, where it was needed for efficiency and scalability motivations, were directly used by some higher-level software tools and applications. In the middle layer we put all the new tools proposed in this project, including structured programming environments, application-specific workflow management tools, automatic code generation tools for exploitation of mixed-precision variable storage and computations, tools and libraries for modeling and measurements of energy consumption. Finally, on the top of this infrastructure, we put the applications, which in turn can be classified as general-purpose mathematical libraries, domain-specific codes and mini applications. All the above applications are representatives of different computational needs and were developed as significant benchmarks of some of the components

of the above architecture. Their development and testing provided useful feedback for the development and setup of all the basic components of the TEXTAROSSA architecture and significant guidelines, not only for the tools developed in this project, but more generally, for efficient and scalable exploitations of current heterogeneous computing nodes in relevant application domains.

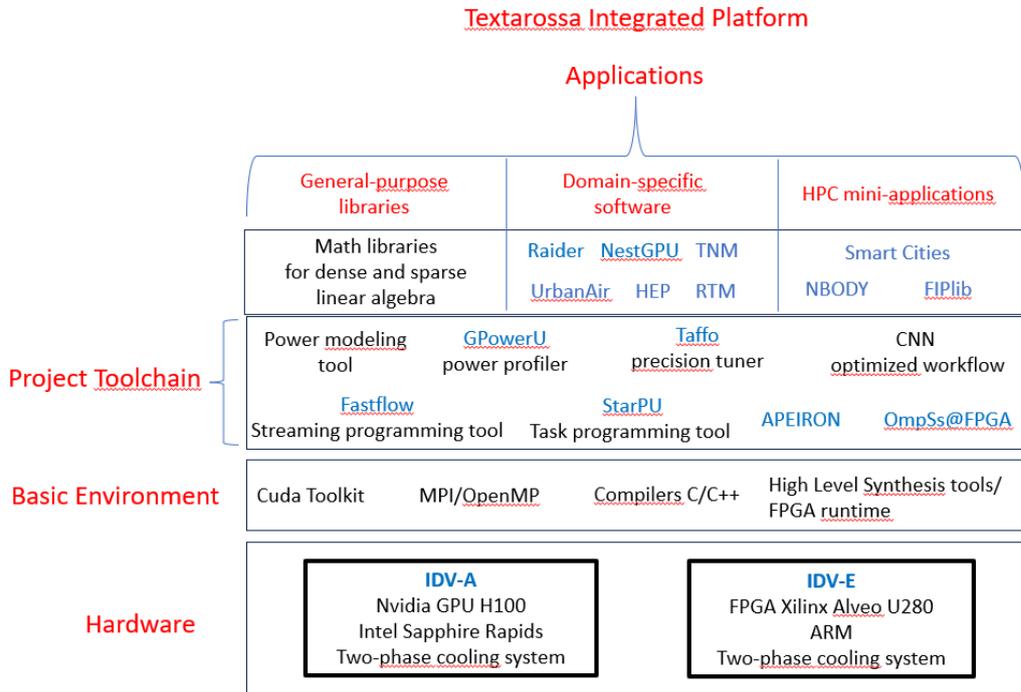


Figure 1: TEXTAROSSA Integrated Platform

This document is organized as follows. In Section 2 a brief reminder is given on hardware being used for evaluation of applications, as well as on benchmarking methodology. In Section 3, the final results of each use case are presented. Section 4 provides general summary and outline of recommendations. Section 5 crisply discusses achieved project objectives, lesson learned and guidelines.

2 Methodology

2.1 Hardware

2.1.1 IDV-A

IDV-A is based on one of Atos' 1U blade servers commercialized inside the Sequana3 architecture. It consists of one Nvidia Redstone-Next GPU board equipped with four Nvidia H100 GPUs, piloted by an Atos C4E CPU board equipped with two Intel Sapphire Rapids CPUs. It is currently cooled down by a combination of water-blocks (single-phase water heat spreaders) for the CPUs and GPUs, heat pipes and a cold-plate which serves as a chassis for the boards and distributes cold water to the various heat-spreaders. More details are described in D5.1 and D1.4.

2.1.2 IDV-E

IDV-E system is delivered by E4 and it is currently available with remote access to project partners. The nodes are equipped with ARM64 and FPGAs. The choice of the system to which to apply the two-phase cooling system fell on the Ampere Mt.Collins 2U system with Ampere Altra Max processor; the main reasons are: (i) it supports a number of PCIe slots providing the possibility of adding FPGA boards (up to 3) and/or other boards if needed, (ii) it has the physical space for adding the cooling system, (iii) it presents a good match between the amount of heat to be removed and the design point of the cooling system developed in the project, (iv) it has an architecture (ARM) compatible with that of the EPI project, (v) the possibility of receiving the system in times compatible with the project (an aspect not taken for granted given the current state of shortage worldwide). As for the FPGA, the choice fell on the U280 Xilinx Passive Model, it is able to provide significant computing power and the flexibility of memory access via HBM2 or DDR protocol with a maximum consumption of 225W. This device also guarantees the use of the Vitis High Level Synthesis software stack. More details are described in D5.2 and D1.4.

2.1.3 Other architectures

In order to assess applications results and demonstrate the potential of node-level efficiency and scalability, when more than 1 computing node are used, some other tests have been done on some available heterogeneous computers, operated by project partners or accessed by some other grants.

- For CNR-MathLib, Leonardo is used, operated by the CINECA Italian supercomputing center and granted by an Early Access Project. The booster module used for the tests is based on BullSequana XH2000 compute blade, it has 3456 nodes equipped with Intel Xeon Platinum 8358 32C 2.6GHz and 4 NVIDIA A100 SXM4 64GB GPUs. The system is interconnected through a Quad-rail NVIDIA HDR100 Infiniband. It is ranked 6th in the November 2023 Top 500 list.
- For Inria-MathLib, computing nodes are used from the plafrim cluster (www.plafrim.fr), which is a scientific instrument located in Bordeaux (France) and designed to support experiment-driven research in all areas of applied mathematics

related to modeling and high-performance computing. We used computing nodes with NVIDIA A100 or V100 GPUs.

- For UrbanAir kernels, Altair – HPC machine operated by Poznan Supercomputing and Networking Center - is used to perform additional tests in multimode environment. The nodes are equipped with Intel Xeon Gold 6242 2.8GHz and NVIDIA V100 SXM2 GPUs.
- For n-body scalability kernel CPU implementation baseline, the Marenstrum 4 machine operated by BSC was used. Each of the computing nodes are equipped with 2 Intel Xeon Platinum 8160 CPU with 24 cores each @ 2.10GHz and 96 GB of memory. Nodes are connected using Intel OmniPath 100Gb/s network
- For multi-node FPGA measurements, the MEEP machine, operated by BSC was used. This machine is formed by 96 compute nodes. Each one is equipped with two Intel Xeon Gold 6330 CPU @ 2.00GHz, 256 GB of main memory and 8 AMD Alveo U55c FPGA accelerator cards. FPGA accelerators are connected using 100Gb/s Ethernet.
- For the NEST-GPU, HEP and TNM (Quantum matcha TEA) applications we used the BSC DIBONA partition node with a dual-socket system with two AMD EPYC 7402 24-Core Processor and four NVIDIA A100-SXM4-40GB to collect results to use as baseline to compare against performance on the IDV-A node.
- For the TNM (Quantum green TEA) the CINECA Leonardo booster node and the Justus2 cluster node (bwHPC): 2xIntel Xeon 6252 Gold with 2x24 cores
- For the RAIDER HPDA multi-FPGA applications we used the INFN APE Lab small development cluster made of 4 single Intel(R) Xeon(R) Silver 4410T CPU nodes, each of them hosting one Xilinx® Alveo U200 FPGA board, with the four FPGAs interconnected in a ring topology.

2.2 Application tuning, performance and energy measurements

In order to have meaningful and consistent benchmarking results across different use cases, a common methodology for measuring the performance and energy efficiency is proposed and discussed in detail in D1.4, and followed in D6.2. Some of the approaches are described in this document in Section 3, where each use case presents in detail the approach to obtain results and measure application-specific KPIs (proposed in D6.1). Where it was needed to follow an architecture-specific pathway and/or adaptation at the algorithmic and/or implementation level, this is discussed in details so that this Deliverable can be read as a set of guidelines for potential users of all the HW/SW IPs developed in the project.

3 Results

3.1 Smart cities – CINI-UNIFI

3.1.1 Smart Cities algorithm description and implementation assumptions

CINI-UNIFI provides example results of HPC services for smart cities. The proposed application refers to processing of video-surveillance applications for smart cities, where data acquired by telecameras are first processed with YOLO (You Only Look Once) to detect scenes containing people. Then people are automatically counted and their position tracked using DeepSORT, which is an evolution of the tracking algorithm SORT: Simple Online and Realtime Tracking. YoloV5 is an algorithm that ensures low-latency in extracting the object of interest in a complex scene. YoloV5 was trained on a custom labelled dataset (to recognize specifically people and people laying down). Moreover, DeepSORT is designed to achieve real-time. Hence, the combination of both algorithms ensures a rapidity that is essential in mission-critical applications. DeepSORT takes advantage of OSNet x1.0, a convolutional neural network used for person re-identification, to identify the same person in the video stream.

Thanks to a geometrical check of the bounding boxes generated by Yolo to highlight the target of interest in the scene (people), the algorithm can also detect if people are laying down or not. This is important when monitoring zones of the city after a natural disaster (e.g. earthquake, landslide) or in a war scenario.

Profiling of the algorithm on several platforms (those with Arm N1 processors are representative of the IDV-E) allows assessing the usefulness of different technologies. Results on RISC-V based platforms are also reported.

To be noted that as far as the computation partitioning between general-purpose processor and Hardware (HW) accelerator is concerned, Deep-SORT has an irregular data flow and hence is less suited for HW acceleration on FPGA than Yolo. For YoloV5 we considered both the case it is implemented SW-wise on the general-purpose processor or HW-wise via an FPGA acceleration. In this case the complete flow can be pipelined with the FPGA processing with Yolo the input video flows and then the general-purpose computing cores applying DeepSORT and other control rules.

In all cases the General Purpose CPU also implements some pre-processing of the input videos acquired from the surveillance camera.

As far as mixed-precision is concerned, the Posits were tested through the CppPosit software library developed by UNIFI [2].

Indeed, in TEXTAROSSA it is not foreseen designing a custom processor as an ASIC with dedicated HW support for new formats like Posits.

Moreover, some of the selected platforms (e.g. those with NVIDIA GPUs) consider automatically an optimization of multiple arithmetic precisions among float32, bfloat/float16 and integers (8 or lower). The native mixed-precision is more effective than emulating via software the CppPosit library on a given platform.

In the reported Figure 2 and Table 1 the results refer to the best mixed-precision configuration.

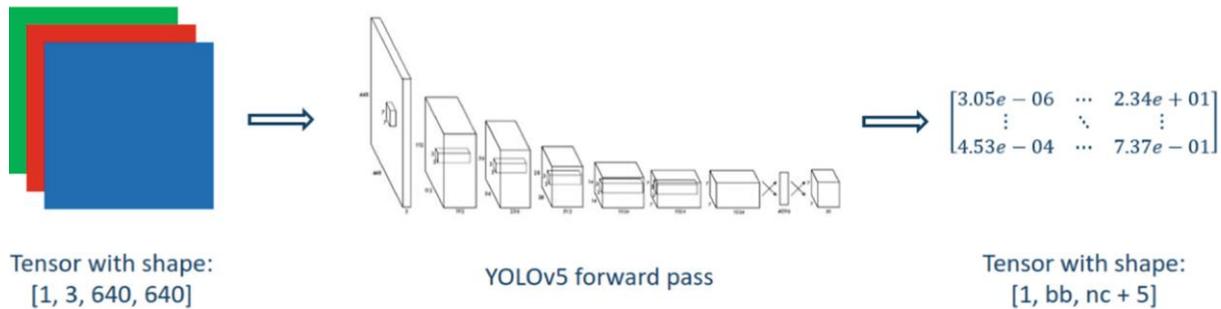


Figure 2: [SmartCities] Computation flow for the YOLOv5 step

3.1.2 Target platform and implementation results

We evaluated the aforementioned deep learning application through Pytorch on several architectures. We employed architectures with and without accelerators (GPUs, FPGA), to evaluate also the performance of vector units inside the CPUs, when available. The tested CPU-only architectures are based on ARM, RISC-V and Intel.

The platforms used in the test are available at the Green Data Center of University of Pisa ([GREEN DATA CENTER \(unipi.it\)](https://www.unipi.it/green-data-center)).

Most of the reported HPC platforms, including those with ARM N1 processors, are installed and maintained at UNIFI Green Data Center by E4 and hence the results achieved are coherent with those achieved on machines made available via VPN by E4.

For ARM we have considered ARM 64 bit architectures and particularly: the ARM Neoverse N1 that is the one available in IDV-E, plus that in the Fujitsu A64FX CPU that is the ARM SVE (scalable vector extension) and finally the ARM Cortex-A72 available in Xavier NVIDIA SoC (system-on-chip) computing board.

- ARM Cortex-A72 CPU (with ARMv8.0-A 64-bit instruction set) with four cores running at 1.5 GHz, with support for the 128-bit Neon SIMD extension and 4 Gbyte of RAM.
- ARM Cortex-A78 CPU (with ARMv8.2-A 64-bit instruction set) has been also tested since is inside the AGX Orin with 12 cores
- HPE Apollo80 system, powered by a Fujitsu A64FX, running at 2.2 GHz with support for the ARM 512-bit Scalable Vector Extension (SVE) SIMD unit.
- Ampere Altra Max with ARM Neoverse N1 CPU with support to 128-bit NEON SIMD instructions (IDV-E platform). ARM Neoverse N1 CPU being derived from Cortex-A76 supports the ARMv8.2-A 64-bit instruction set as in the Cortex-A78 mentioned above
- HiFive Unmatched board powered by a SiFive U740 SoC running at 1.2 GHz for the scalar RISC-V unit called Arriesgado.
- Cascade Lake-based Intel Xeon Silver CPU running at 2.2 GHz with support for the 512-bit AVX SIMD instructions.
- I7-107650H processor from Intel, implemented in 14nm is a 10th generation x86 I7 family CPU with AVX2 and SSE4.2 instruction extensions (256-bit SIMD instructions), with 2.6 GHz basic frequency (boot up to 4.8 GHz)

The tested GPU-based architectures are:

- NVIDIA Jetson Orin SoC with an Ampere-based GPU
- NVIDIA Tesla T4
- NVIDIA A100

Hereafter, we report the benchmark results for the platforms listed above. In Table 1 we report the average number of frames processed per second for each platform on the whole application (preprocessing and YoloV5 + people down check + DeepSORT) application. Furthermore, we detail, for each architecture, the timing performance for the three different components of the application.

Platform	FPS	YOLO (ms)	People-down(ms)	DeepSORT (ms)
ARM Cortex-A72	0.11	7493	1.2	1107
ARM NeoverseN1	0.73	858	0.6	503
IDV-E: ARM N1+FPGA	2	(FPGA)	(FPGA)	503 (ARM)
ARM A64FX	0.43	1292	1.1	1054
RISC-V SIFIVE U7540 emulating Arriesgado#	0.006	148987	2.6	13835
AGX Orin (12 Cortex-A78 cores CPU + GPU Ampere) *	7.7	38.8	0.5	54.9
Intel i7	0.98	807	0.2	207
Intel i7 + GPU GeFore 1650Ti (Turing)*	7.3	85.4	0.3	11,9
Intel Xeon	1.88	335	0.3	197
Intel Xeon + GPU Tesla T4 *	12.8	37.8	0.3	15.1
Intel Xeon + GPU A100*	21.3	10.9	0.3	13.9

Table 1: [SmartCities] Implementation results and benchmarks for different platforms

Please note that for the 4 cases with * where we have a CPU plus a GPU accelerator, due to bottlenecks in the data transfer between the CPU part and the GPU part the overall frame-per-second (FPS) is lower than the inverse of the sum of the individual frame processing time (i.e. for Intel Xeon + GPU A100* the sum of frame processing time would be 25.1 ms, i.e. its inverse would be a theoretical FPS of 39.84 FPS but due to data transfer bottlenecks among CPU-GPU we measured 21.3 FPS).

Please also note that for Arriesgado scalar RISC-V emulated on SI-Five u7540 the very low FPS performance achieved when compared to other solutions is due also to the fact that the architecture is emulated.

It is worth remembering that the Tesla A100 has fully native support of mixed precision. A100 Tensor Core GPU performance specs [3]:

- Peak FP64 9.7 TFLOPS
- Peak FP32 19.5 TFLOPS

- Peak FP16 78 TFLOPS
- Peak BF16 39 TFLOPS
- Peak TF32 Tensor Core 156 TFLOPS
- Peak FP16 Tensor Core 312 TFLOPS
- Peak BF16 Tensor Core 312 TFLOPS
- Peak INT8 Tensor Core 624 TOPS
- Peak INT4 Tensor Core 1,248 TOPS

In the Figure 3 we analyze the power efficiency for all platforms with an FPS value higher than 7, that is the minimum considered useful for a real application. The power efficiency is measured at application level, so the KPI is Frame-Per-Second/Watt or Frame/Joule being 1 Joule equal to 1 Watt * 1 second.

The results achieved show that the most efficient solution is the one with multicore Cortex-A78 64b plus NVIDIA GPU, although optimizations towards an increases power efficiency are needed since the cost is still 2 Joule for each processed Frame.

If we need to increase the performance to a real-time value for more than 20 FPS the best solution is the Intel XEON + A100 GPU with an energy cost of 2.5 Joule per processed frame.

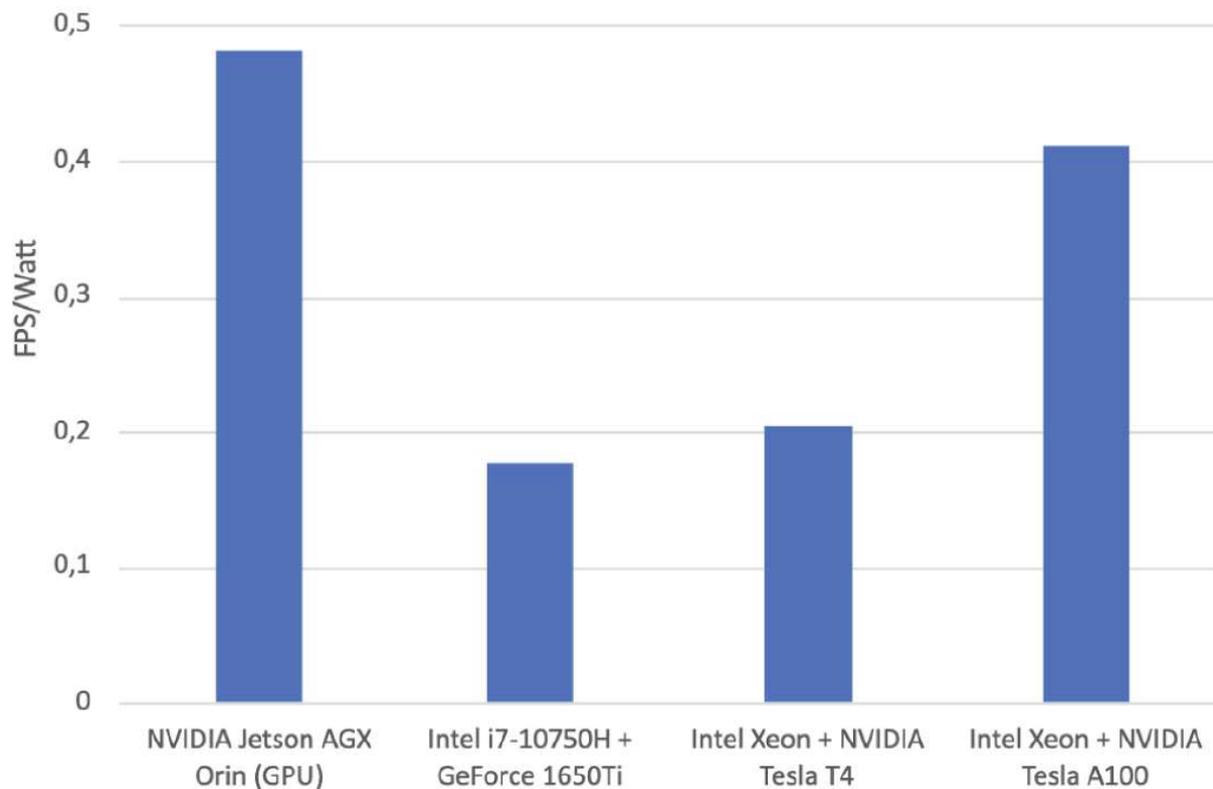


Figure 3: SmartCities] Power efficiency measured as FPS/Watt

the higher is the better, for the platforms with FPS higher than 7 (that is the minimum considered for a practical use)

From the results achieved, the key lesson learned to set-up an edge server for surveillance applications in smart cities are:

- RISC-V computational capabilities, particularly in scalar version, and considering commercially available solutions are still far in performance from other platforms like those based on Intel and/or ARM

- For ARM the SVE version in Fujitsu performed worse than the ARM Neoverse N1 in Ampera Altra Max (the one in IDV-E) that is not using SVE; this can be justified by the fact that the algorithms were not optimized for the scalar vector extension of the processor as implemented in Fujitsu.
- Multicore ARM Cortex-A 64b (like in the A78 with ARM v8.2 instruction set) sustained by a GPU acceleration are the most efficient solution for edge –AI/Video applications.
- Looking to edge server performance, the combination of Intel Xeon plus a GPU Like Tesla A100 is the best in terms of pure performance, FPS as KPI, but is also a good solution in terms of energy efficiency, FPS/Watt KPI. Since the new ARM Neoverse V2 is appearing in the market with ARM 64b v8.4 instruction set and vectorized instruction, we expect that ARM Neoverse V2 plus GPU can be as effective as Intel Xeon plus A100 GPU.
- The availability of FPGA Xilinx Alveo in IDV-E can be exploited by porting the Yolo calculation on it while keeping the DeepSort on the ARM processor and hence a gain in speed can be achieved roughly by a factor of 2. However, the performance is still below those that can be achieved using Intel Xeon plus a GPU Like Tesla A100.
- Modern GPUs like A100 natively support mixed-precision from very low-format (from INT4 to FP64 including new FP16 or Bfloat16)
- New chip evolutions expected from main US companies will extend this mixed-precision support to microscaled Floating Point formats having also FP4 and FP8 as alternative to IEEE 754 standard, see this joint paper from Intel, Meta, Nvidia, AMD, Microsoft, Qualcomm [2310.10537.pdf \(arxiv.org\)](https://arxiv.org/abs/2310.10537) entitled “Microscaling Data Formats for Deep Learning” and the open library from Microsoft [GitHub - microsoft/microscaling: PyTorch emulation library for Microscaling \(MX\)-compatible data formats](https://github.com/microsoft/microscaling). This demonstrates the value of an independent research on compact floating –point variants vs the IEEE 754 standard, like Posits and integration for Posits with other formats in tools like TAFFO. An idea on future work can be the combination in TAFFO of Posits16, 32 with FP32/FP64 for high accurate numerical applications and Posit8, with BFloats16, FP8, FP4, INT4, INT8 for video or AI applications where the classification accuracy can tolerate a higher numerical inaccuracy.

3.2 MathLib – CNR

In this section we discuss results obtained by CNR using the mathematical software library for hybrid architectures featuring NVIDIA GPUs at node level. As already described in previous deliverables, the CNR team has developed computational kernels required in sparse matrix computations and iterative linear solvers, which are widely exploited in Scientific Computing and Data Analysis. Our library was thought both as main component of the TEXTAROSSA platform and as benchmark tool for IDV-A. Focus was on GPU-kernels efficiency and on scalability when multiple GPUs, also on different computing nodes, are needed for computations because dimensions largely exceed the memory resources of a single GPU. Therefore, on IDV-A, the library stresses the GPU operation capabilities and memory/communication channels bandwidth at the node level. The Mathlib development toolchain includes C compilers, the Cuda Toolkit and the MPI library available as basic environment of the TEXTAROSSA platform. Our choice of the above basic tools was motivated by the need to re-use some very efficient GPU kernels already available in the library baseline (for 1 NVIDIA GPU) and to focus on algorithms and implementation scalability when

very large number of nodes could be needed. Moreover, extensive use of the GPowerU project tool¹ has been made for monitoring GPU kernels energy consumption. In the following we first discuss some adaptation and approximations needed at the algorithmic and implementation levels to have benefits from the use of heterogeneous systems and then discuss main results.

The MathLib computational kernels developed and tested in the project are the following:

- Sparse matrix – vector multiplication (SpMV);
- Sparse matrix – matrix multiplication (SpMM);
- Maximum Weight Matching in undirected graphs (MWM);
- Preconditioned Conjugate Gradient (PCG) method coupled with a matching-based Algebraic MultiGrid preconditioner. Note that in the following we use the acronym BCMG to refer to the complete solver.

We point out that the BCMG solver is implemented on the base of all the other computational kernels, which are the main blocks for AMG setup (SpMM and MWM) and for solving by PCG (SpMV). We present performance and energy consumption results obtained on IDV-A and then, in order to analyze the scalability potential of our main sparse linear solver based on the AMG-preconditioned Conjugate Gradient method (BCMG), we also show performance, in a weak scalability regime, obtained on the Leonardo supercomputer, whose access was granted through a Leonardo Early Access Project.

As benchmark data sets, we consider matrices and right-hand sides of algebraic systems required for the solution of the Poisson equation in 3D with homogeneous Dirichlet boundary conditions and unitary right-hand side. This is a standard benchmark test case for sparse matrix computations because it represents the computational kernel of many scientific and engineering applications, and indeed is also used in the HPCG benchmark [4]. In our case, the discretization of the problem is obtained by the classic 7-points finite-difference stencil for the left-hand side operator (the Laplacian operator), which results in a symmetric positive definite (s.p.d) matrix of coefficients well suited for PCG. Note that, in all the experiments made with the BCMG solver and discussed in the following, we stop PCG iterations when the relative residual in the Euclidean norm is less than 10^{-6} or the number of iterations reaches the maximum value fixed to 1000 (actually, in all the experiments with the linear solver, the required accuracy is obtained with no need to stop for the limit on the maximum number of iterations). In the case of the SpMV kernel, we consider, as vector operand, a vector of all ones, whereas in the case of the SpMM kernel, both the operands are the same, so that we compute the square of a Laplacian matrix. Finally, for the MWM kernel, we consider the undirected adjacency graphs of the Laplacian matrices to which suitable real weights are associated, as applied in the BCMG aggregation algorithm for the preconditioner setup.

We note that our baseline was a mono-GPU version of all the kernels, while in this project we focused on a hybrid parallel version leveraging multi-GPUs computing nodes. We are interested in analyzing both strong scalability, i.e., the reduction in the execution times when a problem with a fixed size is considered on an increasing number of parallel resources, and weak scalability properties of our kernels, i.e. when dealing with problems of increasing size, while parallel resources increase. The parallel design pattern is based on Single Program Multiple Data (SPMD) programming model relying on a row-block distribution of the system matrix and the related right-hand sides among the parallel tasks. Blocks of contiguous rows are assigned to each task. We observe that each task is associated to one GPU accelerator which oversees all the computation phases.

¹ <https://github.com/crossi/GPowerU>

Details on the algorithms and parallel design patterns implemented for all kernels of the MathLib are discussed in [5] and reported in Deliverable 6.2. Here we mention, as general guidelines for users and software developers, that some approximations in original (mono-GPU) numerical algorithms have been required, so that the combination of communication-avoiding algorithms, fine-grained parallelism, and overlapping between computation and data communication could allow us to design a scalable version of the BCMG linear solver. Indeed, In the context of heterogeneous computing, where we want to use accelerators like GPUs to maximize performance, developing effective approaches for the iterative solution of linear systems brings the need for new methods, algorithms, and implementations capable of exploiting the underlying hardware and basic software. Regardless of the method adopted for solving the system, the computation on sparse matrices is particularly challenging with respect to its dense counterpart due to the irregular memory access pattern and intrinsic load imbalance caused by the sparsity pattern of the matrix rows. GPUs rely on fine-grained parallelism and access to medium size memories with high bandwidth, but also not negligible latency. Therefore, although the current software stack makes available programming environments which provide a clear interface to the features of the hardware, it is still challenging to use the above accelerators efficiently, especially for memory-bound kernels, like sparse matrix-vector products involved in iterative linear solvers. For example, the NVIDIA GPUs are built in terms of arrays of multithreaded, streaming multiprocessors, and each multiprocessor is composed of a fixed number of scalar processors. The CUDA programming paradigm is based on the concept of blocks of threads which share data. Therefore, having a regular density in the rows of a matrix is a favorable situation for high-throughput SIMD operations, whereas the irregular structure of general sparse matrices poses some limitations for efficient usage of the architecture which often are smoothed by organizing matrices in suitable data structures. Furthermore, a sequence of SIMD operations applied to the same data allows realizing the so-called data/thread-locality which makes GPU exploitation very efficient. For the same reasons, basic iterative algorithms which express a high-level of data parallelism are preferred to more complex algorithms which induce data dependency although, generally, the former may have worse convergence properties. Indeed, extra computation is often well tolerated and balanced by very efficient execution on the GPU. On the other hand, the ability to exploit the scalability of large clusters of heterogeneous nodes largely depends on appropriate coordination among multiple levels of computation so that data partitioning, data/workload balancing, data communication between GPUs, CPUs, and among distributed nodes, do not penalize in a significant way the final performance. In the paper we also included details on the algorithmic parameters which characterize the AMG preconditioner which, for the sake of brevity, we omit in this deliverable.

3.2.1 Performance and Energy Consumption on IDV-A

In this section we discuss results of different kernels, in terms of both parallel performance and power/energy consumption. For both strong scalability and weak scalability analysis we consider different matrix dimensions per each kernel, so that each GPU can be used at full load, as detailed in the tables. Power consumption measures have been obtained for all the kernels varying the number of GPUs. As also described in deliverable D6.1, our main KPIs are grouped into 3 main categories, as classified in Table 2. We point out that in this deliverable, after some discussion among the project partners, with the aim to follow a common practice, we modified KPIs for energy, adding a measure considering energy consumption and a measure whose ratio has an average power to the denominator. Finally, we observe that for all the kernels we analyze results obtained with the final configuration of IDV-A, including all the HW/SW toolchain driving the 2-phase cooling system, while for the BCMG solver we were able to make a comparison between and after the installation of the 2-phase cooling system.

KPI for energy	KPI for computational efficiency	KPI for accuracy
Iterations/Joule Iterations/AvgWatt (only for iterative linear solver)	Execution time	Yes (User's parameter dependent for iterative linear solver)
Dofs (Degrees of Freedom or unknowns)/Joule Dofs/avgWatt	(strong and weak) speedup	
	Number of iterations/time per iteration (only for iterative linear solver)	
	Accuracy	

Table 2: [MathLib CNR] KPIs

Accuracy of the results of all kernels, except the linear solver, is up to the machine precision in double precision floating-point arithmetic. In the case of the BCMG linear solver, as already said, an accuracy on the solution up to 6 digits is achieved. Different tolerances in the stopping criterion can be set up to reduce or increase solution accuracy. For the sake of completeness, we include in the tables for every kernel, also the values of peak and average power as well as the total energy measured during the execution of the kernels. Note that for the above measures we follow a so-called dynamic approach, that is we neglect power and energy of idle GPUs state and we only consider values when the kernels start their execution, including all phases of I/O from CPU to the GPU and till the GPUs go back to the idle state. In Table 3, Table 4 and Table 5 we report the performance of the SpMV, SpMM and MWM kernels, respectively.

GPUs	Dofs	Time (sec.)	Sp	Power (peak)	Power (avg)	Dofs/avgWatt	Dofs/Joule	Total Energy
Strong Scalability								
1	3.1×10^8	0.30	1	157	116.0	2.6×10^6	6.17×10^5	496.10
2	3.1×10^8	0.26	1.2	138	116.5	2.6×10^6	4.16×10^5	736.43
4	3.1×10^8	0.20	1.5	126	116.3	2.6×10^6	2.43×10^5	1262.31
Weak scalability								
2	6.2×10^8	0.34	1.8	138	116.5	5.3×10^6	3.41×10^5	1794.62
4	12.4×10^8	0.48	2.5	159	117.5	10.4×10^6	7.52×10^5	1629.43

Table 3: [MathLib CNR] Performance and Energy Consumption of the SpMV kernel

We observe that, as expected for such a communication-bound kernel, parallel efficiency degrades when GPU number increases also in a weak scaling setting; a larger degradation can be observed for the SpMM kernel in Table 4. Furthermore, while total energy consumption is almost preserved in the weak scaling setting going from 2 GPUs to 4 GPUs, in the strong scaling, doubling number of GPUs from 2 to 4 produces a significant increasing ratio in energy consumption. Similar behavior in terms of energy consumption is observed for the SpMM kernel in the strong scaling setting, indeed the energy consumption increases for increasing GPU number, an increase, as expected from the increase in execution times, is also observed in the weak scaling setting for this kernel. The behaviour observed for energy consumption of the MWM kernel is very similar to that of the SpMV kernel, where an increase for increasing

GPUs is observed in the strong scaling while in the weak scaling setting it remains almost stable going from 2 to 4 GPUs. This behavior of the energy consumption of the above kernels, in the weak scaling setting, has to be better investigated.

GPUs	Dofs	Time (sec.)	Sp	Power (peak)	Power (avg)	Dofs/avgWatt	Dofs/Joule	Total Energy
Strong Scalability								
1	3.8x10 ⁷	0.38	1	189	119.0	3.2x10 ⁵	6.72x10 ⁴	569.56
2	3.8x10 ⁷	1.46	0.26	158	119.0	3.2x10 ⁵	5.03x10 ⁴	761.00
4	3.8x10 ⁷	1.56	0.24	138	118.8	3.2x10 ⁵	3.31 x10 ⁴	1155.85
Weak scalability								
2	7.6x10 ⁷	0.59	1.25	191	121.0	6.3x10 ⁵	5.13x10 ⁴	1493.37
4	15.2x10 ⁷	2.00	0.74	174	120.8	12.7x10 ⁵	7.69x10 ⁴	1991.30

Table 4: [MathLib CNR] Performance and Energy Consumption of the SpMM kernel.

GPUs	Dofs	Time (sec.)	Sp	Power (peak)	Power (avg)	Dofs/avgWatt	Dofs/Joule	Total Energy
Strong Scalability								
1	2.7x10 ⁸	0.257	1	191	118.0	2.3x10 ⁶	5.21x10 ⁵	515.20
2	2.7x10 ⁸	0.128	2.0	164	118.5	2.3x10 ⁶	3.41x10 ⁵	788.02
4	2.7x10 ⁸	0.064	4.0	143	116.8	2.3x10 ⁶	2.45x10 ⁵	1093.98
Weak scalability								
2	5.4x10 ⁸	0.257	2.0	201	119.5	4.5x10 ⁶	6.98x10 ⁴	1758.77
4	10.8x10 ⁸	0.257	4.0	200	120.0	8.9x10 ⁶	9.68x10 ⁴	1878.84

Table 5: [MathLib CNR] Performance and Energy Consumption of the MWM kernel

In the following we put results obtained with the BCMG solver before and after installation of the 2-phase cooling system. In the figures we report pictures of the power consumption behavior of the execution of the BCMG solver on 4 GPUs, when each GPU oversees a local matrix with about 61x10⁶ Dofs for a linear system with a total of 2.4x10⁸ Dofs.

In Figure 4 we report the behavior before the 2-phase cooling installation and in Figure 5 we show the same behavior after installation. We observe very similar power behavior but in the right side of the picture, where we see that GPUs go toward the idle state power rapidly drops in Figure 4 while in Figure 5 the same time was not sufficient to reach the idle state power, showing that the cooling controller has some impacts on the transition of the GPUs to the idle state. If we go deep inside the numbers reported in Table 6 and Table 7 we can better observe that, while performance of the solver shows a small degradation after the installation of the 2-phase cooling system, the energy consumption shows a slight reduction.

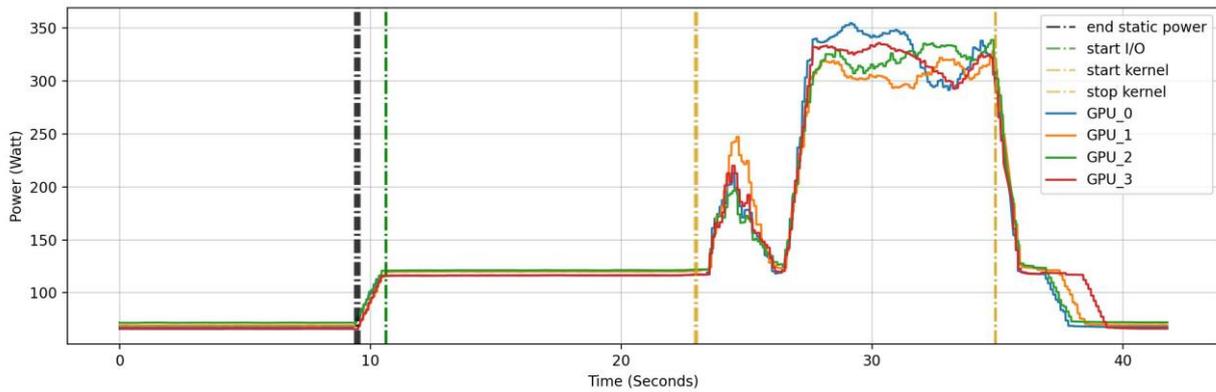


Figure 4: [MathLib CNR] Power behavior of BCMG on IDV-A before the 2-phase cooling system installation

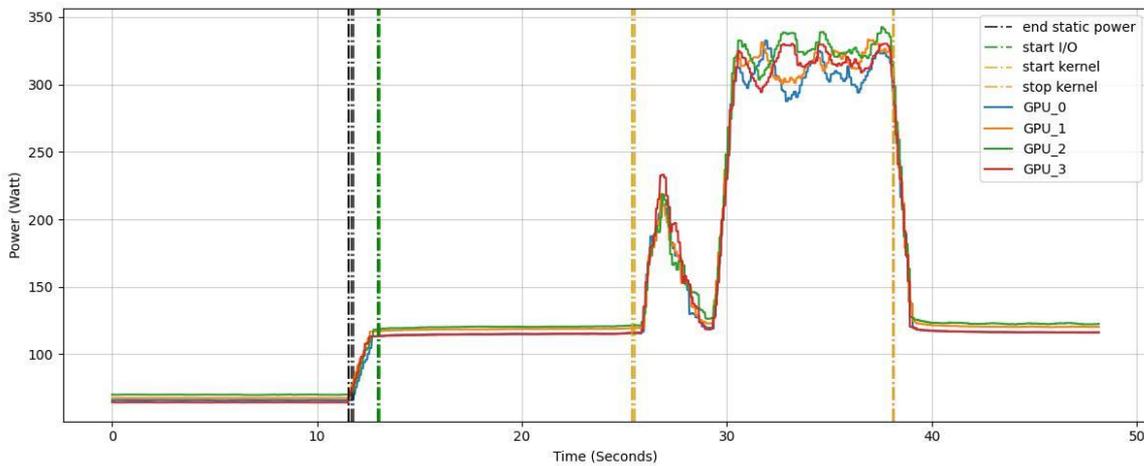


Figure 5: [MathLib CNR] Power behavior of BCMG on IDV-A after the 2-phase cooling system installation

In Table 6 and Table 7 summarize results for the BCMG solver, as for the other kernels, before and after installation of the 2-phase cooling system.

GPUs	Dofs	Time (sec.)	Sp	Power (peak)	Power (avg)	Dofs/avg Watt	Dofs/Joule	Total Energy
Strong Scalability								
1	6.1×10^7	8.59	1	399	178.0	3.4×10^5	2.85×10^4	2147.72
2	6.1×10^7	5.25	1.6	338	166.5	3.7×10^5	2.37×10^4	2578.74
4	6.1×10^7	4.45	1.9	247	148.3	4.1×10^5	1.86×10^4	3283.00
Weak scalability								
2	12.2×10^7	9.96	1.7	352	162.0	7.6×10^5	2.03×10^4	6011.36
4	24.4×10^7	11.27	3.1	354	153.0	16.0×10^5	1.90×10^4	12898.43

Table 6: [MathLib CNR] Performance and Energy Consumption of the BCMG solver before installation of the 2-phase cooling system

GPUs	Dofs	Time (sec.)	Sp	Power (peak)	Power (avg)	Dofs/avgWatt	Dofs/Joule	Total Energy
Strong Scalability								
1	6.1×10^7	9.77	1	400	169.0	3.6×10^5	3.18×10^4	1920.95
2	6.1×10^7	6.14	1.6	319	156.5	3.9×10^5	2.70×10^4	2268.53
4	6.1×10^7	5.68	1.7	235	140.3	4.4×10^5	2.15×10^4	2843.50
Weak scalability								
2	12.2×10^7	10.70	1.8	352	156.0	7.8×10^5	2.09×10^4	5863.95
4	24.4×10^7	11.98	3.3	342	152.5	16.0×10^5	1.95×10^4	12524.63

Table 7: [MathLib CNR] Performance and Energy Consumption of the BCMG solver after installation of the 2-phase cooling system

3.2.2 Weak Scalability results of BCMG up to 64 billion dofs on the Leonardo Supercomputer

In this section, we analyse the scalability potential of our BCMG sparse linear solver when the number of GPUs largely increases; the fixed matrix size per node is equal to $250^3 \sim 15.6\text{M}$ dofs, going from 1 to 4096 GPUs of the general-purpose module of the Leonardo Italian Supercomputer operated by Cineca. Therefore, we solve problems up to 64 billion Dofs. We analyse the performance of the linear solver looking at the number of iterations to obtain the desired accuracy, the execution time to solve the system and the execution time per each solver iteration, to characterize the ability of the solver in leveraging large-scale resources for solving problems of increasing size. In the following figures, indeed, we report all KPIs which characterize the application phase of a linear solver. In Figure 6 we show the number of iterations needed to reach the required accuracy for increasing dimension and number of GPUs. We can observe that, as expected, due to a general increase in the conditioning number of the system matrices and to the uncoupled aggregation approach in the setup of the AMG hierarchy, number of iterations changes for increasing number of GPUs; after a rapid increase from 1 to 4 GPUs, we observe a good algorithmic scalability, with an almost stable number of iterations, up to 2048 GPU, while a significant increase is observed going from 2048 to 4096 GPUs showing a degradation of the quality of the preconditioner. This degradation is also responsible for the increase of the solve time, as observed from Figure 7. Moreover, we point out that we were able to solve a linear system with 64 billion dofs in less than 3.8 seconds, with an increase in solve time of less than 6.5 times to solve a system whose size increases by more than 4 thousand times, going from 1 to 4096 GPUs. In Figure 8 we show the execution time per iteration which is a KPI showing efficiency and scalability of the implementation of BCMG and all its computational building blocks. We can observe an increase in the time per iteration which is expected for such a type of communication/memory bound problems, however we observe an increase of only 3.4 times going from 1 to 4096 GPUs, which indicates a very good implementation scalability and leads to a scaled speedup of the solve phase of 632.5 on 4096 GPUs (see Figure 9).

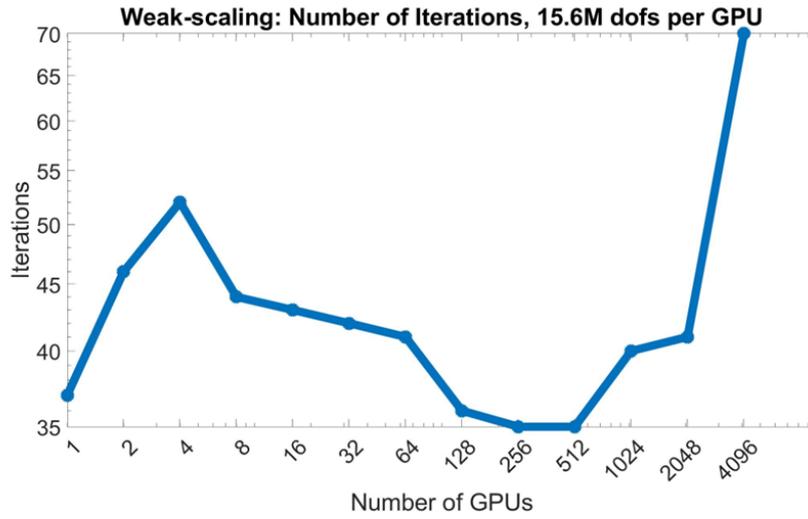


Figure 6: [MathLib CNR] Weak scalability on Leonardo up to 64 billion dofs, number of iterations

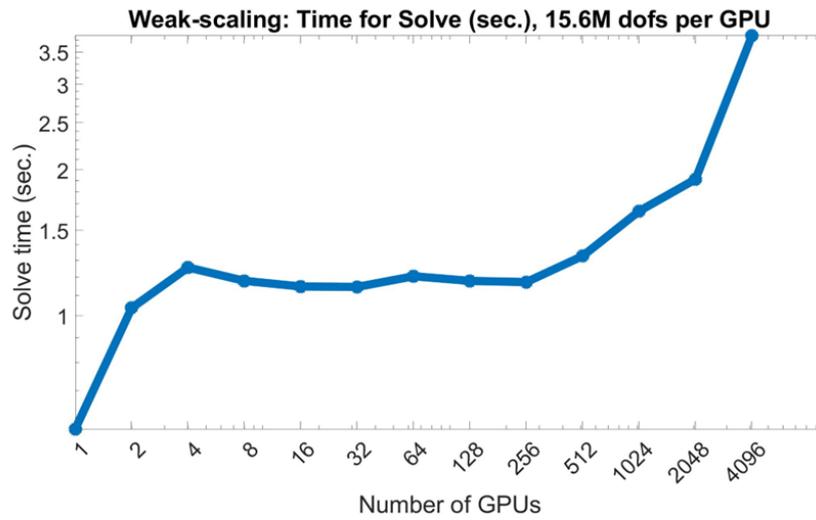


Figure 7: [MathLib CNR] Weak scalability on Leonardo up to 64 billion dofs, solve time

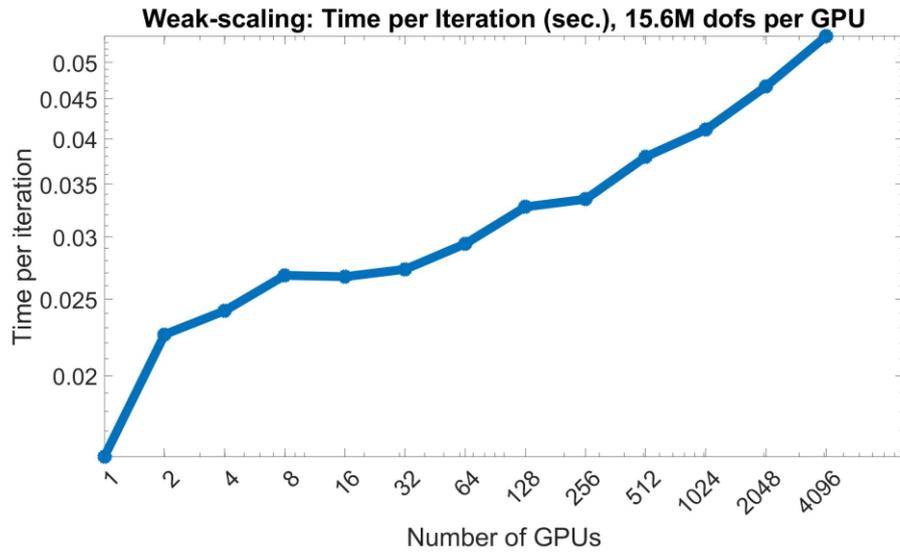


Figure 8: [MathLib CNR] Weak scalability on Leonardo up to 64 billion dofs, execution time

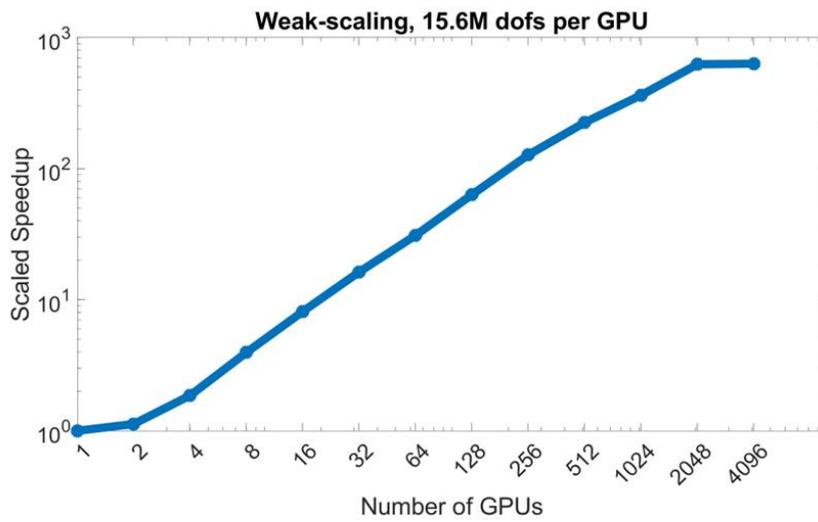


Figure 9: [MathLib CNR] Weak scalability on Leonardo up to 64 billion dofs, speedup

3.2.3 Comparison with NVIDIA AmgX library

In this section we show comparisons with the hybrid version of NVIDIA AmgX [6, 7], that is the state of the art for AMG-preconditioned sparse linear solver on NVIDIA GPUs. AmgX makes available various AMG preconditioners, based on different well-known coarsening approaches already available in other libraries, and producing AMG hierarchies with different computational complexities. For a fair comparison, we selected two input configurations for the AMGX solver: AMGX-C refers to a preconditioned Conjugate Gradient solver based on a robust preconditioner having very good algorithmic scalability at the cost of a large computational complexity, while AMGX-AGG refers to the PCG solver coupled with an aggregation-based preconditioner, which defines AMG hierarchies based on a similar coarsening approach and having complexities comparable with our preconditioner. Due to the limitation in using more than 128 nodes, we were able to run our tests on up to 512 GPUs of the Leonardo supercomputer. Stopping criteria are the same for all the solvers and are the same we used for the tests with BCMGX discussed in the previous sections. For these tests we were able to use a problem size per GPU equal to $200^3 = 8$ million Dofs up to more than 4 billion dofs on 512 GPUs. In Figure 10 we show number of iterations needed to reach the desired accuracy. We can see that, as expected, the best behavior is obtained by AMGX-C, while our BCMGX shows largely better algorithmic scalability with respect to AMGX-AGG. In Figure 11 we show total solve time. Here we observe that the smallest number of iterations obtained by AMGX-C does not balance its large cost per iteration (see Figure 12) for a number of GPUs larger than 8 and then for increasing number of GPUs and problem size, our BCMGX obtains the smallest execution time with respect to the AMGX solvers. The above results are very promising and demonstrate the benefits of our new parallel algorithms and implementation design patterns for heterogeneous architectures. Our MathLib can be considered, indeed, as a main tool of the TEXTAROSSA Integrated Platform for Scientific Computing leveraging heterogeneous architectures exploiting NVIDIA GPUs.

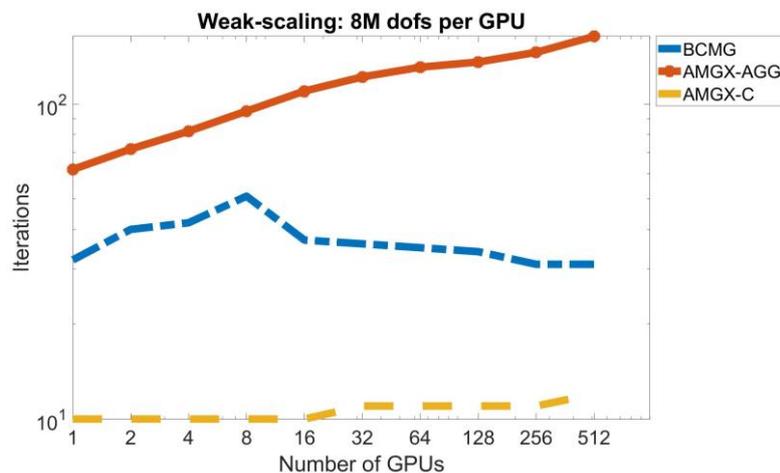


Figure 10: [MathLib CNR] Weak scalability on Leonardo. BCMGX vs AMGX, number of iterations

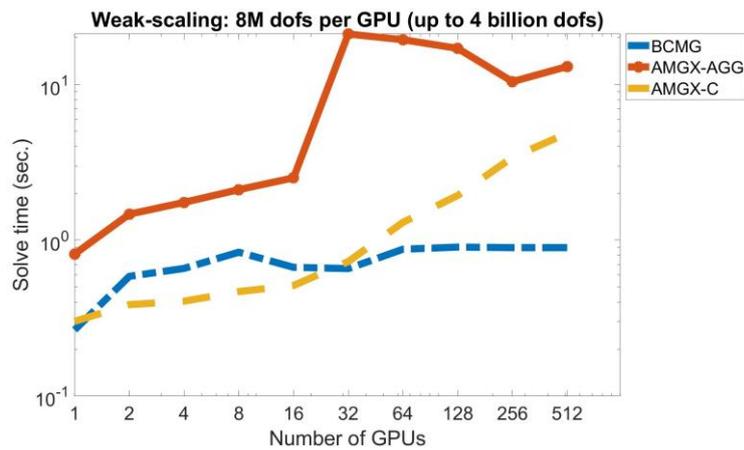


Figure 11: [MathLib CNR] Weak scalability on Leonardo. BCMGX vs AMGX, solve time

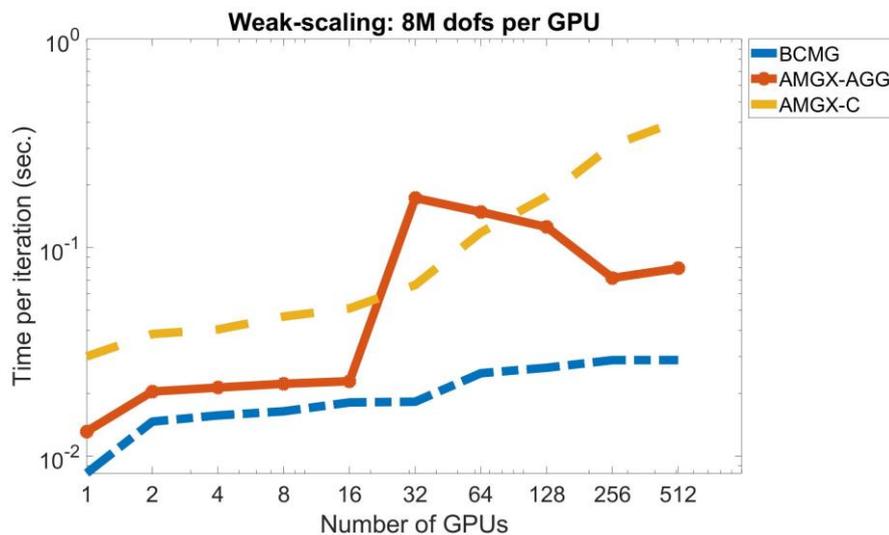


Figure 12: [MathLib CNR] Weak scalability on Leonardo. BCMGX vs AMGX, time per iteration

3.3 RTM – FRAUNHOFER

Reverse Time Migration (RTM) is used in the field of seismic for oil and gas exploration. Among other migration algorithms, RTM is able to visualize more details compared with more simple algorithms as e.g. Kirchhoff migration that uses high frequency ray approximation. Especially steep dips can't be mapped well by Kirchhoff migration. On the other hand, the computational effort is higher because a simulation of the wave equation has to be done for the source and the receiver data followed by the correlation and stacking of the matching source and receiver domains (imaging condition). This means that the computational effort is especially critical with RTM algorithms.

Depending on the complexity of the velocity model RTM algorithms come at least in three flavours, the isotropic (ISO), the vertical transverse isotropy and the tilted transverse isotropy (TTI). The latter two have an anisotropy axis which can be tilted versus the vertical axis in TTI.

Especially the ISO and VTI form are typically used in practical applications. These algorithms are mostly memory bound. For this reason, increasing the pure FLOP rate is not helpful to decrease the time to solution. On the other hand, decreasing the precision of the Floating Point format might exploit better the available memory bandwidth and increase the FLOP rate under constant or even shrinking power envelope. The question to explore is if lower precision is sufficient to keep up the image quality of RTM migration algorithms calculated in full precision. As the typical Floating Point format in seismic is 32 bit IEEE Float we evaluate different 16 bit Floating Point formats. Among them 16 bit Posit format and 16 bit Floating Point format.

Previous experiments have demonstrated that RTM images calculated in reduced precision deliver competitive images compared with calculations in 32 bit Float. Simple stability tests also demonstrated positive outcomes. However, these images were based on a single shot and were far from modeling a realistic seismic scenario. So the next step is to model a more realistic case. To limit the calculation time to an applicable amount we switch to 2D images and use the well-known Marmousi model. The Marmousi model is an artificially created reference model based on a section in the Cuanza Basin. It is very often used in literature to compare the performance of different migration algorithms.

3.3.1 Design of experiments

A small reference implementation is created that is able to calculate the isotropic RTM Marmousi example. The source and receiver wavefields can be stored in different Floating Point standards, the Float32, the Float16 and the Posit format. The RTM kernel can be calculated in the same Floating Point format matching the storage or it can be calculated in Float32 while the Floating Point format is converted back and forth. Another part is the Floating Point format of the calculated Image from the two wavefields. This format can be identical to the reduced precision Floating Point format or it can be stored in Float32. The Floating Point format of the image can be mixed arbitrarily with the calculation format of the kernel.

Furthermore, the calculation of the imaging condition in reduced precision can be done including the Kahan summation. The image is calculated by calculating the correlation between the source wavefield und the receiver wavefield. The different correlations which are in the order of magnitude of the number of calculated timesteps are summed up to yield the final image. The precision of this sum can be improved using the Kahan summation. The advantage of the Kahan summation is a higher precision than using normal summation. The disadvantage of the Kahan summation is that two numbers per pixel are needed. Thus the bandwidth needed to store the image in 16 bit format is the same as the number of bit to store a 32bit image in. The Floating Point effort is even higher than calculating in Float32. However, this approach can be useful because a very lean compute core can be used to calculate the imaging condition having only ALUs in reduced precision.

3.3.2 Results Posit format

The Marmousi model is reduced to save computational time. The model has 240 shots that are reduced to 40 shots around the center (shots #101 to #140). Furthermore, the simulated time is reduced from 2.7 seconds to 1.62 seconds. The grid spacing is 16m in depth direction and 25m in lateral direction. The timestep is 1.405 milliseconds. The number of pixels in the domain is

185x223 (lateral times depth). The full domain would be 367x187 but is adjusted in these experiments. As only 40 shots are calculated the simulated domain can be reduced in lateral direction. In depth direction that size has been slightly extended beyond the full domain. The reason is that no specific boundary condition has been implemented. Thus, padding the domain reduces the amount of artifacts in the images created by reflected waves at the artificial boundary of the domain. The sources and receivers are coupled in 18 pixels beyond the upper boundary. In the calculation campaign different combinations of Floating Point formats in the domain, the kernel and the image are tested. All the Posit16 values use one exponent. The source signal is modeled as a Gabor wavelet.

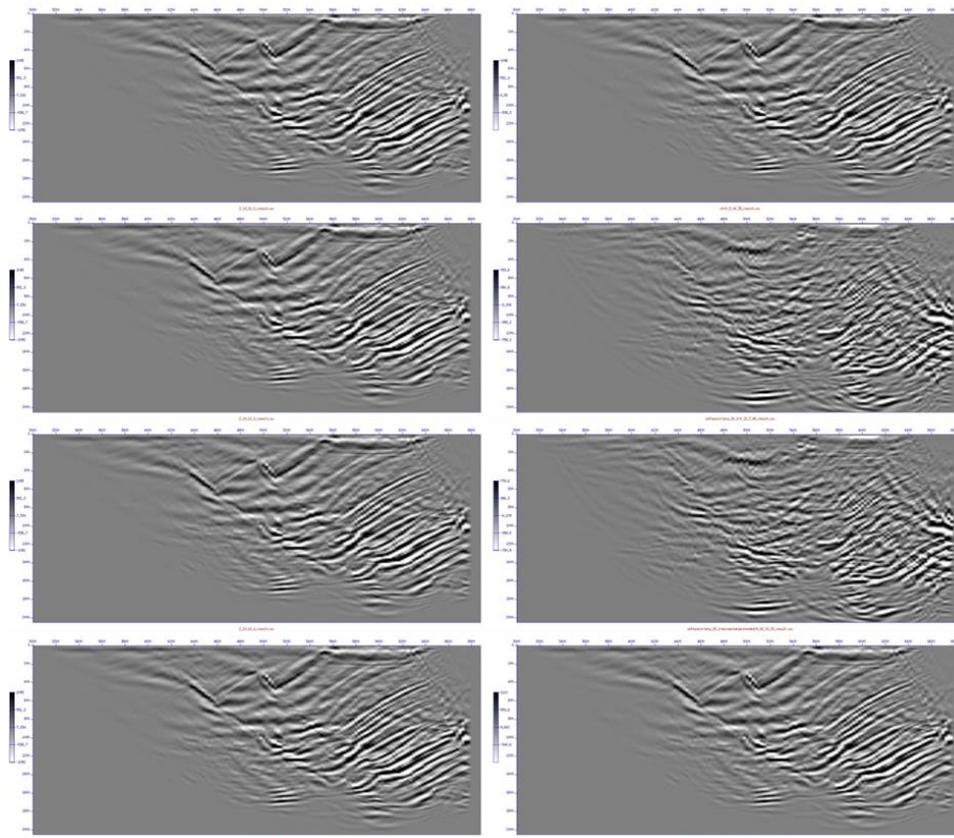


Figure 13: [RTM] Marmousi results with different combinations of Floating Point formats

Left: Reference simulation in Float32 Right: specific combination of Posit with Float32 formats

The combinations are from top to bottom:

Label “p16.1_storage_only”: Domain Posit, kernel Float32, image Float32

Label “p16.1”: Domain Posit, kernel Posit, image Float32

Label “p16.1_accumulate”: Domain Posit, kernel Posit using quire for dot products, image Float32

Label “p16.1_accumulate_scale”: Domain Posit, kernel Posit using quire while scaling model and traces, image Float32

Figure 13 and Figure 14 demonstrate that using Posit 16 with one exponent as a storage format for the domain only creates very similar images compared to the pure Float32 reference image. Using Posit for calculation of the kernel too yields unacceptable results (p16.1). The reason for the poor performance might be catastrophic cancellation in calculation of the stencil inside the kernel. To mitigate a possible cancellation a Posit quire is used to sum up the stencil. The quire has a high number of bits as an intermediate storage of the sum to avoid cancellation in intermediate results. Only the final result is rounded to the available mantissa. Image “p16.1_accumulate” demonstrates that the quire doesn’t improve the image quality significantly. Another source of low precision could be a low mantissa because of large Floating

Point values. The mantissa of Posit values decreases the more the absolute value deviates from 1. Figure 15 demonstrates the distribution of absolute values in the velocity model and the receiver data. The median of the velocity model is 2638.28. The average of the two medians (negative values and positive values) is 30.981. The receiver data and the velocity data are scaled so that the median values map to 1. The output images are scaled back so that the scaling is cancelled out. Simulation “p16.1_accumulate_scale” demonstrates that the scaling fixes the broken image. In Figure 13 simulation “p16.1_scale” demonstrates that the scaling alone is sufficient to create a good image and that the quire is not needed. Simulation ”p16.1_scale_img_p16.1” and ”p16.1_scale_img_p16.1_kahan” demonstrate that an image stored in Float16 doesn’t make the image unusable

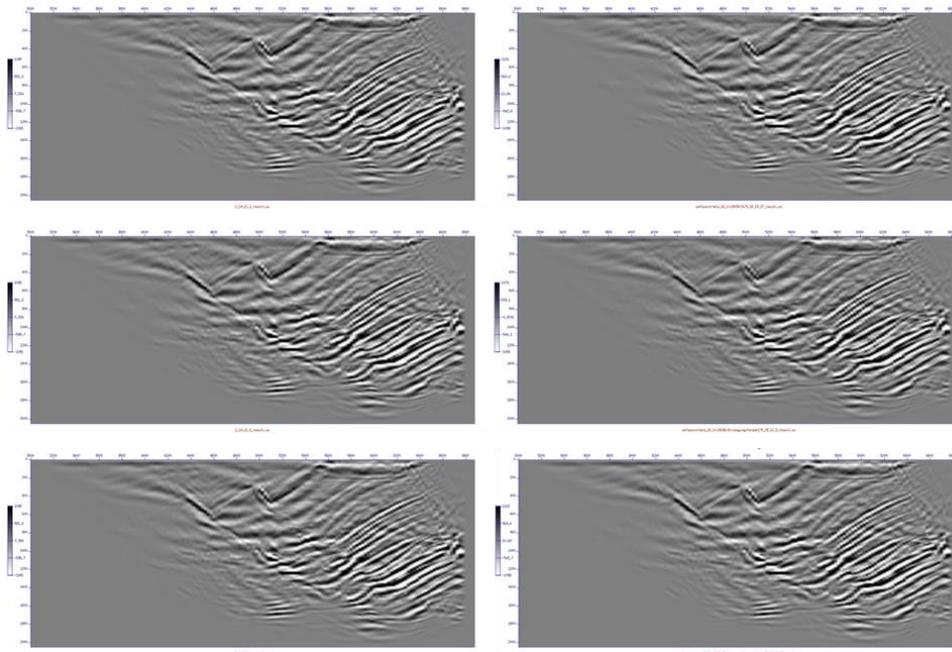


Figure 14: [RTM] Marmousi results with different combinations of Floating Point formats continued
 Left: Reference simulation in Float32 Right: specific combination of Posit with Float32 formats
 The combinations are from top to bottom:
 Label “p16.1_scale”: Domain Posit, kernel Posit scaling model and traces, image Float32
 Label ”p16.1_scale_img_p16.1”: Domain posit, kernel Posit scaling model and traces, image Posit
 Label ”p16.1_scale_img_p16.1_kahan”: Domain posit, kernel Posit scaling model and traces, image Posit Kahan

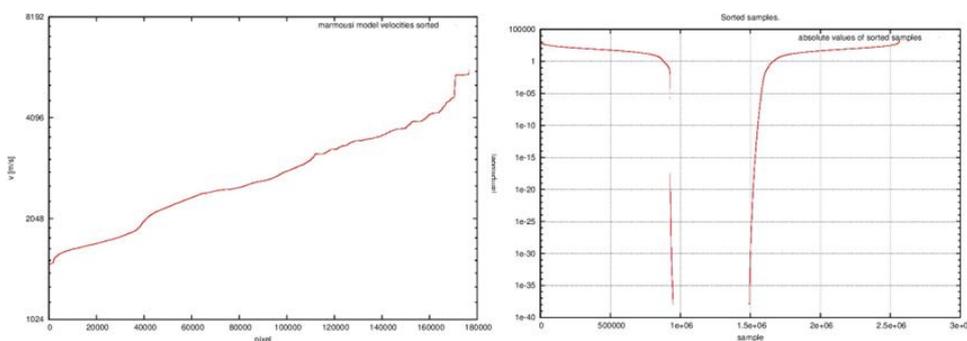


Figure 15: [RTM] Distribution of absolute values in velocity model (left) and all the receiver traces (right)
 Input data is sorted ascending for size and plotted as absolute value.

Figure 16 plots the infinity norm of the difference between Float32 and Posit normalized by the infinity norm of the Posit image. The simulation using Posit as storage format for the image deviates significantly more from the reference than all the other useful cases. The usage of Kahan summation brings back the deviation to similar values as other useful cases (approx. 4%).

To have close to zero difference between Float32 and Posit is not necessarily crucial. Most important is to have no shifts of the reflectors especially in-depth direction. On the other hand, a difference close to zero implies few shifts. But although the case "p16.1_scale_img_p16.1" has significantly higher deviation towards Float32 as e.g. "p16.1_scale" the eye cannot see relevant shifts of reflectors in both cases.

3.3.2.1 Numerical Stability

To save computational time the simulation has been restricted to a subset of the model. Especially the simulated time has been reduced. To verify that the Marmousi example is stable also for longer simulated times the simulated times was increased in a further reduced test example. Here the development of the total energy of the source wavefields is plotted over time. A proxy for the total energy is the summed squares of all the pixels. According to physics, a conservation of energy is expected over time. So, this property is a must have of the algorithm and a violation of this property reveals numerical instability.

Figure 17 demonstrates a significant increase of the energy beyond 1.6 seconds simulated time. The same simulation with 17 bit Posits demonstrates that this explosion of energy can be avoided by using a slightly higher precision. This means that the precision of the Posit 16 bit is very close to sufficient to avoid this explosion. In the end a reorganization of the kernel was sufficient to stabilize the algorithm even if 16 bit Posits.

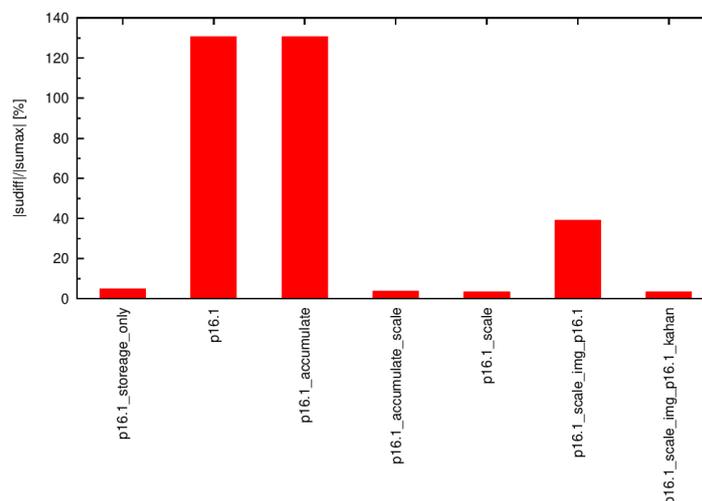


Figure 16: [RTM] Distribution of absolute values in velocity model (left) and all the receiver traces (right)

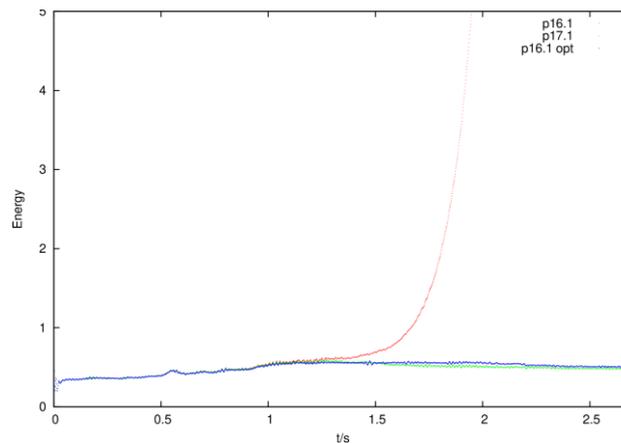


Figure 17: [RTM] Distribution of absolute values in velocity model (left) and all the receiver traces (right)
 The energy of shot number 120 of 240 in the Marmousi model is plotted over time. The Floating Point format is Posit with 16 bit (red), 17 bit (green) and 16 bit with reorganized kernel (blue). Always one exponent is used. Parameters: Grid spacing is 25x16 meters, timestep is 1.405ms, the simulated time is from 0 to 2.7s. The number of pixels is 185x223 for the red and the green curve and 185x373 for the blue curve.

3.3.2.2 Conclusion

To use the Posit 16 one exponent Floating Point to store the source and receiver wavefields delivered acceptable images more or less out of the box. Using the same format for calculation of the kernel raised some issues. These issues could be solved by scaling the model and the receiver samples combined with a numerical reorganization of the kernel. Using the Posit as format to aggregate the correlated images too increases the deviation towards Float 32 bit significantly but the results are still acceptable. Additionally, the Kahan summation can be used to rebuild a low deviation towards Foat32. However, this mitigates the bandwidth savings of the reduced precision format. On the other hand this is only a minor drawback as the number of memory accesses to the image during the stacking process is small compared to the memory accesses to the wavefields during the stacking process and during the kernel calculation. The advantage to use Kahan summation might be that the compute core can be designed more simply by dropping the support for Float32 completely.

3.3.3 Results Float16 format

The simulation in Float 16 bit Floating Point format is less compute time consuming so that the full Marmousi model can be calculated. The grid spacing is 16m in depth direction and 25m in lateral direction. The simulated time is 2.7 seconds, the timestep is 1.405 milliseconds. The number of pixels in the domain is 439x373 (lateral times depth). The full domain would be 367x187 but is padded to avoid artifacts caused by reflections on the artificial boundaries. The sources and receivers are coupled in at 93 pixels beyond the upper boundary. Sources and Receivers are scaled to avoid overflow of exponents in the Float 16 data type. The maximum representable value is 65504 in Float16. In the calculation campaign different combinations of Floating Point formats of the domain, the kernel and the image are tested. The source signal is modeled as a Gabor wavelet.

Figure 18 depicts the calculated images. The first model “f16” shows storage and calculation of the wavefields in Float16 while the image is stacked in Float32. The next case uses Float 16 bit for stacking of the image too while the case “f16_img_f16_kahan” uses Kahan summation additionally to stack the images. The remaining two cases store the wavefields in -Float16 and calculate the kernel in Float32. The “f16_storage_only” case stacks the image in Float32 while

the “f16_storage_only_img_f16” stacks the image in Float16. All the images show an acceptable quality. The eye can detect very subtle differences between images where the kernel has been calculated in Float32 compared to kernels calculated in Float16.

Figure 19 plots the relative norm of difference versus the pure Float32 case as reference. The cases where the kernel was calculated in Float16 are all very similar around 7% deviation. The deviation of the respective Posit16 cases towards Float32 (3%) was lower (better). Calculating the kernel in Float32 reduces the deviation to 0.5% or 4% if the image is stacked in Float16. Stacking the image in Float16 makes a significant difference here.

3.3.3.1 Numerical stability

To verify that the Marmousi example is stable while calculated in Float16 the development of the total energy of the source wavefields is observed over time (0s-2.7s). A proxy for the total energy is the summed squares of all the pixels. In none of the cases the plotted energy over simulated time did show any explosive increase (data not shown). So, no signs of instability have been detected.

3.3.3.2 Conclusion

All the different combinations of Float16 with Float32 delivered acceptable images and numerically stable simulations. If the kernel was calculated in Float16 the deviation versus the reference was double compared to the respective Posit16 calculations. On the other hand, if the kernel was calculated in Float32 the Float16 simulation had smaller deviations towards the reference. Scaling the input data was necessary with Float16 as well as with Posit16.

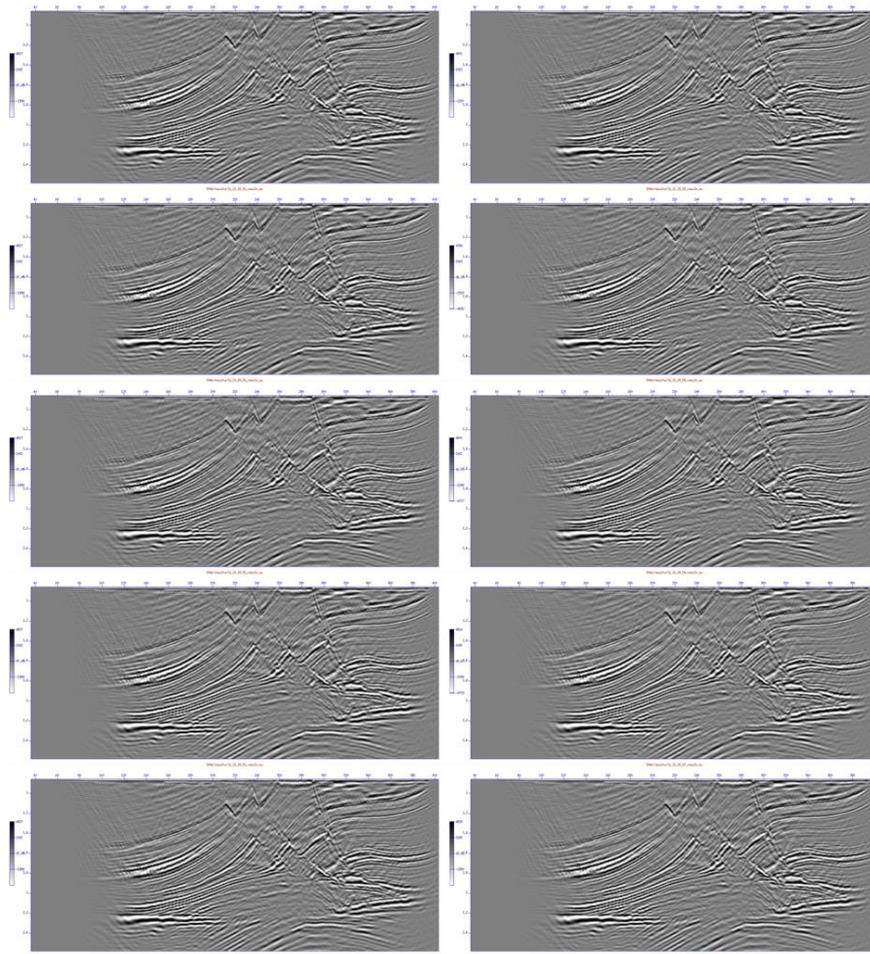


Figure 18: [RTM] Marmousi results with different combinations of Floating Point formats

Left: Reference simulation in Float32 Right: specific combination of Float16 with Float32 formats
The sources and receivers are scaled to avoid exponent overflow in Float16 values.

The combinations are from top to bottom:

Label "f16": Domain Float16, kernel Float16, image Float32

Label "f16_img_f16": Domain Float16, kernel Float16, image Float16

Label "f16_img_f16_kahan": Domain Float16, kernel Float16, image Float16 with Kahan summation

Label "f16_storage_only": Domain Float16, kernel Float32, image Float32

Label "f16_storage_only_img_f16": Domain Float16, kernel Float32, image Float16

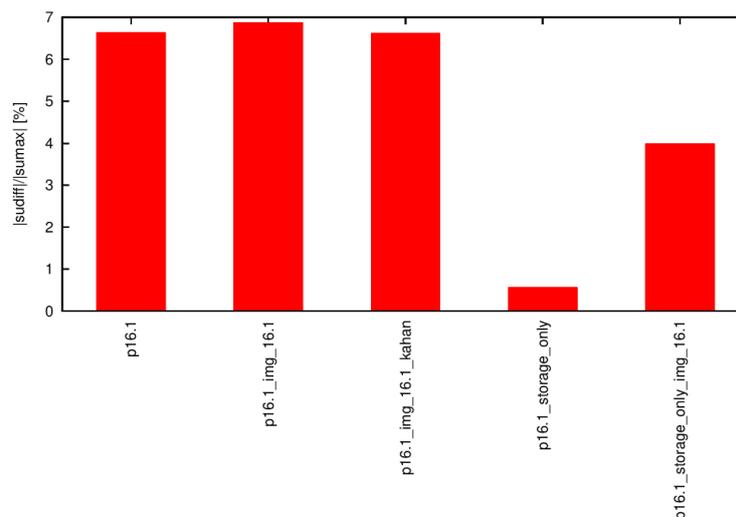


Figure 19: [RTM] Infinity norm of difference between Float16 and Float32 as percentage of the infinity norm of the Float16 example

3.4 HEP - INFN

High-Energy Physics (HEP) experiments are increasingly using a mix of processing hardware (heterogeneous architectures) to keep up with computing solutions offered by accelerator technologies.

The TEXTAROSSA project focused on two applications from CERN's Patatrack² team: Pixeltrack (a pixel tracking algorithm used for CMS detector) and CLUE (a cluster algorithm for high-granularity calorimeters). Both applications provide several portability layers. In this way it is possible to efficiently run a single codebase on various hardware types from different vendors, since recoding for each specific architecture is impractical.

Our objective was to evaluate the performance of these applications in two key areas, enabling the evaluation of the IDV architectures developed in the project:

- throughput: number of reconstructed events per second;
- energy efficiency: number of reconstructed events per Joule, obtained from the ratio between throughput and power.

As reported on deliverable D6.2, for the TEXTAROSSA project we concluded that while SYCL is a promising approach, Intel's oneAPI implementation isn't stable or mature enough for our needs. Specifically, we require running the applications used to collect and process data from HEP detectors on various mixed hardware platforms. For that reason, we opted to use a different abstraction layer called Alpaka [1] to collect data for the TEXTAROSSA project. As quoted from the Alpaka documentation: "The Alpaka library is a C++14 header-only abstraction library designed for accelerator development. Its goal is to ensure performance portability across accelerators by abstracting (not hiding!) the underlying layers of parallelism." While coding with Alpaka is more complex than with SYCL, this library supports a wider range of devices from different vendors compared to oneAPI, making it better suited for our project.

² <https://patatrack.web.cern.ch/patatrack/>

Therefore, we opted to conduct all measurements leveraging the Alpaka version of the application, in this way we were able to run the same source code on various devices without any modifications. We used the applications to characterize the IDV-A (CPU only, CPU+GPU) and IDV-E (CPU only) nodes' performance.

3.4.1 Baseline KPIs assessment on Dibona node

We repeated the tests and measurements on the Dibona node, that were already reported in Deliverable 6.2, to be more confident in assessing the baseline for the applications' KPIs to be used in the comparison with the IDV-A and IDV-E nodes' performance.

3.4.1.1 Tests on CPU

Our focus on the CPU is twofold: firstly, to verify that our code efficiently scales increasing the number of processing cores. Secondly, we want to see how energy efficiency is affected by this scaling. We achieved this by running both applications with a varying number of cores and threads. To measure power consumption, we employed the *likwid-perfctr* tool. This tool also assisted us in accurately specifying the number of cores to be used.

The number of threads could be indicated in the command we use to run the application, instead. In the same command we also specified how long the applications must run. Every test was 2-minute long and the *likwid-perfctr* tool extracted the power consumption values every second. An example of one of those execution command is:

```
likwid-perfctr -f -C S<socket_id>:<cores> -g ENERGY -t 1s -O -o
<output-file>.csv ./alpaka --serial --runForMinutes 2 --numberOfThreads
<threads>
```

The results we achieved are reported below. Table 8 and Table 9 show the measures of throughput (events/second), Power (W), Energy Efficiency (events/J) and the Energy Efficiency wrt single core (i.e., normalized to the events/J obtained on a single core and measured in a.u.) for the CLUE application and for the Pixeltrack application when they run on CPU, scaling the number of logic cores.

cores	thread	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	1	4.341	81.114	0.053	1.0
2	2	8.574	85.074	0.100	1.748
6	6	16.878	92.266	0.183	3.077
12	12	33.720	107.233	0.314	5.169
16	16	66.586	136.053	0.489	8.049
24	24	98.840	166.268	0.594	9.676
48	48	189.239	162.402	1.165	18.980
64	64	232.670	172.709	1.347	21.917
72	72	253.44	181.032	1.400	22.777
96	96	296.225	173.950	1.703	27.745

Table 8: [HEP] KPIs for CLUE on Dibona CPU

cores	thread	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	1	30.118	82.116	0.367	1.0
2	2	60.702	86.439	0.629	1.890
6	6	178.051	103.708	1.717	4.354
12	12	350.269	129.147	2.712	6.835
24	24	705.473	182.601	3.863	9.556
48	48	1365.25	177.424	7.695	13.306
64	64	1573.48	190.017	8.281	20.419
72	72	1604.44	192.939	8.316	20.535
96	96	1353.11	178.625	7.575	18.733

Table 9: [HEP] KPIs for Pixeltrack on Dibona CPU

The scaling of throughput for both applications is illustrated on Figure 20 and Figure 21.

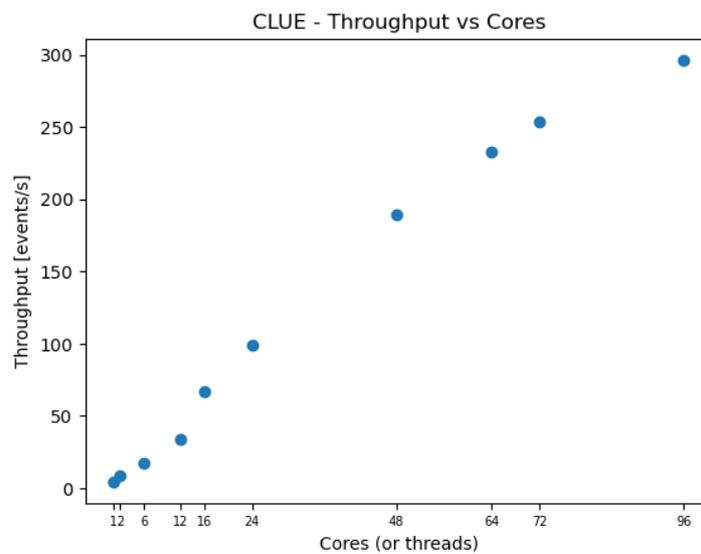


Figure 20: [HEP] Throughput vs cores for CLUE on Dibona CPU

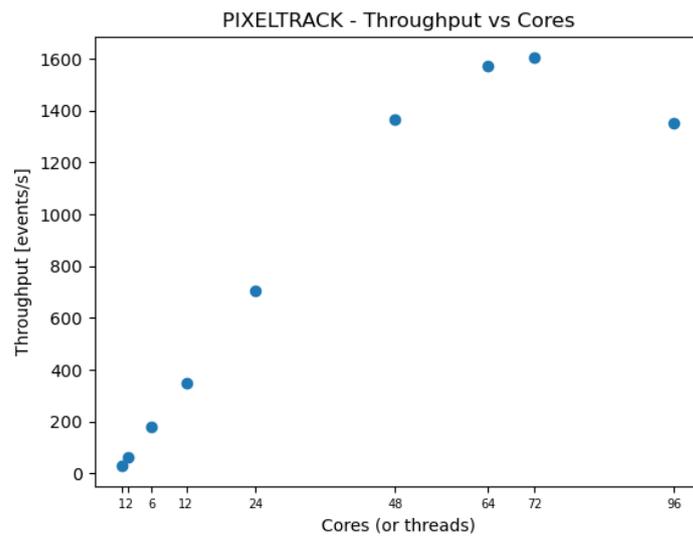


Figure 21: [HEP] Throughput vs cores for Pixeltrack on Dibona CPU

Figure 22 and Figure 23 explore the relationship between the number of cores used and the applications' energy efficiency. It's important to note that the data in both these figures is normalized, meaning it's been adjusted relative to the performance achieved with a single thread on a single core.

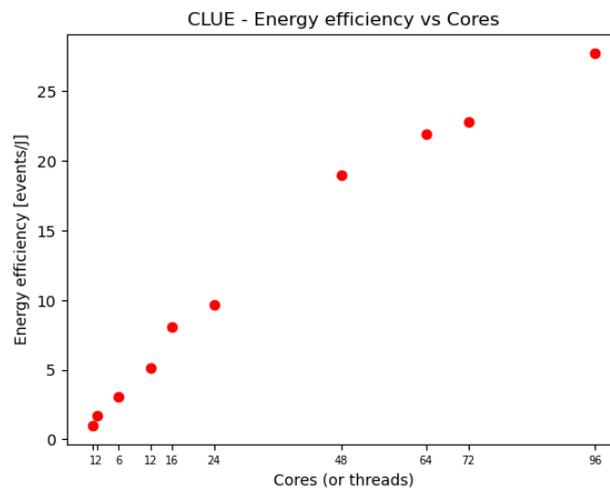


Figure 22: [HEP] Energy efficiency vs cores for CLUE on Dibona CPU

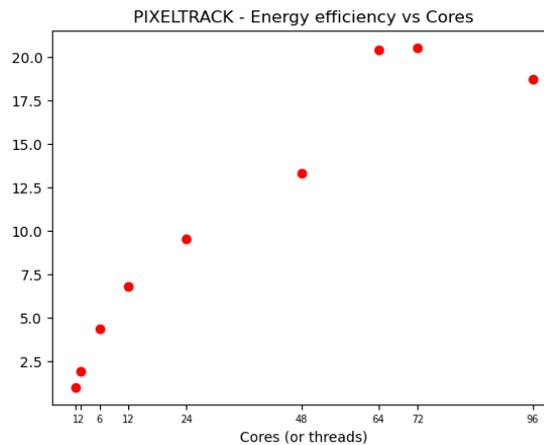


Figure 23: [HEP] Energy efficiency vs cores for Pixeltrack on Dibona CPU

Figure 20 and Figure 21 show that the throughput KPI grows linearly when we run one thread per core using one socket (up to 48 cores). This trend remains almost constant for both the applications, but for the Pixeltrack algorithm, above 72 threads, the performance gets worse. A similar behaviour can be noticed in Figure 22 and Figure 23, showing the scaling of the Energy Efficiency KPI with the number of CPU cores used by the applications. This difference might be caused by the Pixeltrack implementation possibly requiring more memory accesses, incurring in resource contention.

3.4.1.2 Tests on GPU

Our focus on the GPU is similar to the CPU analysis: demonstrating how well our code utilizes an increasing number of GPUs and how energy efficiency is affected.

To achieve this, we fixed the number of CPU threads per GPU at 12, ensuring each thread runs on a separate core. This ensures all 48 physical cores of a single CPU are utilized when using 4 GPUs. As with the CPU tests, the number of CPU threads can be specified in the application's launch command. Additionally, this command allows configuring the mapping between GPUs, CPU threads, and cores. All our tests ran for a duration of 2 minutes.

For illustration purposes, here's an example command used to run the application on 2 GPUs:

```
(CUDA_VISIBLE_DEVICES=0 taskset -c 0-11 ./alpaka --cuda --runForMinutes 2 --numberOfThreads 12) & (CUDA_VISIBLE_DEVICES=1 taskset -c 12-23 ./alpaka --cuda --runForMinutes 2 --numberOfThreads 12)
```

For the power consumption measurements, we used the GPowerU tool.

The results we achieved are reported below. Table 10 and Table 11 show the measures of throughput (events/second), Power (W) and Energy Efficiency (events/J) for the CLUE application and for the Pixeltrack application when they run on one or more GPUs.

CLUE on GPU – A100				
GPU	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)

1	1401.940	175.500	7.991	1.0
2	2937.820	368.261	7.979	0.998
3	4369.630	551.833	7.922	0.991
4	5651.190	728.849	7.756	0.970

Table 10: [HEP] KPIs for CLUE on GPU on Dibona node

Pixeltrack on GPU – A100				
GPU	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	2174.430	149.333	14.59	1.0
2	4302.790	300.804	14.34	0.983
3	6370.730	445.749	14.29	0.980
4	8591.570	595.272	14.43	0.990

Table 11: [HEP] KPIs for Pixeltrack on GPU on Dibona node

Figure 24 and Figure 25 show for both applications the close-to-perfect linear scaling of throughput KPI with the number of the used GPUS on the Dibona node.

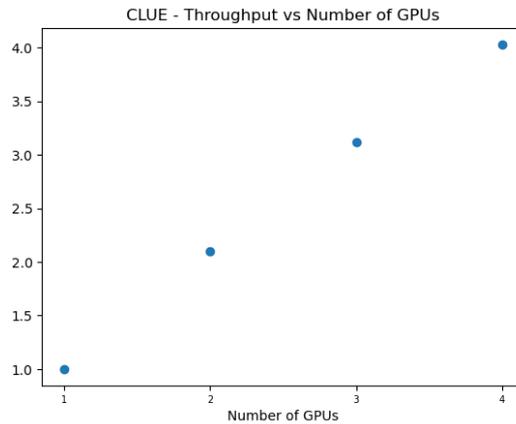


Figure 24: [HEP] Throughput (a.u.) vs number of used GPUS for CLUE on Dibona node

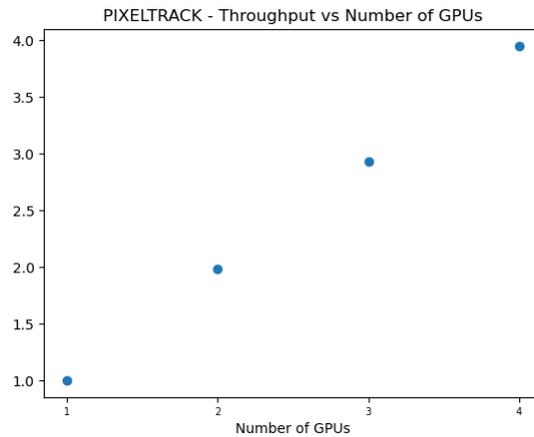


Figure 25: [HEP] Throughput (a.u.) vs number of used GPUs for Pixeltrack on Dibona node

Figure 26 and Figure 27, on the other hand, show for both applications an almost constant trend of the energy efficiency KPI (normalized to what is obtained on one GPU) scaling the number of the used GPUs on the Dibona node from one to four.

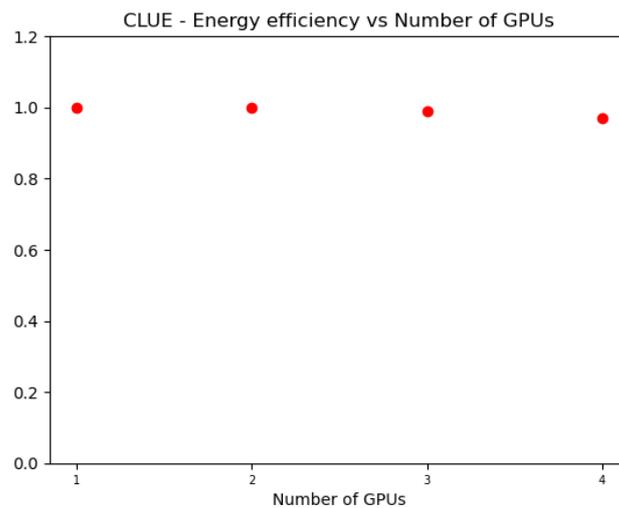


Figure 26: [HEP] Energy efficiency (normalized to one GPU) vs number of used GPUs for CLUE on Dibona node

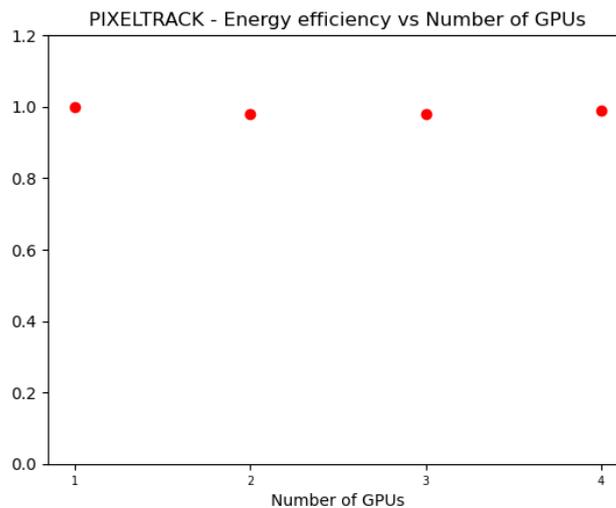


Figure 27: [HEP] Energy efficiency (normalized to one GPU) vs number of used GPUs for Pixeltrack on Dibona node

3.4.2 Tests on IDV-A

Being available the IDV-A, a dual-socket server equipped with two Intel(R) Xeon(R) Platinum 8470 (Sapphire Rapids generation with 52 *Physical* cores Total Threads.) and 4 NVIDIA H100-64MB, we had the chance to test the performance of the CLUE and Pixeltrack benchmarking applications on this node implemented with more recent technologies. In this section we report the results of the measured applications' KPIs on the IDV-A node, for both CPU only and CPU+GPU configurations.

3.4.2.1 Tests on CPU

In Table 12 and Table 13 we report the measurements of throughput (events/second), Power (W), Energy Efficiency (events/J) and Energy Efficiency normalized to single core (a.u.) for both CLUE and Pixeltrack applications run on their CPU versions.

cores	thread	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	1	3.198	405.81	0.008	1.0
2	2	6.364	408.87	0.015	1.975
6	4	18.865	417.10	0.045	5.740
12	8	37.462	428.19	0.088	11.102
24	16	75.028	451.89	0.166	21.068
52	52	161.582	504.82	0.320	40.616
64	64	196.779	529.61	0.371	47.148
72	72	220.188	546.27	0.403	51.148
96	96	314.112	607.79	0.517	65.580

Table 12 [HEP] KPIs for CLUE on CPU for IDV-A node

cores	thread	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	1	21.469	407.23	0.053	1.0
2	2	44.915	409.62	0.110	2.080
6	6	132.096	420.32	0.314	5.961
12	12	261.728	431.91	0.606	11.494
24	24	527.451	466.24	0.884	21.458
52	52	1138.95	535.87	2.125	40.316
64	64	1376.83	568.78	2.420	45.916
72	72	1528.82	591.81	2.583	49.001
96	96	2060.88	672.98	3.062	58.086

Table 13: [HEP] KPIs for Pixeltrack on CPU for IDV-A node

Figure 28 and Figure 29 show, as in the case of the Dibona node, a close-to-perfect linear trend of the throughput KPI with the scaling of the number of used CPU cores, with a marginal improvement in absolute terms on IDV-A for this KPI for the CLUE application, while there is a noticeable improvement in scaling approaching the maximum number of available cores for the Pixeltrack application, with a better scaling behaviour on IDV-A and a noticeable improvement in the maximum attainable throughput in the new architecture.

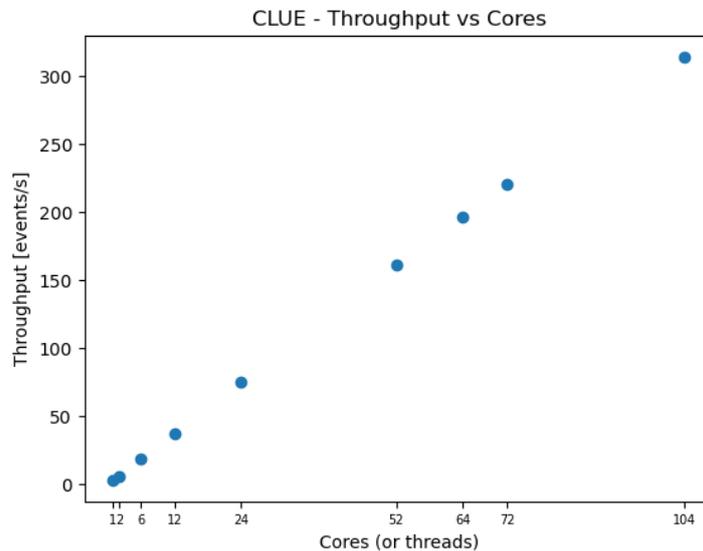


Figure 28: [HEP] Throughput vs number of used cores for CLUE on IDV-A CPU

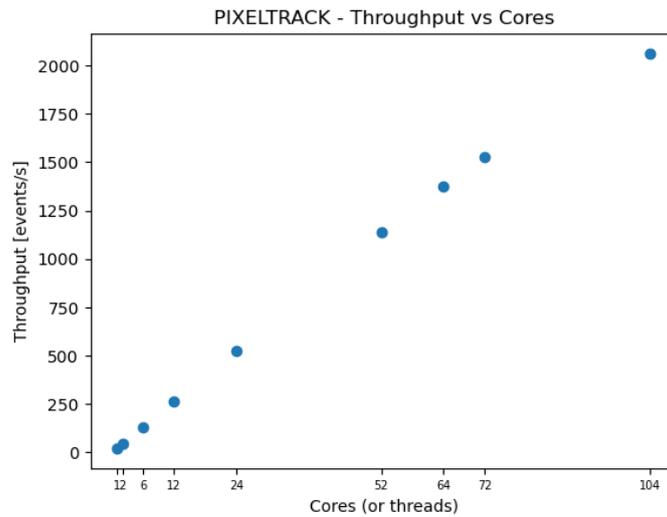


Figure 29: [HEP] Throughput vs number of used cores for Pixeltrack on IDV-A CPU

The scaling of the energy efficiency with the number of used CPU cores for the CLUE application on the IDV-E CPU (Figure 30) shows a similar linear trend to that on the Dibona CPU (Figure 22). Instead, for what regards the scaling on the energy efficiency on the IDV-A CPU of the Pixeltrack application, Figure 31 shows a quite relevant improvement of the IDV-A CPU with respect to the corresponding result for the Dibona node (Figure 23). Focusing on the absolutes of the energy efficiency KPIs shows a significant degradation of performance on IDV-A though. This is mainly due to the ~400 W power consumption in idle state of the node CPUs, with a significant contribution accountable to the thermal control daemon (in its current version) for the two-phase cooling system.

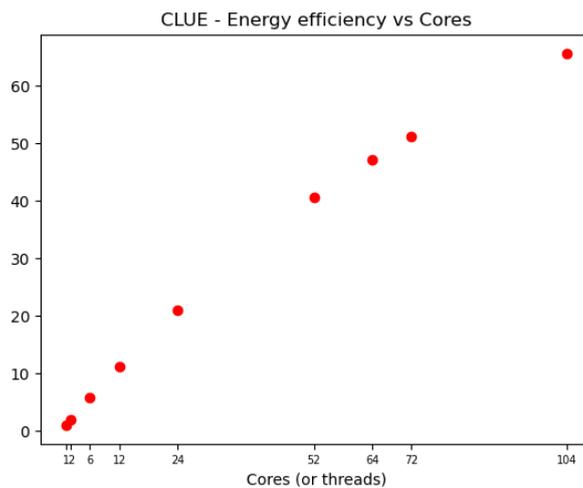


Figure 30: [HEP] Energy efficiency vs number of used cores for CLUE on IDV-A CPU

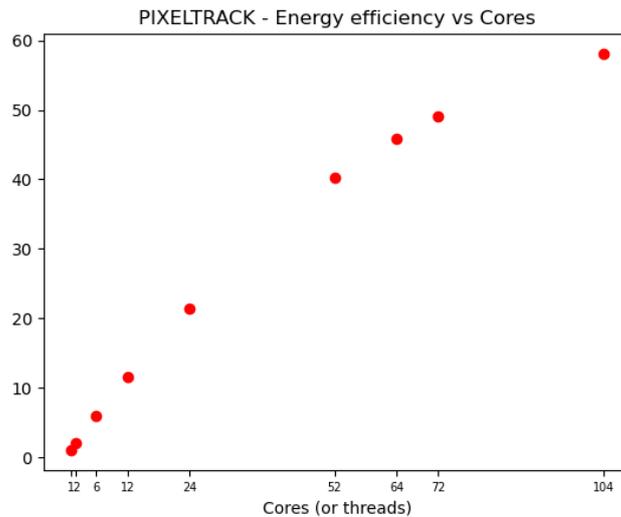


Figure 31: [HEP] Energy efficiency vs number of used cores for Pixeltrack on IDV-A CPU

3.4.2.2 Tests on GPU

We run the CLUE and Pixeltrack applications on the IDV-A node in the same configurations used to measure the baseline KPIs on the Dibona node and described in Sec.3.4.1.2 .

Measured KPIs scaling the number of used GPUs are reported below.

CLUE on GPU - H100				
GPU	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	1102.1	181.40	6,075	1.0
2	2194.63	364.50	6,021	0.991
3	3314.11	552.47	5,999	0.987
4	4300.63	728.38	5.90	0.971

Table 14: [HEP] KPIs for CLUE on GPU on IDV-A node

Pixeltrack on GPU - H100				
GPU	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	2507.31	182.678	13,725	1.0
2	4999.05	367.347	13,608	0.991
3	7212.48	546.521	13,197	0.961
4	9514.89	723.583	13,150	0.958

Table 15: [HEP] KPIs for Pixeltrack on GPU on IDV-A node

Moving from Dibona node (Table 11) the improvement in Pixeltrack throughput KPI is noticeable and expected in the IDV-A node (Table 15). The slightly worse performance of CLUE on IDV-A on throughput KPI (Table 14) compared to that measured on Dibona node (Table 10) needs further investigation and code tuning for the Nvidia H100 GPU.

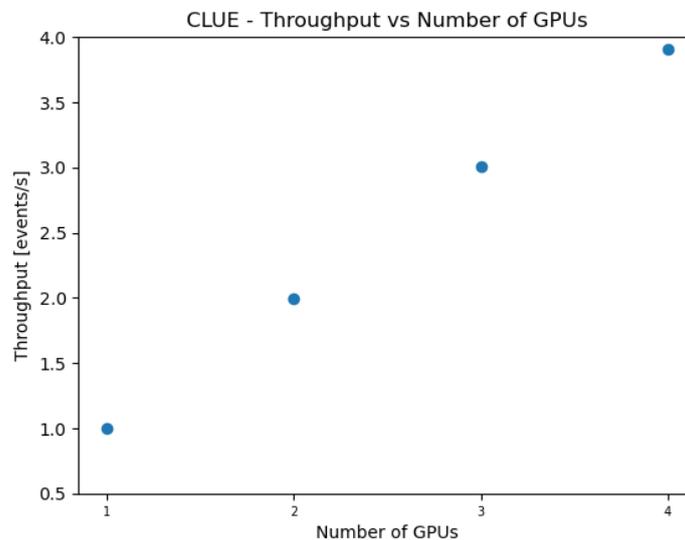


Figure 32: [HEP] Throughput vs number of used GPUs for CLUE on IDV-A node

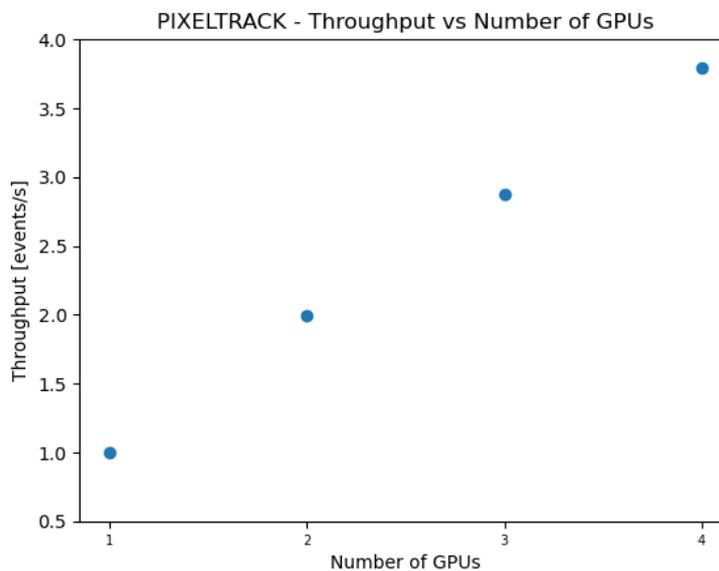


Figure 33: [HEP] Throughput vs number of used GPUs for Pixeltrack on IDV-A node

Figure 34 and Figure 35, on the other hand, show for both applications a slightly decreasing trend of the energy efficiency KPI (normalized to what is obtained on one GPU) scaling the number of the used GPUs on the IDV-A node.

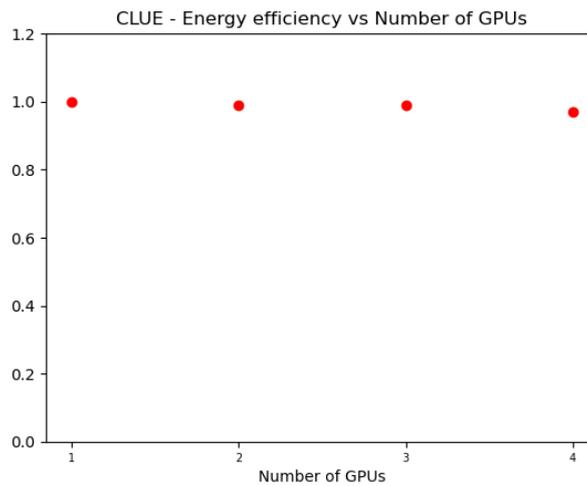


Figure 34: [HEP] Energy efficiency vs number of used GPUs for CLUE on IDV-A node

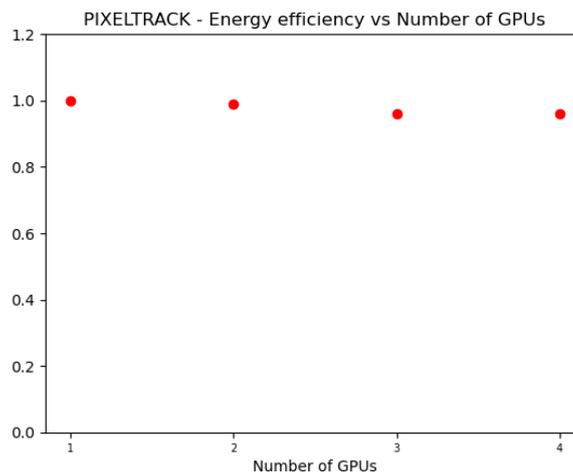


Figure 35: [HEP] Energy efficiency vs number of used GPUs for Pixeltrack on IDV-A node

3.4.3 Tests on IDV-E

Having the opportunity to run the benchmarks on the IDV-E nodes equipped with ARM64 Ampere Altra CPUs, we can report the performance of the aforementioned Pixeltrack and CLUE applications on this server.

Pixeltrack on ARM64					
cores	thread	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	1	24.863	64.80	0.383	1.0
2	2	58.582	63.38	0.924	2.409
4	4	112.156	71.35	1.572	4.097
8	8	212.776	74.65	2.850	7.429

16	16	410.435	89.20	4.601	11.992
32	32	778.272	109.64	7.098	18.500
64	64	1339.22	149.24	8.973	23.387
96	92	1691.69	166.86	10.138	26.423
128	128	1842.98	195.75	9.415	32.955
160	160	2011.04	212.74	9.453	24.637
192	192	2200.17	223.35	9.851	25.673
224	224	2017.87	208.94	9.657	25.170
256	256	2016.09	220.33	9.150	23.848

Table 16: [HEP] KPIs for Pixeltrack on CPU on IDV-E node

CLUE on ARM64					
cores	thread	Throughput (events/sec)	Power (W)	Energy efficiency (events/J)	Energy efficiency (ref single core)
1	1	5.340	55.73	0.10	1.0
2	2	10.618	58.61	0.18	1.890
4	4	21.015	64.07	0.33	3.422
8	8	41.656	69.72	0.60	6.235
16	16	79.596	80.77	0.99	10.283
32	32	139.758	99.32	1.41	14.683
64	64	186.721	120.24	1.55	16.204
96	92	186.628	134.48	1.39	14.481
128	128	180.006	153.75	1.17	12.217
160	160	302.695	195.35	1.55	16.170
192	192	367.203	221.70	1.66	17.283
224	224	363.727	241.47	1.51	15.718
256	256	355.141	257.32	1.38	14.401

Table 17: [HEP] Throughput vs number of used cores for CLUE on IDV-E CPU

The scaling of throughput with the number of core used for Pixeltrack (**Errore. L'origine riferimento non è stata trovata.**Figure 36) and for CLUE (Figure 37) shows that the optimal working point for both applications is reached using 192 out of the 256 cores available. This corresponds exactly also to the absolute maximum of the energy efficiency KPI for the CLUE application, as shown inFigure 39, and to maximum value attainable when using both IDV-E CPUs (more than 128 cores) for the Pixeltrack application, as plotted in Figure 38.

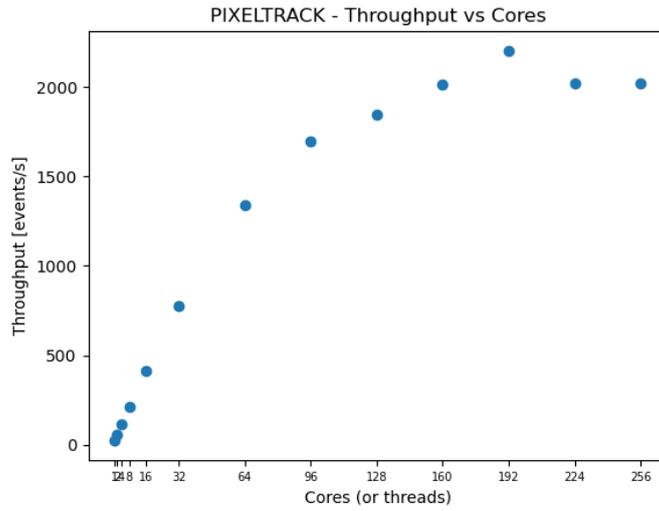


Figure 36: [HEP] Throughput vs number of used cores for Pixeltrack on IDV-E CPU

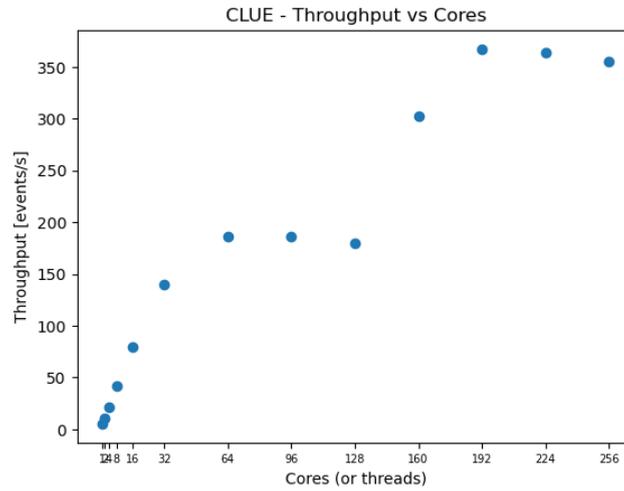


Figure 37: [HEP] Throughput vs number of used cores for CLUE on IDV-E CPU

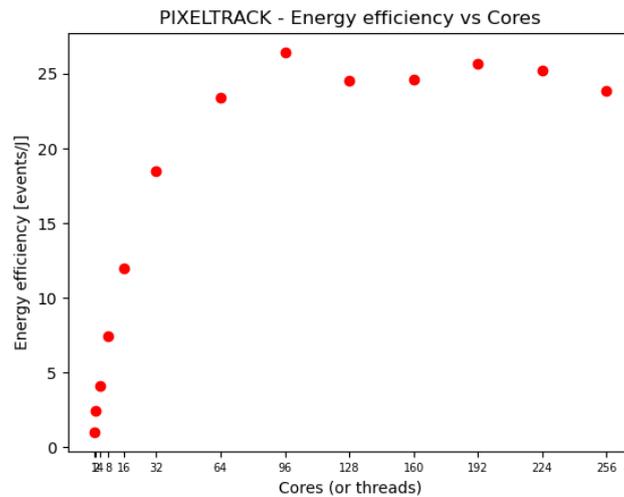


Figure 38: [HEP] Energy efficiency vs cores for Pixeltrack on IDV-E CPU

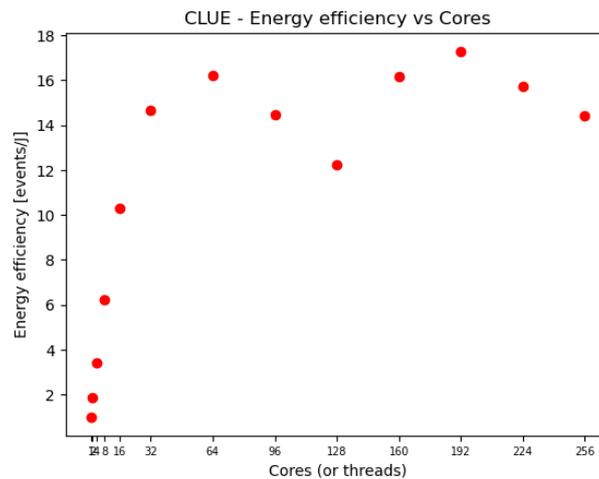


Figure 39: [HEP] Energy efficiency vs cores for CLUE on IDV-E CPU

For both applications, the IDV-E ARM64 architecture reaches, or surpasses, the throughput of the more performant of the two X86_64 ones, the IDV-A. Pixeltrack implementation surpasses by roughly a 50% the more energy efficient of them (Dibona) in the energy efficiency KPI while CLUE shows similar results.

3.5 NEST-GPU – INFN

As regards performance benchmarking in neuronal network simulations we mention the ongoing effort undertaken by the NEST team in the way of a standardization process [8] that designed a conceptual, generic workflow and produced a reference implementation, the beNNch framework. This is a set of modules for configuration, execution and analysis of benchmarks for NN simulations recording data and metadata in a unified way to foster reproducibility. Among the different modules we pick the test called *hpc_benchmark* [9]: it is a weak scaling benchmark employing a two populations network of excitatory and inhibitory leaky integrate-and-fire neurons, each with a fixed number of randomly drawn incoming connections (independent of the network size). This already has a multi-process implementation in the NEST repository as a Python script and we use an adapted version to perform the same test using the NEST-GPU application as a driver for the GPU as simulation engine. We mention in passing that the NEST-GPU has also received a new, GPU-accelerated implementation of the setup phase where the data structures pertaining to connections among neurons are created directly in GPU memory – a phase which at the time of D6.2 was still performed on CPU only and was quite time-consuming – that is described and benchmarked in detail in [10].

The power measurements were performed on the IDV-A platform using the GPowerU tool developed in TEXTAROSSA to encapsulate the launch of the *hpc_benchmark* script when run on the NVIDIA H100 GPUs (the same was done on the A100 GPUs equipped in the TEXTAROSSA partition in Dibona cluster for reference) and on the IDV-E platform by reporting at the same time the output of the *sensors* command – that returns the instantaneous power per CPU socket – together with the output of a Power Distribution Unit connected to the IDV-E that returns the whole power draw.

An interesting artifact of sampling the power output of an idle GPU with GPowerU is shown in Figure 40; here it can be clearly seen that the wattage sensor output is in discrete steps between 0.061W and 0.064W and that idle power draw of a GPU hovers on a value over around 67W.

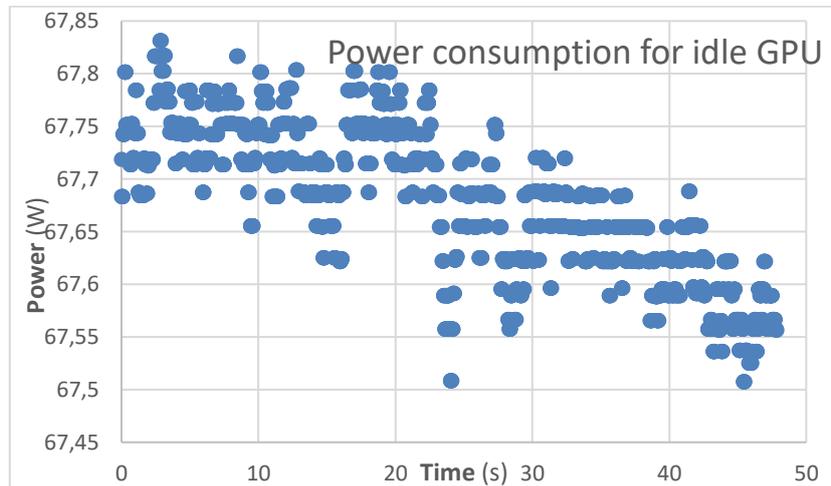


Figure 40: [NEST-GPU] Example of idle GPU power envelope

The test is configured to set up the network and perform a ‘warm-up’ 50ms pre-simulation in order to put the system into a steady state, followed by a 10.000ms simulation which is then timed; the results are in the following plots.

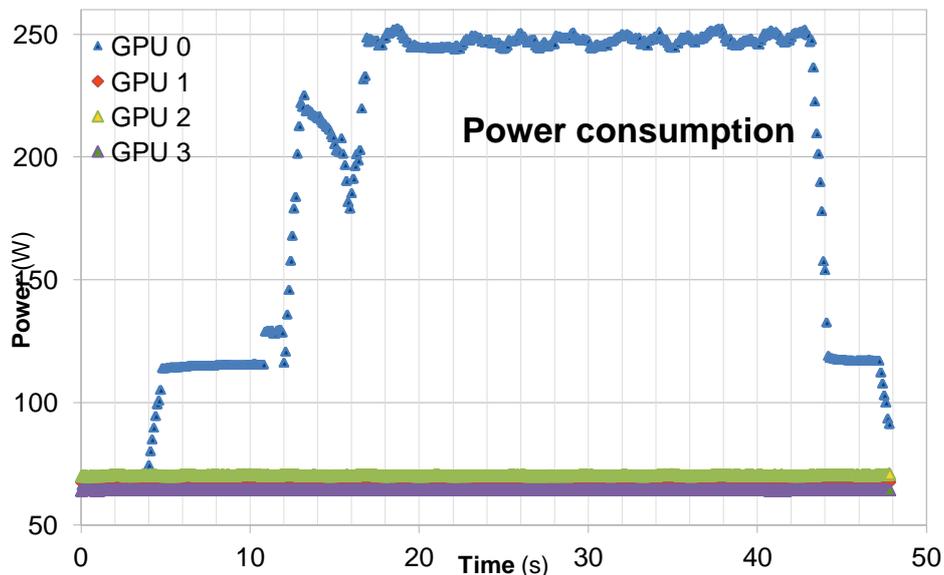


Figure 41: [NEST-GPU] Single active GPU power reading on IDV-A

From the plot in Figure 41 we can see that there is a network building phase – the first plateau – that pulls the power draw up from the idle state, then a ‘pop’ for the very first actual simulated milliseconds – 50ms which are required to put the network into a working ‘steady state’ and

usually discarded in regular data taking –, a ‘drop’ when this is finished and then a more or less smooth line – the second plateau – for the remaining 10.000 simulated milliseconds.

When in this regime the average power draw that NEST-GPU can push a single H100 to reach is around **247W** with an absolute maximum of **252W**. If we take into consideration all GPUs the total draw is **449.5W**. Runtimes are **3.8s** for building and **27.3** for the actual simulation.

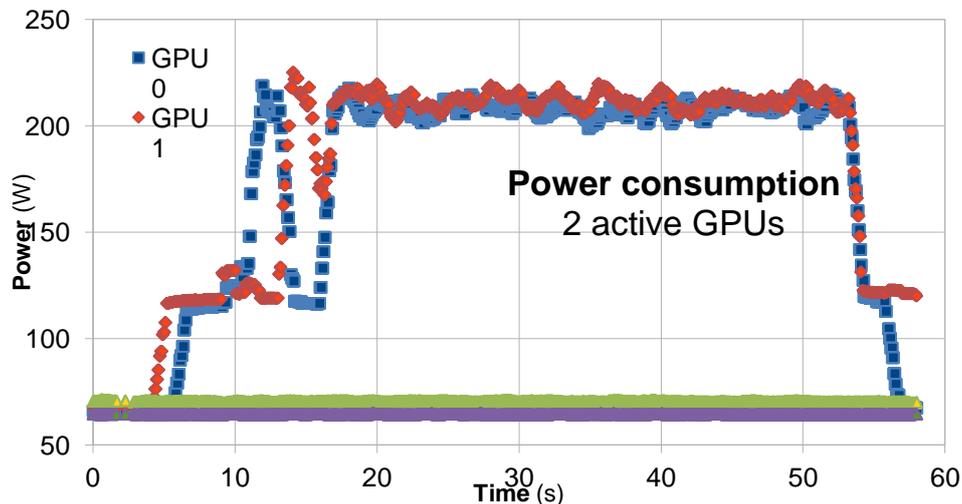


Figure 42: [NEST-GPU] 2 active GPUs power reading on IDV-A

As said, current testing is with a weak scaling benchmark; the same run is performed replicating a network of the same size on a second GPU and randomly connecting a fixed fraction of the neurons between the two. Here the simulation steps are necessarily interleaved with the exchange of messages between the GPU memories, which seems the reason why the throttling of the GPU frequencies by the thermal controller daemon is able to keep the average power draw in the steady state between **208W** and **212W**, as shown in Figure 42, with an absolute maximum of **219W** and a grand total of **554W** including the idle ones. Here the building time is between **3.5s** and **4.4s** and the actual simulation takes **37.4s** on one GPU and **38.2s** on the other. It is apparent that adding inter-GPU communications – which are absent in the single GPU run – amounts to a significant 37% more runtime. It is being investigated how much this can be improved by optimizing the MPI calls – that currently require staging from the GPU memory to the CPU and back – with the GPU-aware ones that include support for direct accessing the GPU memory.

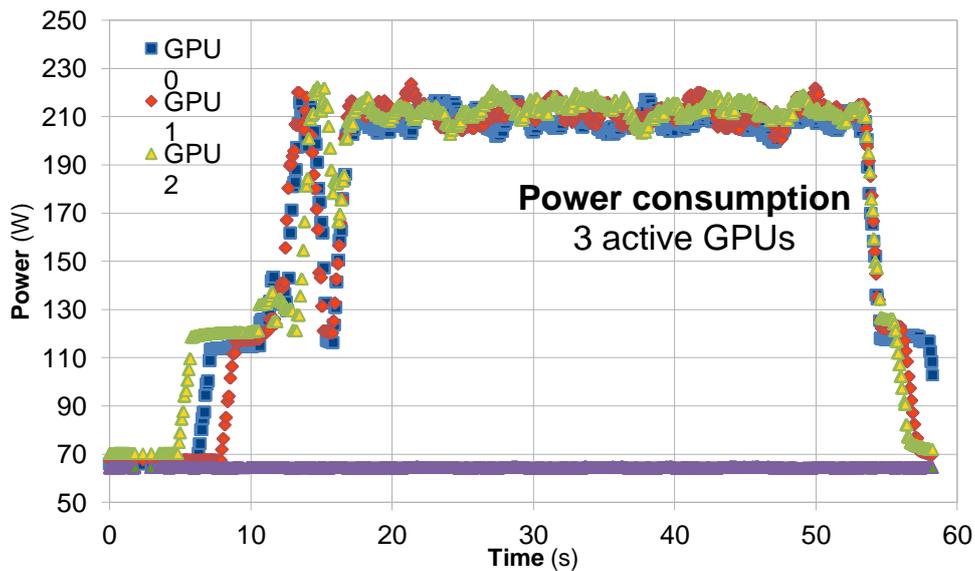


Figure 43: [NEST-GPU] 3 active GPUs power reading on IDV-A

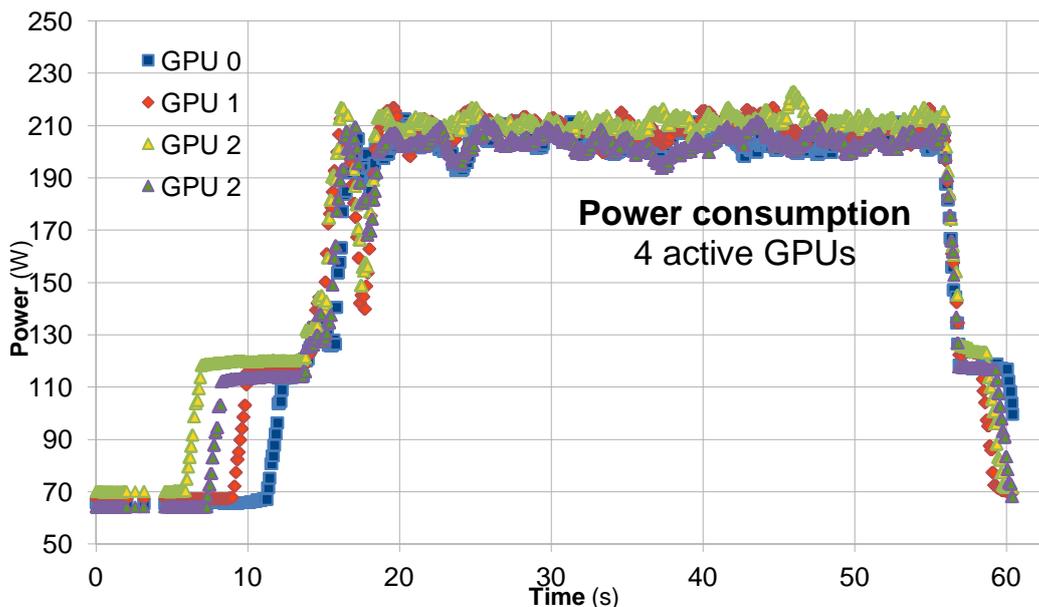


Figure 44: [NEST-GPU] 4 active GPUs power reading on IDV-A

The situation is similar for the 3 GPUs (Figure 43) and the 4 GPUs (Figure 44) cases; having accounted for the advantage of the no-comms situation for the single GPU, we see that the real weak scaling is quite good: building times are still around **3s** or less either for the 3 and 4 GPUs cases while actual simulation runtimes are between **37.7s** and **38s** for the 3 GPUs case and **38.4s** and **38.6s** for the 4 GPUs case, meaning a less than 3% difference in runtime between datasets which are two, three and four times the size of the single GPU case. Of course, the power consumption goes up as well, with 3 GPUs drawing between **208W** and **213W** on average in the steady state with a grand total of **697W** including the idle one and 4 GPUs drawing between **204W** and **211W** on average in the steady state and a grand total of **829W**.

For reference, we found that the very same problem run on the 4 NVIDIA A100 GPUs of the TextaRossa partition of the Dibona cluster runs consistently slower in the 2, 3 and 4 GPU cases – between **39s** and **40s** in the first case, between **40.5s** and **41.5s** in the second and between **42s** and **43s** in the third – with a power draw that does not exceed **156W** per active GPU in the steady state – while the single GPU case is actually a little faster, with **24.7s** and a **184W** average power draw in the steady state. The results for the Dibona A100 GPUs are plotted in Figure 45, Figure 46, Figure 47 and Figure 48 and the runtimes – estimated with the longest among all GPUs for the multi-GPU runs – for the simulation part of the run and both systems are in Table 18.

# of GPUs	Runtimes on A100 (s)	Runtimes on H100 (s)	Wattage on A100 (W)	Wattage on H100 (W)
1	24.72	27.31	184.36	247.28
2	40.35	38.25	308.86	419.41
3	41.52	38.04	460.15	632.64
4	43.13	38.66	611.65	828.65

Table 18: NEST-GPU comparison between A100 and H100 GPUs

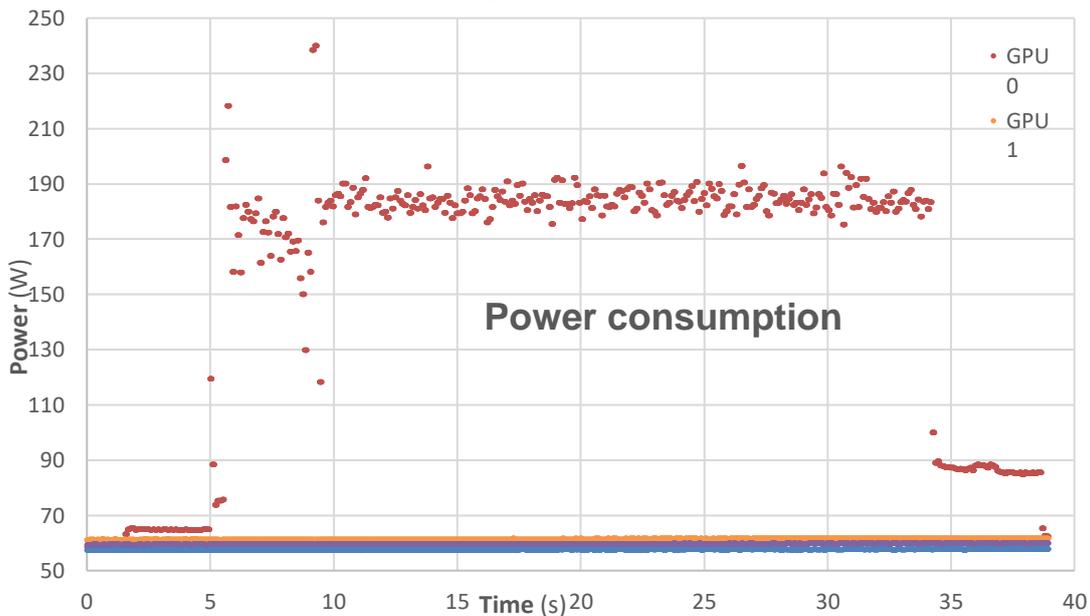


Figure 45: [NEST-GPU] Single active GPU power reading on Dibona

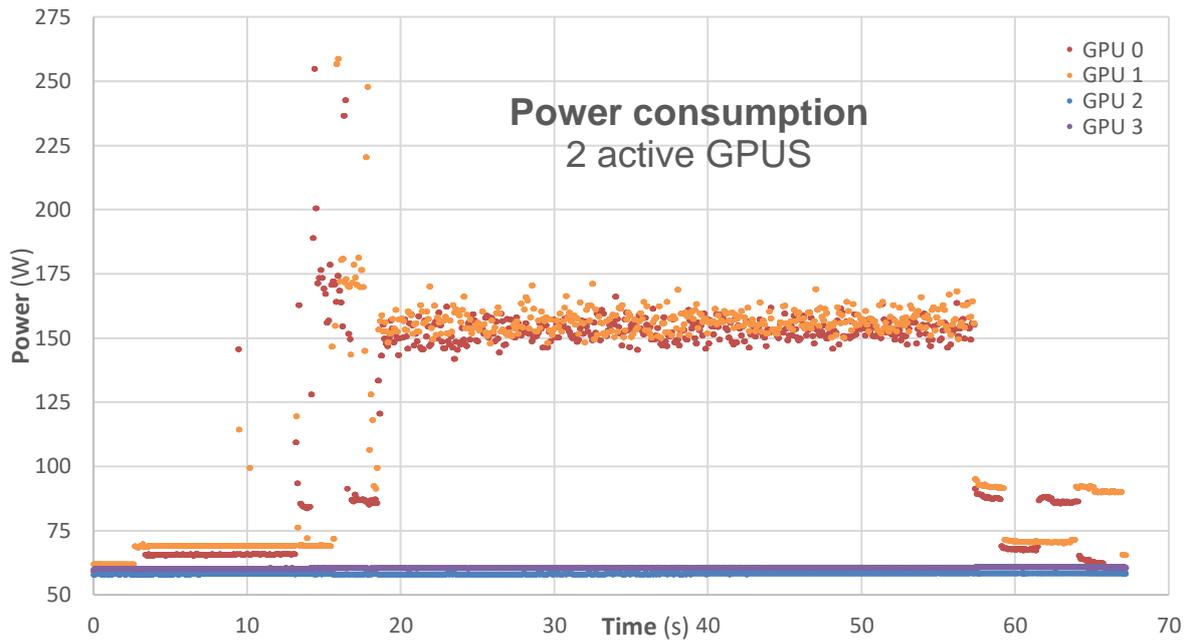


Figure 46: [NEST-GPU] 2 active GPUs power reading on Dibona

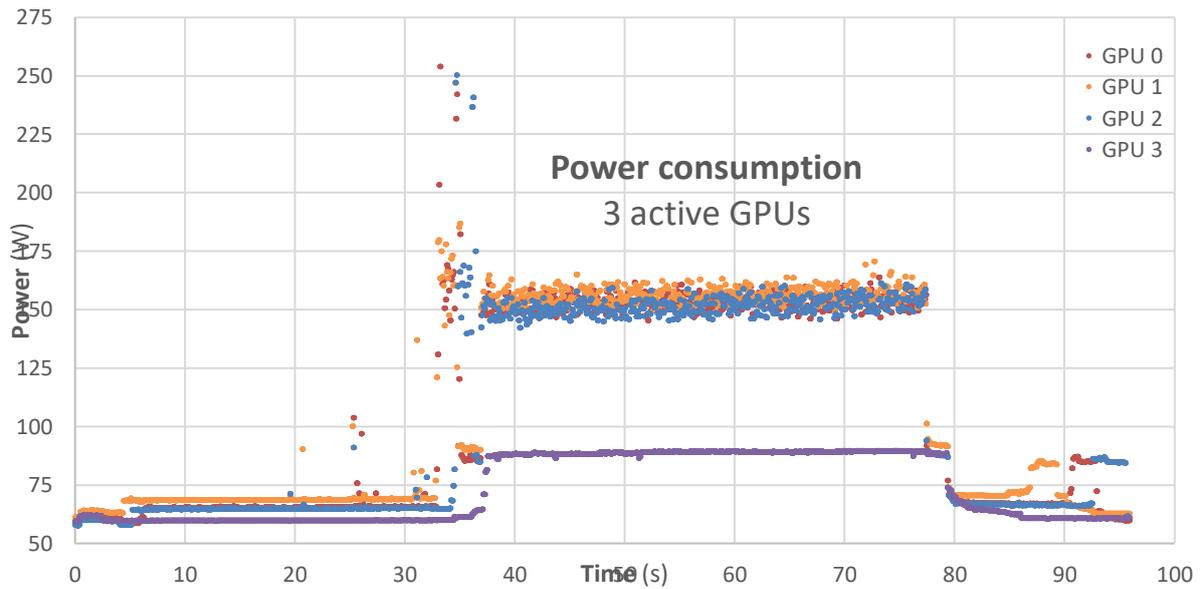


Figure 47: [NEST-GPU] 3 active GPUs power reading on Dibona

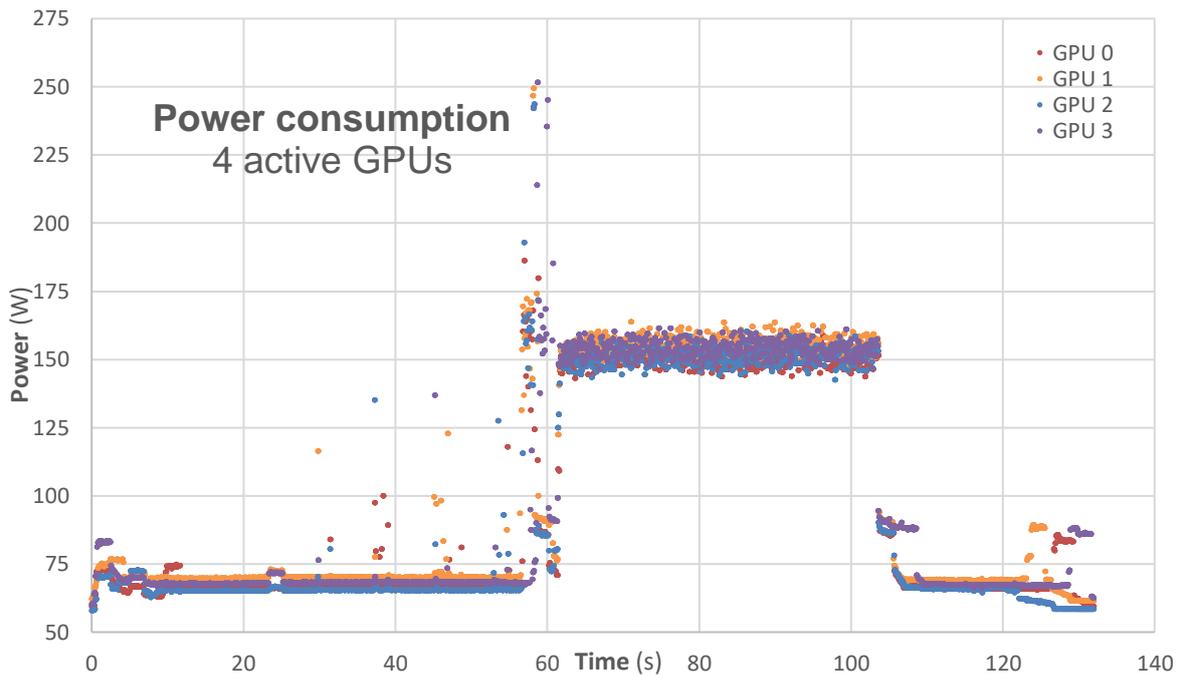


Figure 48: [NEST-GPU] 4 active GPUs power reading on Dibona

The very same benchmark was performed with the CPU-only version of the NEST code running on the ARM cores on the IDV-E platform. Following the investigation done in D6.2 we directly chose the optimal process/cores layout for the platform, 16 MPI processes using 8 OpenMP threads each for a grand total of 128 busy cores and fixed the size to be the same of the single GPU run on the IDV-A. In the plot in Figure 49 is the polling of the internal probe for power consumption of both CPU sockets as returned by the *sensors* command with a 0.7s time interval between readings.

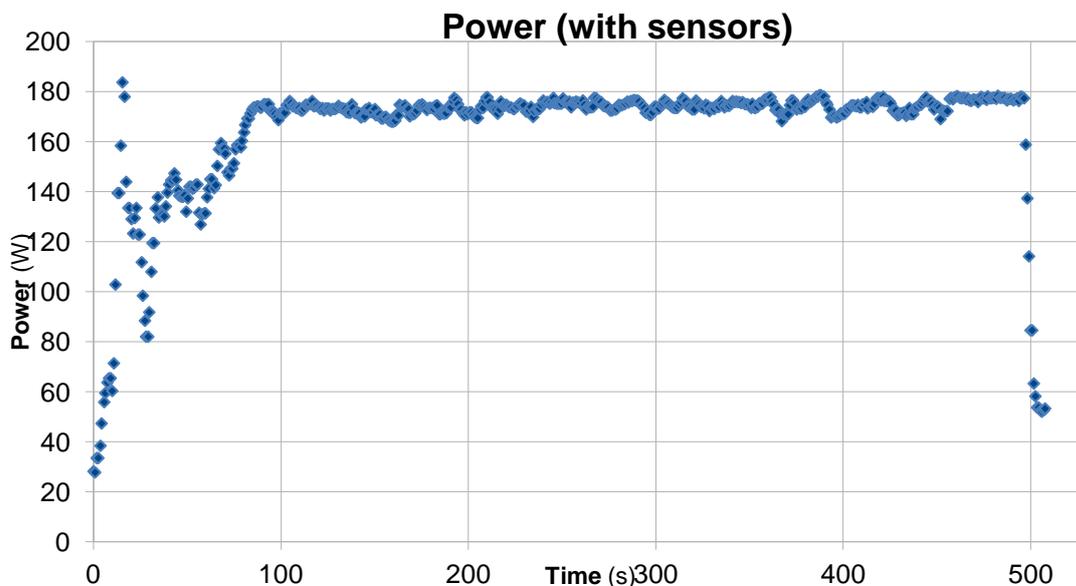


Figure 49: [NEST-GPU] IDV-E power reading

Here the building and pre-simulation phase are less distinguishable, appearing as an irregular slope up to the steady state plateau; maximum power draw is **184W** while average power in the steady state is **174W**, with **466s** for the 10.000 simulated milliseconds of network activity. For completeness, we report in Figure 50 the output of the Power Distribution Unit connected to the IDV-E node as queried by the system scripts available on the system. The slope is consistent with the previous plot, reaching a steady state after a little less that 100s; the total power output is reported as **416W** in the steady state (with a **467W** spike) which with an estimated 206W in idle (the lowest point in the plot) is consistent with more or less 200W more when running at full load.

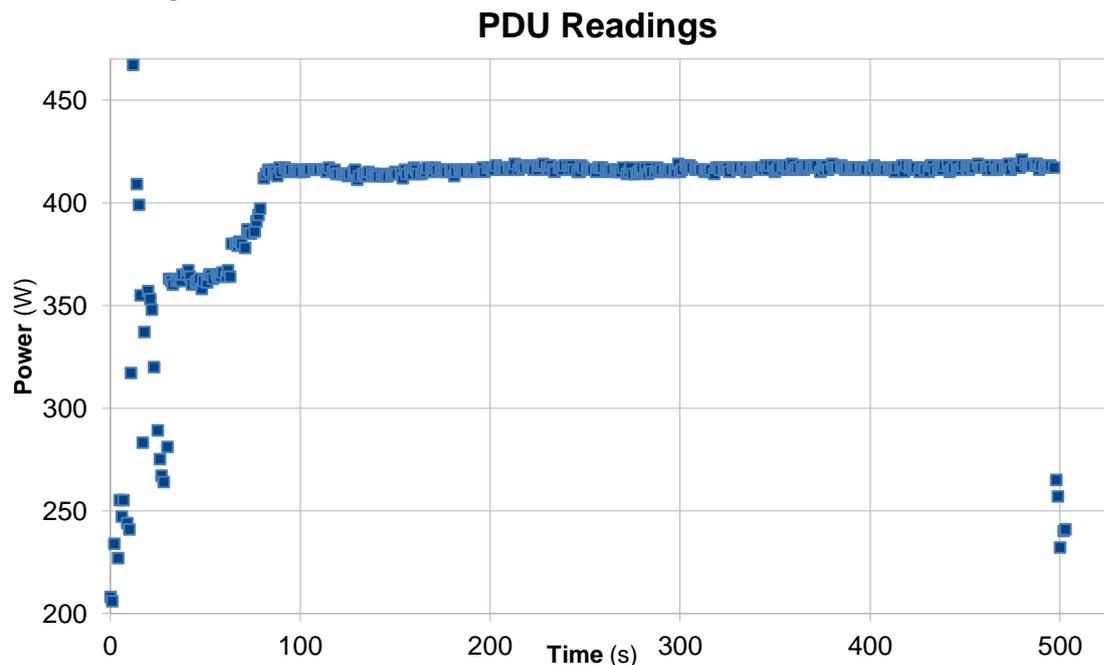


Figure 50: [NEST-GPU] IDV-E PDU reading

As a final comparison, we tried running the CPU-only variant of the neural network application on the CPUs of the IDV-A. These CPUs belong to the Sapphire Rapids family by Intel with 52 cores (doubled when hyperthreading is active) and on the IDV-A are in a dual socket configuration for a grand total of 208 cores. We tried simulating the same network sizes also achievable by 1, 2, 3 and all 4 GPUs together while using the *turbostat* utility to perform the power readings with a 0.5s interval between them; the results are plotted in Figure 51, Figure 52, Figure 53 and Figure 54 and in Table 19 and correspond to the optimal process/cores layout for running the neural simulation on the IDV-A: 8 OpenMP threads per MPI process, 24 total MPI processes yielding a total of 192 used cores.

It can be seen that the ‘idle’ plateau for the power draw (it can be spotted more clearly as a ‘drop’ at the end of the runtime in the plot) is around **400W**; we remark that this is probably a significant overestimation since the power regulation and cooling system on the IDV-A requires a software daemon with a very aggressive scheduling which keeps some cores busy running throughout the run.

It can also be seen that the runs become slower and slower when increasing the size of the network, which is expected, but disproportionally with regard to the increased factor: for example, the network size fitting a single GPU memory can be simulated in around **170s** while twice as large a network takes more than twice the time at **424s**. This is due to the fact that the communications between larger numbers of spiking neurons grow and become prevalent

against the compute; this interference by necessary communication also appears not to let the processors reach higher power states and therefore yields lower maximum and average wattages, as clearly exhibited in the plots.

Sim. size equivalent to that on # of GPUs	Runtimes (CPU-only) on IDV-A (s)	Wattage (CPU-only) on IDV-A (W)
1	170.27	678.87
2	424.37	642.96
3	775.27	627.74
4	1182.85	623.17

Table 19: [NEST-GPU] neural simulation on IDV-A CPU-only – runtimes and power

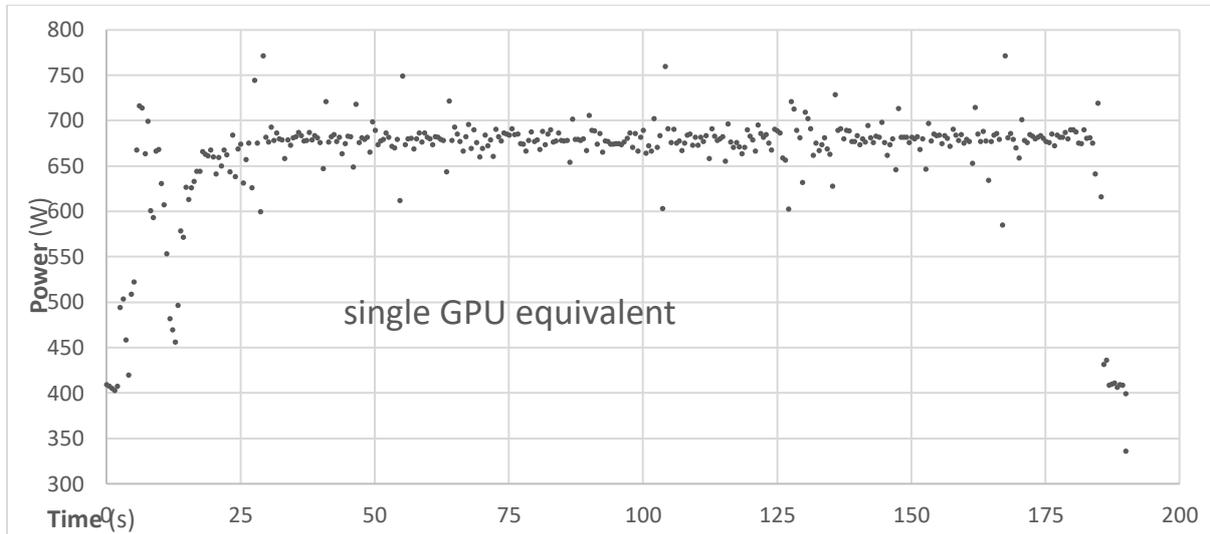


Figure 51: [NEST-GPU] CPU-only run power reading on IDV-A – 1 GPU size

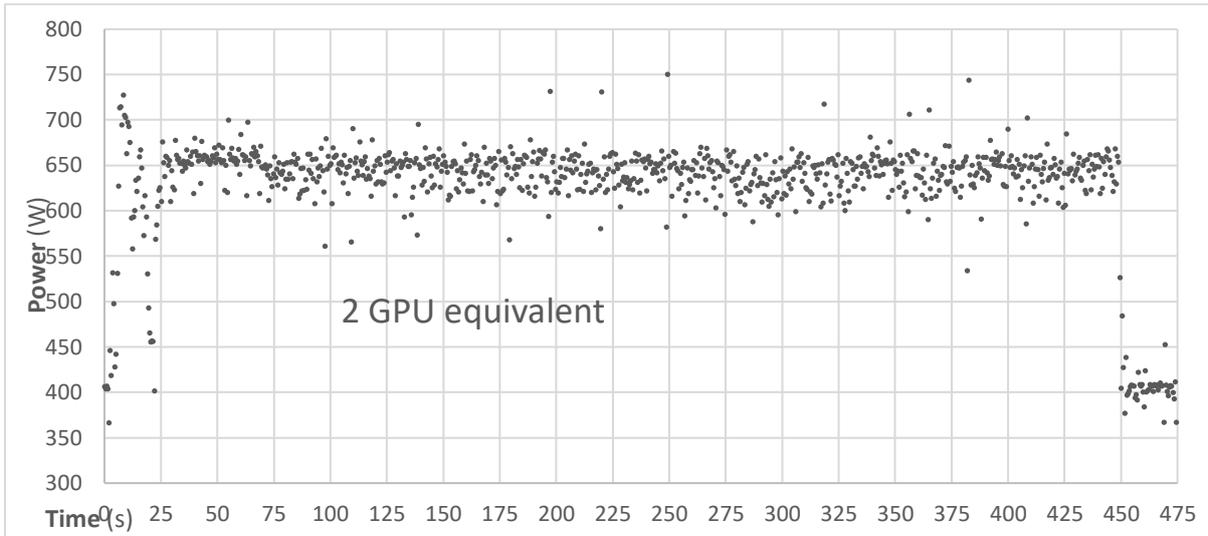


Figure 52: [NEST-GPU] CPU-only run power reading on IDV-A – 2 GPUs size

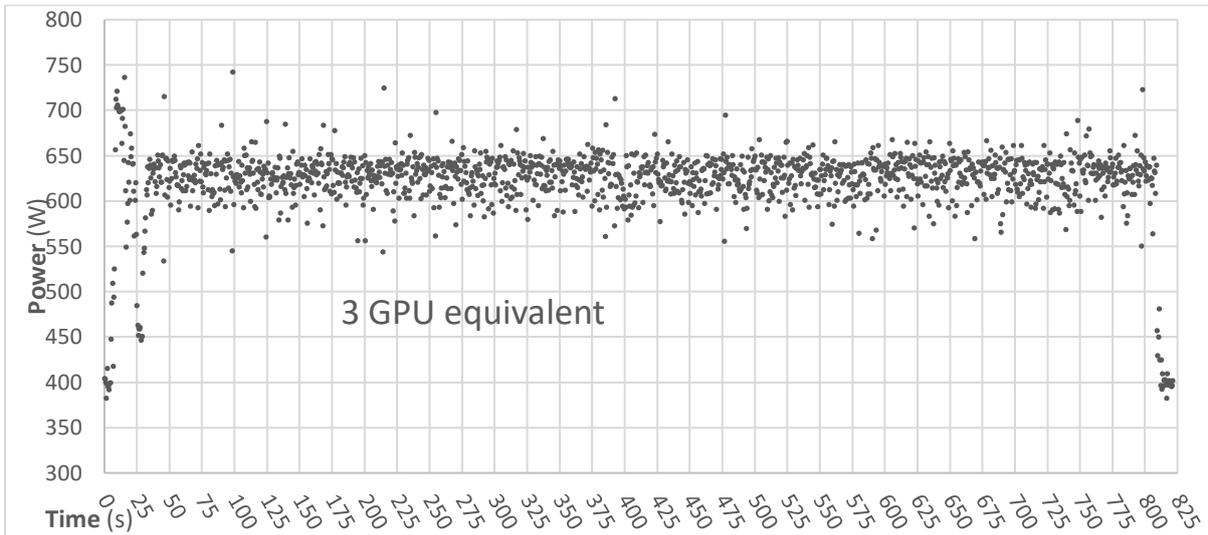


Figure 53: [NEST-GPU] CPU-only run power reading on IDV-A – 3 GPUs size

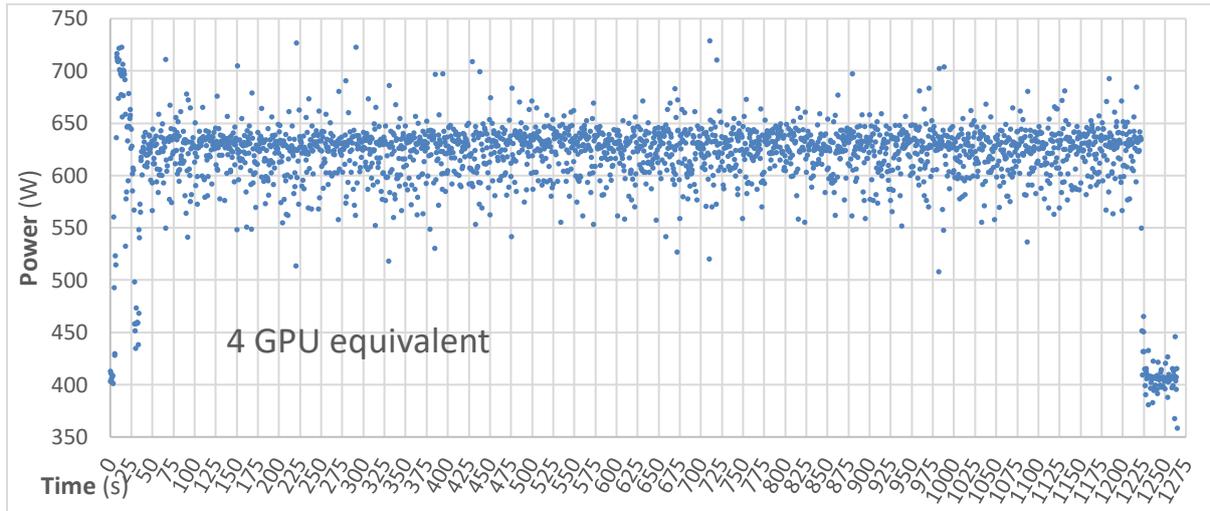


Figure 54: [NEST-GPU] CPU-only run power reading on IDV-A – 4 GPUs size

To conclude this section, we report in Table 20 a summary of the measured KPIs for the neural network simulation, on the different configurations examined for the three computing platforms used (IDV-A, IDV-E, Dibona).

KPI	IDV-A single GPU	Dibona single GPU	IDV-A multi-GPU	Dibona multi-GPU	IDV-A CPU-only	IDV-E CPU-only
Simulated milliseconds/s	366.3	404.6	259.1÷267.4	231.8÷ 247.8	58.7	21.5
Energy to solution (kJ)	6.7	4.6	7.6÷8.0 (per GPU)	6.2÷6.6 (per GPU)	115.6	81.1
Synaptic updates/J	1814.7	2643.5	1600.0÷1520.0	1842.4÷1961.3	105.2	149.9

Table 20: [NEST-GPU] KPIs for the neural network simulation application

3.6 RAIDER – INFN

In this section we report the results obtained from the execution of the RAIDER application implemented with APEIRON framework on the twin-FPGA IDV-E node first, and then on a couple of IDV-E nodes (four interconnected Xilinx U280 FPGA boards in total).

RAIDER is a high throughput online streaming processing application implemented on FPGA and its task is to perform particle identification (PID) on the stream of events generated by the RICH (Ring Imaging CHerenkov) detector in the CERN NA62 experiment, using neural networks.

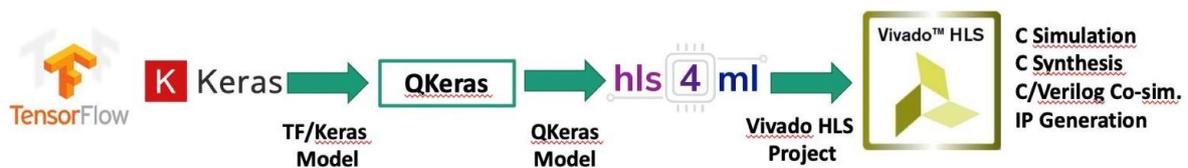


Figure 55: [RAIDER] The workflow for the generation of CNN kernels in RAIDER

In this case we adopt a Convolutional Neural Network (CNN) developed and trained offline by using Tensorflow/Keras, the steps of the workflow for the CNN generation are shown in Figure 55.

The model receives as input a 16x16 image of the hit photomultipliers (PMTs) map - depicted in Figure 56 - for each physics event and its goal is to produce an estimate for the number of charged particles (0, 1, 2, >=3) for any RICH detector event. This corresponds to the number of ring tracks that can be reconstructed from the pattern of PMTs that have been illuminated by the Cherenkov light cone emitted by a charged particle traversing the detector.

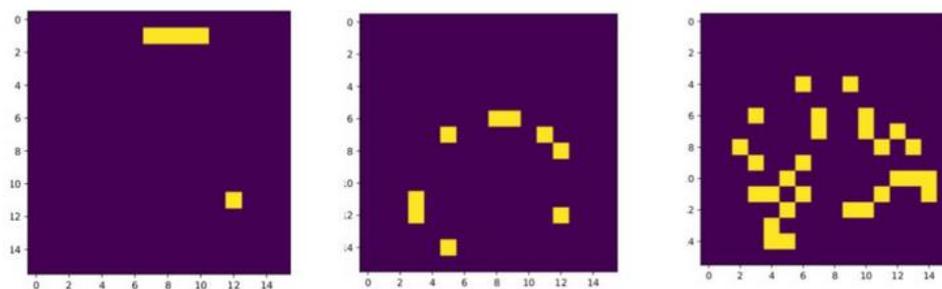


Figure 56: [RAIDER] Example of input images for the CNN

(left class 0, center class 1, right class 2).

To produce the training and validation data for the CNN, we prepared different data sets composed by events extracted directly from NA62 database using the experiment analysis framework. The ground truth, used for training, was provided by the seedless offline reconstruction method.

To limit the FPGA resources footprint, as second stage, we performed a quantization step on the model using QKeras, resulting in two different fixed-point representations: $\langle 8, 1 \rangle$ for weights and biases and $\langle 16, 5 \rangle$ for activations.

As last step, the quantized model is translated into the corresponding Vitis HLS implementation using HLS4ML tool; implementation validated with Vivado C/Verilog co-simulation in terms of resources usage, performances (throughput and latency) and efficiency (referring to the classification accuracy of the CNN model). Lastly, the model can be synthesized as kernel IP to be integrated in the APEIRON framework and deployed on the FPGA.

Since the instantiation interval of the CNN obtained from HLS4ML scales with the size of the image, we expect to have a throughput for a single kernel implementation much smaller than

the 10 MHz required by the NA62 L0 trigger. So, in order to improve the performances, RAIDER application has been designed with multiple processing kernels displaced in different nodes, capable to receive data/events from the network thanks the HAPECOM communication APIs.

We used two different setups to test the performance of this RAIDER implementation depicted in Figure 57 and Figure 66:

- Four interconnected Xilinx® Alveo U200 installed in the INFN Roma APE Lab in a ring topology, one FPGA for each of the four single Intel(R) Xeon(R) Silver 4410T CPU nodes, used for the co-design and development of the APEIRON framework and of the RAIDER application.
- Four interconnected Xilinx® Alveo U280 installed on the IDV-E node in a ring topology, 2 FPGAs for each of the two IDV-E nodes, to validate the full functionality of the APEIRON framework on the node and assess its performance using the RAIDER KPIs defined in D6.1: throughput, expressed as the number of reconstructed events per second, and energy efficiency, expressed as the number of reconstructed events per Joule consumed by the processing node.

In particular, we were interested in studying the throughput of the overall system, since it is one of the most important requirements for the integration in trigger and data acquisition systems in HEP experiments. The directly interconnected FPGA boards are used as nodes of a RAIDER deployment with the APEIRON framework, with two distinct roles:

- **Preprocessing node:** data (events containing a list of illuminated photomultipliers-PMT hitlist) are loaded from Host memory and forwarded to the network via an HLS kernel (*krnl_sender*). Data are received and processed by a variable number of Imagifier HLS kernels that transform the PMT hitlist information into a 256bit word (16x16 B&W image) that is sent to the Computing node through the internode ports of the INFN Communication IP. Finally, another task (*krnl_receiver*) is in charge of receiving the output of the CNNs processing and of storing it on the Host memory. The overall processing time, from the first data packet sent to the last result packet received, is measured on this node.
- **Computing node:** images coming from external links are taken as input and dispatched to the CNN HLS kernels (each of them connected to a different INFN Communication IP intranode port) to compute the inference task. Results are then sent back to the preprocessing node.

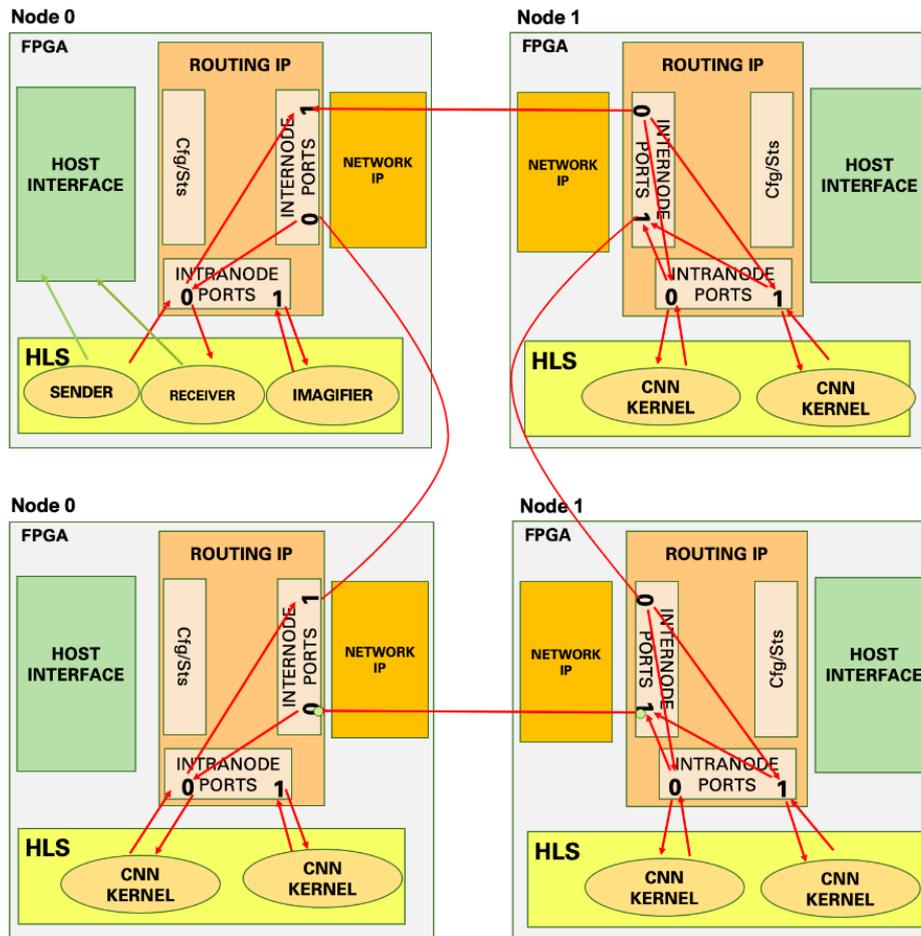


Figure 57: [RAIDER] Test setup on the Xilinx® Alveo U200 installed in the INFN Roma1 APE Lab

All tests performed in both testbeds have been done by using a 200 MHz global clock in the hardware setup. This clock increasing with respect to the one used for the D6.2 measures has been possible thanks to the further improvements in the INFN Communication IP (described in detail in the D2.9). Another difference with respect to the D6.2 application is that the data - to be sent through the network from the preprocessing node – are now loaded from the BRAM instead of the DDR FPGA memory. This change is due to a limit that we’ve seen while testing the application loading events from DDR memory: while increasing the number of CNN kernels in the setup to test RAIDER scaling, we noticed that, working with 3 CNNs or more, the maximum attainable throughput was of 1.278MHz. This corresponds to the instantiation interval of the *sender* HLS kernel of ~160 clock cycles, as can be seen from the report obtained in output of the compiling Vitis process depicted in Figure 58.

Modules & Loops	Issue Type	Slack	Latency (cycles)	Latency (ns)	Iteration Latency	Interval
+ krnl_sender	-	0.00	432000075	2.160e+09	-	432000076
+ krnl_sender_Pipeline_main_loop	-	0.00	432000001	2.160e+09	-	432000001
o main_loop	-	3.65	432000000	2.160e+09	160	-

Figure 58: [RAIDER] Vitis synthesis report of krnl_sender HLS kernel

For this reason, we decide to work with data loaded from the BRAM by the *krnl_sender* for the tests reported in this deliverable, emulating more closely the input data throughput at the experiment where data from the detector arrives through a network channel.

3.6.1 RESULTS ON APE LAB TESTBED

In this setup, we have scaled the system starting from tests on 2 FPGA nodes (one preprocessing and one computing) up to 4 (adding 2 more computing nodes), setting from the host (at computing kernels launch) the number of running CNN on each computing node. For what concerns energy efficiency, CPU hosts power measurements have been performed using *turbostat*: a Linux command-line utility that reports processor topology, frequency, idle power-state statistics, temperature, and power on X86 processors. The measurement was performed concurrently with application execution on each host in the setup, results are shown in Figure 59.

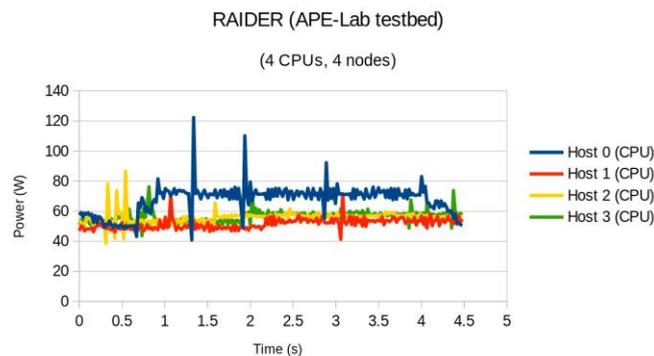


Figure 59: [RAIDER] Application CPU hosts power profiles (4 Intel Sapphire Rapid)

Power measurements for FPGA implementing HLS kernels on both type of nodes, were extrapolated for the XRT summary .csv file obtained in output from the CPU host application.

Following the node and number of CNNs scaling, power profiles have been produced for each of the RAIDER configuration tested and are reported in Figure 60, Figure 61 and Figure 62.

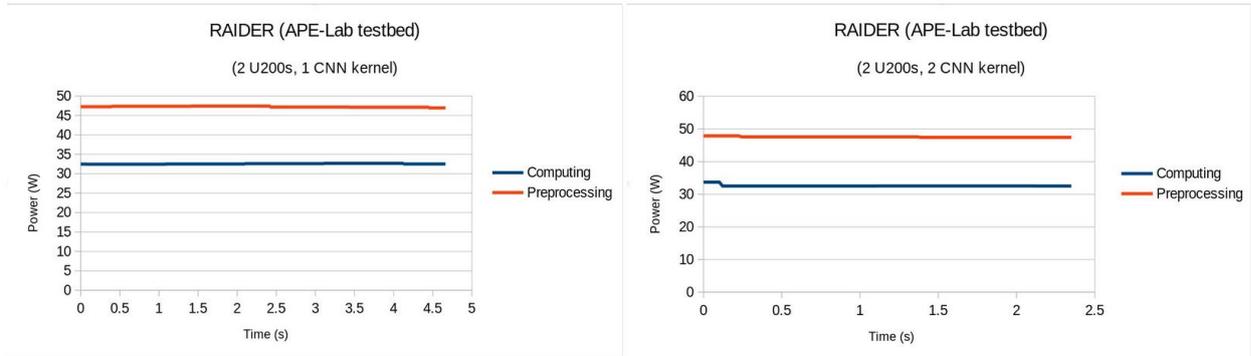


Figure 60: [RAIDER] FPGA power profiles (2 Alveo Xilinx U200 setup)

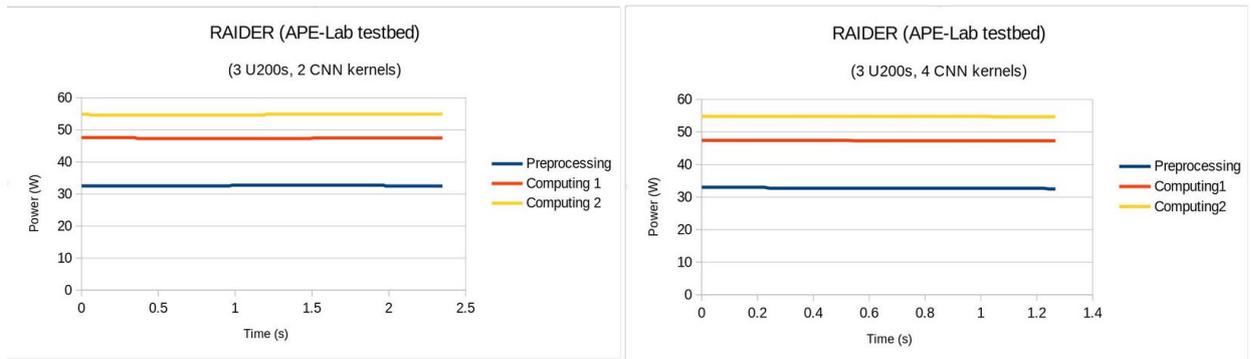


Figure 61: [RAIDER] FPGA power profiles (3 Alveo Xilinx U200 setup)

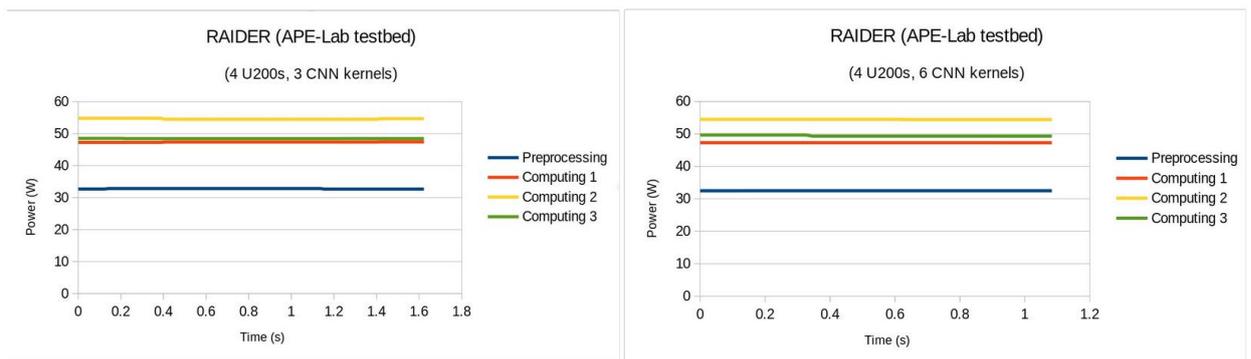


Figure 62: [RAIDER] FPGA power profiles (4 Alveo Xilinx U200 setup)

The integrated processing throughput of the RAIDER application deployed on the APE Lab testbed has been measured host-side by scaling the number of implemented computing nodes. Results are tabulated in Table 21 and plotted in Figure 63, together with the energy efficiency values obtained by integrating the power profiles measured for all the tested setups.

# computing nodes	# CNNs	Throughput (Mevents/s)	Energy efficiency (kevents/J)
1 node	1 CNN	0.581	3.015
	2 CNNs	1.163	6.005
2 nodes	2CNNs	1.163	3.838
	4CNNs	2.325	7.692
3 nodes	3CNNs	1.744	4.298
	6CNNs	2.692	6.626

Table 21: [RAIDER] Processing Throughput with an increasing number of Computing nodes (and CNN HLS kernels)

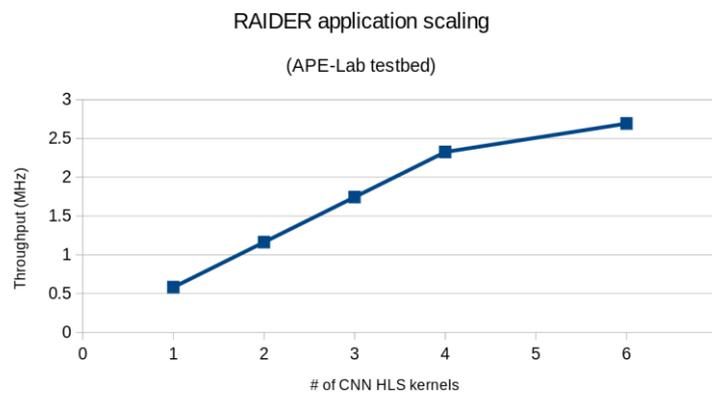


Figure 63: [RAIDER] Throughput scaling with the number of deployed CNN HLS kernels

The presented results show the good scaling of system performance with the number of computing nodes, while the flattening slope of the curve when the number of CNNs goes beyond 4 is mainly due to the saturation of the data injection rate in the *krnl_sender* and the instantiation interval of the *imagifier* (~70 clock cycles) in the preprocessing node.

3.6.2 U280 RESULTS ON IDV-E

In this setup, we have tested the system scaling from 2 FPGAs (one preprocessing and one computing) up to 4 (adding 2 more computing FPGAs) with a slightly different setup with respect to the one used in the U200 APE Lab testbed. Since the Xilinx Alveo U280 has more available resources, e.g. DSPs, with respect to the U200, we could implement 3 CNNs HLS kernels on a single computing U280 FPGA (a Vitis report of the hardware occupation is shown in Figure 64) since a single CNN occupies ~25% of the maximum DSP resources.

Name	LUT	LUTsMem	REG	BRAM	URAM	DSP
Platform	13.40%	3.07%	7.76%	11.61%	0.00%	0.04%
✓ User Budget	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%
Used Resources	11.51%	0.68%	1.93%	0.67%	1.25%	77.26%
Unused Resources	88.49%	99.32%	98.07%	99.33%	98.75%	22.74%
aggregator_0 (1)	0.02%	0.00%	0.04%	0.00%	0.00%	0.00%
aggregator_0_1	0.02%	0.00%	0.04%	0.00%	0.00%	0.00%
aggregator_1 (1)	0.02%	0.00%	0.04%	0.00%	0.00%	0.00%
aggregator_1_1	0.02%	0.00%	0.04%	0.00%	0.00%	0.00%
aggregator_2 (1)	0.02%	0.00%	0.04%	0.00%	0.00%	0.00%
aggregator_2_1	0.02%	0.00%	0.04%	0.00%	0.00%	0.00%
dispatcher_0 (1)	0.03%	0.00%	0.04%	0.00%	0.21%	0.00%
dispatcher_0_1	0.03%	0.00%	0.04%	0.00%	0.21%	0.00%
dispatcher_1 (1)	0.03%	0.00%	0.04%	0.00%	0.21%	0.00%
dispatcher_1_1	0.03%	0.00%	0.04%	0.00%	0.21%	0.00%
dispatcher_2 (1)	0.03%	0.00%	0.04%	0.00%	0.21%	0.00%
dispatcher_2_1	0.03%	0.00%	0.04%	0.00%	0.21%	0.00%
top_nnet (3)	11.34%	0.68%	1.69%	0.67%	0.62%	77.26%
top_nnet_1	3.78%	0.23%	0.56%	0.22%	0.21%	25.75%
top_nnet_2	3.78%	0.23%	0.56%	0.22%	0.21%	25.75%
top_nnet_3	3.78%	0.23%	0.56%	0.22%	0.21%	25.75%

Figure 64: [RAIDER] Vitis synthesis report of hardware project to be deployed on Computing node. The DSP utilization is in the red boxes.

In addition to this, taking into account the limitations in terms of throughput found in the APE Lab testbed, we decided to deploy up to 3 *imagifier* HLS kernel on the single preprocessing FPGA. The resulting hardware setup is so reported in Figure 65 and Figure 66, in which are depicted, respectively, the design used for single and double node RAIDER execution on IDV-E testbed.

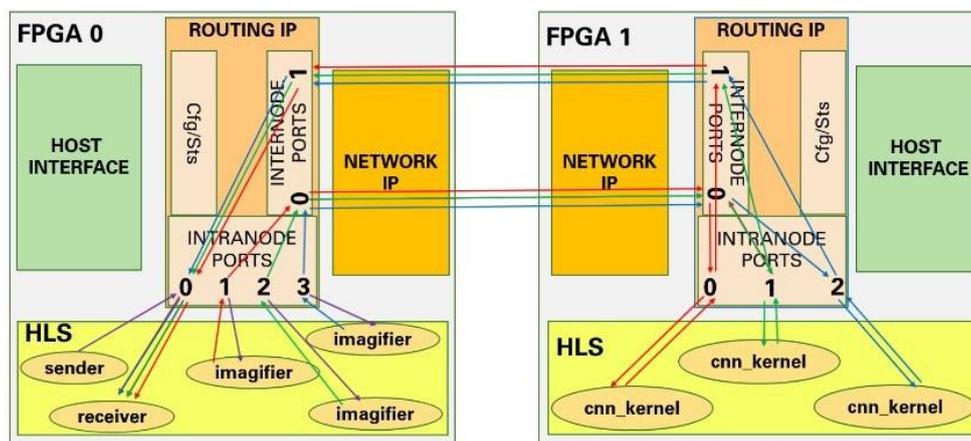


Figure 65: [RAIDER] Test setup on the Xilinx® Alveo U280 installed in the IDV-E testbed (single node)

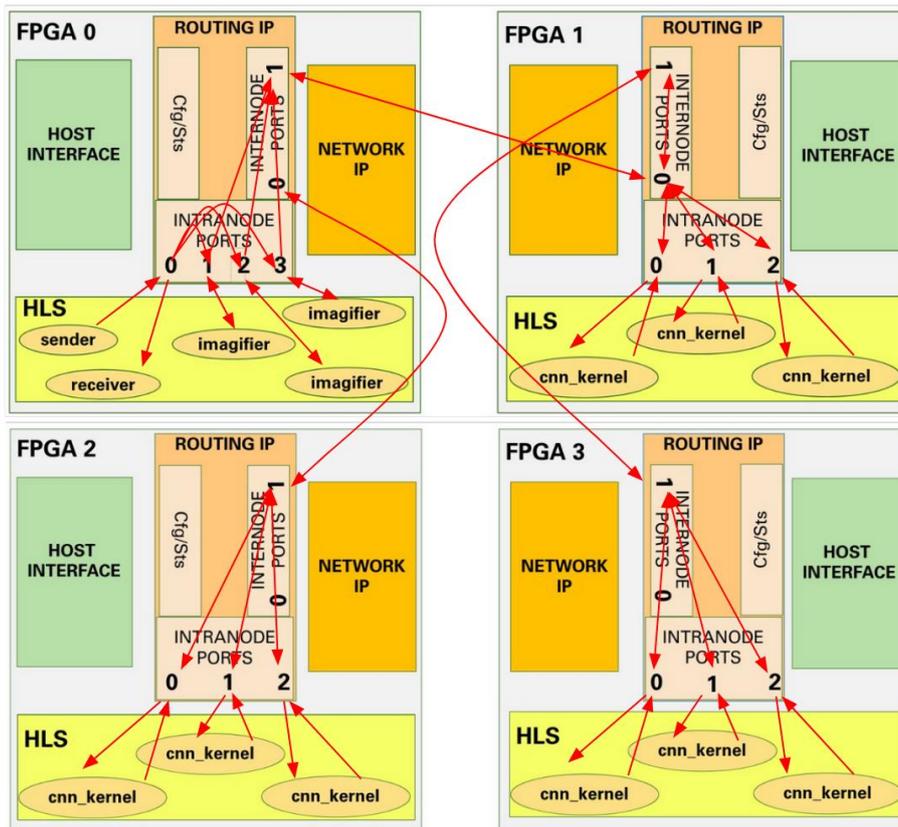


Figure 66: [RAIDER] Test setup on the Xilinx® Alveo U280 installed in the IDV-E testbed (double node)

The integrated processing throughput and energy efficiency of the system has been measured host-side and results are tabulated in Table 22.

# computing FPGAs	# CNNs	Throughput (Mevents/s)	Energy efficiency (kevents/J)
1 FPGA	3 CNNs	1.813	12.490
2 FPGAs	6 CNNs	3.409	13.685
3 FPGAs	9 CNNs	4.874	16.336

Table 22: [RAIDER] processing time per event with an increasing number of Computing FPGAs (and CNN HLS kernels)

For what concerns energy efficiency, CPU hosts power measurements have been performed using *sensors*: a free and open-source application that provides tools and drivers for monitoring temperatures, voltages, and fans speeds of servers. The measurement has been performed concurrently with the application execution on each node present in setup and are reported in Figure 67. From CPU power profiles, it is possible to notice a ~4W difference between the two nodes used in the IDV-E testbed.

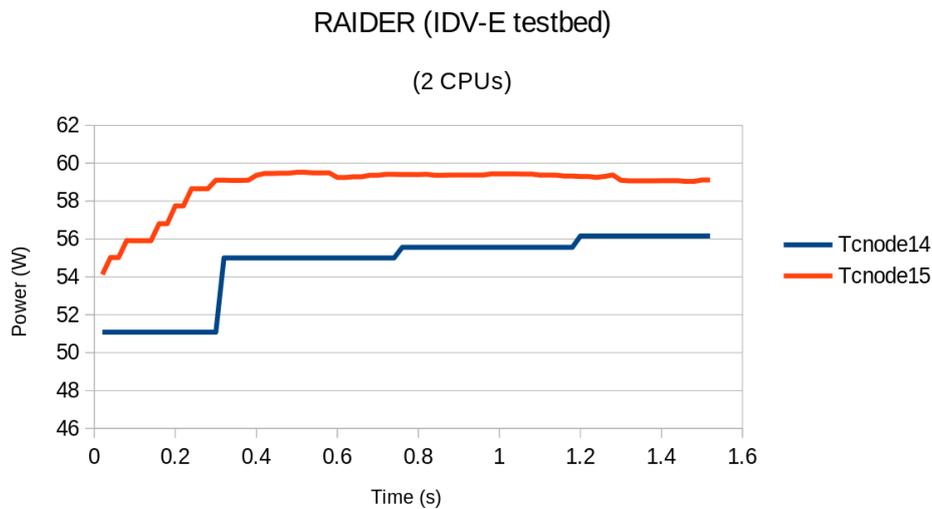


Figure 67: [RAIDER] CPU hosts power profiles (2 Ampere Altra Max processor)

For tracking accelerator FPGA cards power profiles, as for the APE Lab testbed case data were extrapolated for the XRT summary .csv file obtained in output from the CPU host application. As can be seen from Figure 68, Figure 69 and Figure 70, Computing FPGAs 2 and 3 seems to have a power average slightly larger than the Computing FPGA 1: this is probably due to the fact that the Computing FPGA 2 and 3 are placed in a different IDV-E server with respect to the one in which the Preprocessing FPGA and Computing FPGA 1 are. In particular computing FPGA 2 and computing FPGA 3 are hosted in a air-cooled node (tcnode14), while computing FPGA 1 is hosted in node tcnode15 equipped with the two-phase cooling system developed in TEXTAROSSA.

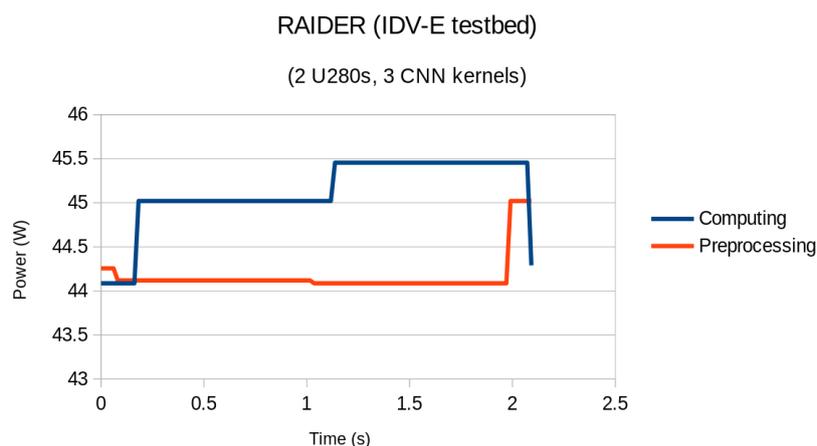


Figure 68: [RAIDER] FPGA power profiles (2 Alveo Xilinx U280 setup)

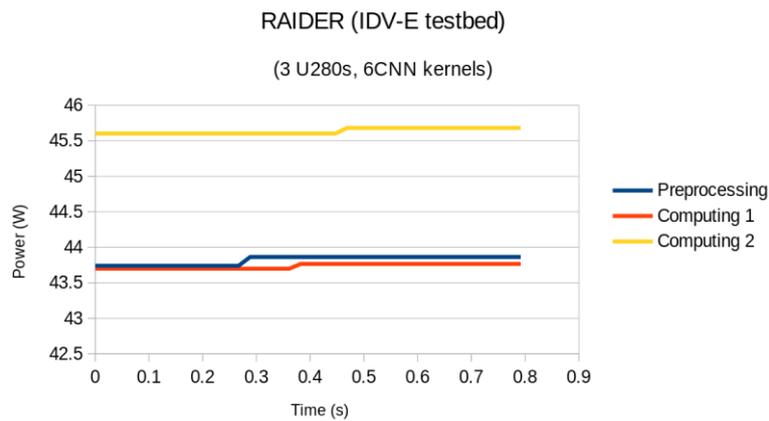


Figure 69: [RAIDER] FPGA power profiles (3 Alveo Xilinx U280 setup)

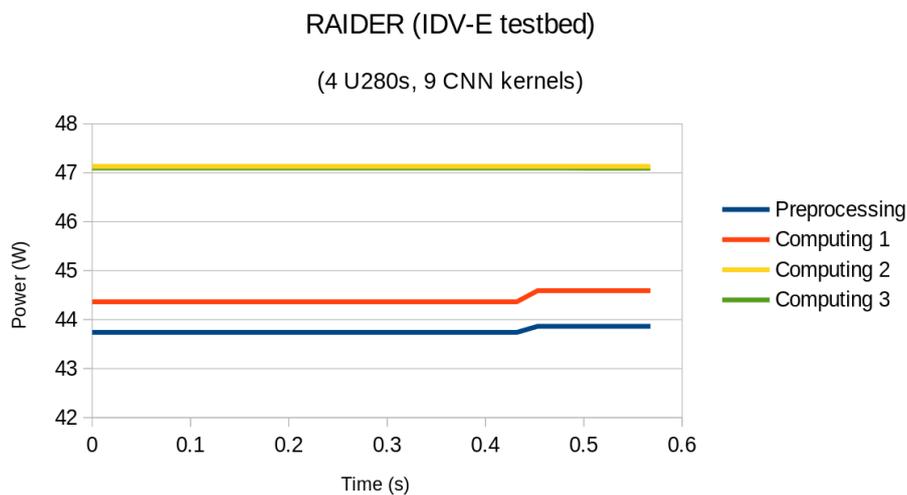


Figure 70: [RAIDER] FPGA power profiles (4 Alveo Xilinx U280 setup)

3.6.2.1 KPIs

To make a comparison with the KPIs presented in D6.2, in Table 23 we report as baseline the performance values obtained for the RAIDER setup in that deliverable. For what concern the energy to solution and energy efficiency measurements, we considered correct to insert the relative FPGA energy consumption together with the node total energy data: this allows us to disentangle the different contributions to the total energy consumption and highlight the performance of the FPGA accelerator board alone.

Our goal is to show the updates of the application with respect to the past one and the scaling results obtained via the capabilities of the APEIRON framework and the adoption of the IDV-E platform.

KPI	RAIDER @100 MHZ [1 FPGA, 1CNN]	RAIDER @100 MHZ [1 FPGA, 2CNNs]
(per class)	efficiency: - 0: 92% - 1: 79% - 2: 75% - 3+: 76% purity: - 0: 83% - 1: 88% - 2: 70% - 3+: 80%	efficiency: - 0: 92% - 1: 79% - 2: 75% - 3+: 76% purity: - 0: 83% - 1: 88% - 2: 70% - 3+: 80%
time to solution [s]	9.701	4.898
throughput [events/s]	278324.152	551245.410
energy to solution [J]	563.174 (262.090 FPGA)	267.831 (137.902 FPGA)
energy efficiency [events/J]	4794.255 (10301.805 FPGA)	10079.328 (19579.121 FPGA)

Table 23: [RAIDER] Baseline KPIs evaluated on the execution of the RAIDER application on a single Xilinx Alveo U200 FPGA, processing 2.7M events with a global clock of 100 MHz

Since the updates of the INFN Communication IP discussed in D2.9, it is now possible within the RAIDER application to synthesize a hardware bitstream capable of working with a global clock of 200MHz. This clock increase has been applied in the test for either testbed described in the previous section.

For the APE-Lab testbed setup, the integrals of both CPUs and FPGAs power profiles -depicted in the previous graphs – have been used to compute the energy-to-solution and energy efficiency values which are reported in Table 24, Table 25 and Table 26.

KPI	RAIDER @200 MHZ [2 FPGA, 1CNN]	RAIDER @200 MHZ [2 FPGA, 2CNNs]
time to solution [s]	4.644	2.332
throughput [events/s]	581385.208	1162762.305
energy to solution [J]	895.422 (370.672 FPGA)	449.622 (186.117 FPGA)
energy efficiency [events/J]	3015.271 (7284.055 FPGA)	6005.044 (14506.937 FPGA)

Table 24: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a two Xilinx Alveo U200 FPGA system, processing 2.7M events with a global clock of 200 MHz

KPI	RAIDER @200 MHZ [3 FPGA, 2CNN]	RAIDER @200 MHZ [3 FPGA, 4CNNs]
time to solution [s]	2.332	1.161
throughput [events/s]	1162762.305	1744120.135
energy to solution [J]	703.460 (312.914 FPGA)	351.029 (156.593 FPGA)
energy efficiency [events/J]	3838.171 (8628.577 FPGA)	7691.672 (17242.109 FPGA)

Table 25: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a three Xilinx Alveo U200 FPGA system, processing 2.7M events with a global clock of 200 MHz

KPI	RAIDER @200 MHZ [4 FPGA, 3CNN]	RAIDER @200 MHZ [4 FPGA, 6CNNs]
time to solution [s]	1.548	1.003
throughput [events/s]	2325478.642	2692072.654
energy to solution [J]	628.258 (283.615 FPGA)	407.512 (184.206 FPGA)
energy efficiency [events/J]	4297.597 (9519.944 FPGA)	6625.572 (14657.541 FPGA)

Table 26: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a four Xilinx Alveo U200 FPGA system, processing 2.7M events with a global clock of 200 MHz

For the IDV-E testbed setup, the integrals of both CPUs and FPGAs power profiles, depicted in Figure 68, Figure 69 and Figure 70, have been used to compute the energy-to-solution and energy efficiency values which are reported in Table 27, Table 28 and Table 29. To note: the energy values reported are obtained - in the case of 3 and 4 FPGAs used - working on a combination of two Ampere Altra nodes with just one provided with the two-phase cooling system developed within the TEXTAROSSA project.

KPI	RAIDER @200 MHZ [2 FPGA, 3CNNs]
time to solution [s]	1.48956
throughput [events/s]	1812622.012
energy to solution [J]	216.170 (133.004 FPGA)
energy efficiency [events/J]	12490.147 (20300.090 FPGA)

Table 27: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a two Xilinx Alveo U280 FPGA IDV-E node, processing 2.7M events with a global clock of 200 MHz

KPI	RAIDER @200 MHZ [3 FPGA, 6CNNs]
time to solution [s]	0.792
throughput [events/s]	3409090.909
energy to solution [J]	197.301 (105.489 FPGA)
energy efficiency [events/J]	13684.675 (25595.171 FPGA)

Table 28: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a three Xilinx Alveo U280 FPGA configuration distributed over 2 IDV-E nodes, processing 2.7M events with a global clock of 200 MHz

KPI	RAIDER @200 MHZ [4 FPGA, 9CNNs]
time to solution [s]	0.554
throughput [events/s]	4873646.209
energy to solution [J]	165.277 (101.055 FPGA)
energy efficiency [events/J]	16336.183 (26718.126 FPGA)

Table 29: [RAIDER] KPIs evaluated on the execution of the RAIDER application on a four Xilinx Alveo U280 FPGA configuration distributed over 2 IDV-E nodes, processing 2.7M events with a global clock of 200 MHz

To summarize and to visualize the trends of the KPIs presented in the previous tables for the two used testbeds, we report throughput and energy efficiency values, respectively on Figure 71 and Figure 72, in order to see at glance the comparison of the performances between APE Lab and IDV-E testbeds. In detail, from the inclination of the linear regression curve obtained from the different setups and for the different KPI, we can notice that the RAIDER application deployed on the U280 multi-FPGA setup (IDV-E) tends to scale better with respect to the U200 one (APE Lab) in both throughput and energy efficiency. This is because, as described in Section 2.6.2, on the IDV-E nodes FPGAs can be implemented more CNN HLS kernels on a single computing node, and so increasing the global computing performance of the hardware.

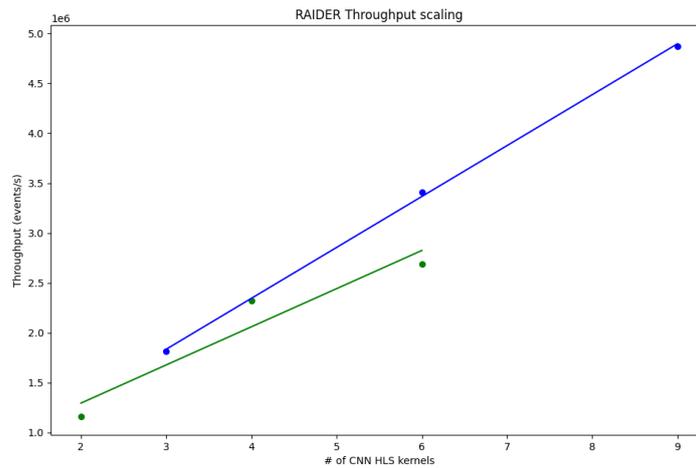


Figure 71:[RAIDER] Throughput scaling trends

(in blue, values from Alveo Xilinx U280 setup; in green, values from Alveo Xilinx U200 setup)

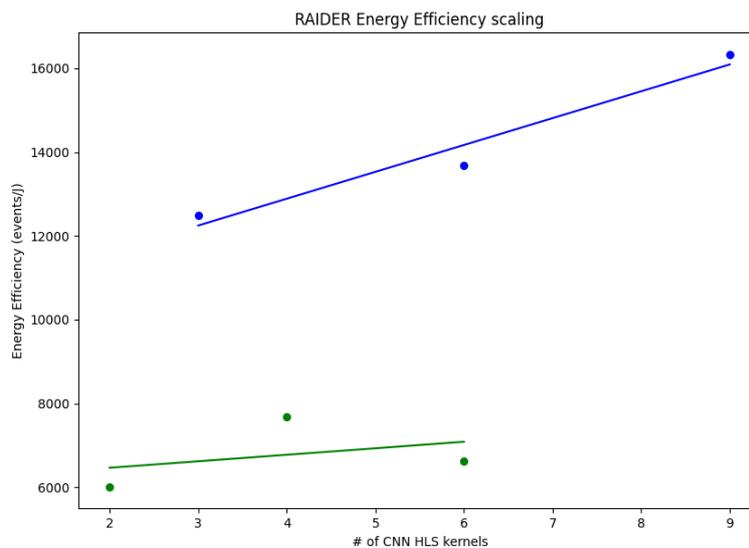


Figure 72: [RAIDER] Energy Efficiency scaling trends

(in blue, values from Alveo Xilinx U280 setup; in green, values from Alveo Xilinx U200 setup)

Finally, in Table 30, we consider the improvement factor over the KPIs baseline (reported in Table 23) with the RAIDER design synthesized at 100 MHz on a single U200 FPGA implementing 2 CNN . For the comparison with both testbeds, we choose to report the values relatives to the higher throughput configuration obtained (since the throughput is one of the main goals of RAIDER application, developed to work in a HEP experiment with a trigger rate of 10MHz).

KPI	Improvement factor over single U200 FPGA (2 CNN) setup	
	APE Lab Testbed (4xU200, 6 CNN)	IDV-E testbed (4xU280, 9 CNN)
throughput	4.884	8.841
energy efficiency	0.657	1.621

Table 30: [RAIDER] Improvement factors of the KPIs for the tested designs over the baseline.

To comment these KPI tables, we can see that a good scaling in term of energy efficiency can be reached only if the added nodes have enough processing capabilities that allow to increase the integrated processing throughput at a higher rate than that of the additional energy consumption they introduce in the system. In fact, as can be seen from Table 30, a quite good improvement factor over the baseline of the RAIDER energy efficiency has been obtained with the IDV-E testbed that integrates Xilinx Alveo U280 FPGA boards.

3.6.3 RESULTS ON TASK+STREAM INTEGRATION ON IDV-E

As part of Task 4.6 SW integration & Optimization we experimented how to integrate a RAIDER 6x6 kernel into OmpSs@FPGA framework. This proof-of-concept implementation was challenging as it implied mixing the task-based model that underpins OmpSs@FPGA with the stream model that is a natural approach to programming the RAIDER kernel. The objective of this exploration was to decide on the feasibility of this mixed model and the potential benefits that it could provide.

Several different implementations were tested. The main results are summarized in Figure 73 that shows the throughput in images processed per second (1 million image per second is equivalent to 1MHz) obtained by the different versions (please note the logarithmic scale of the vertical axis).

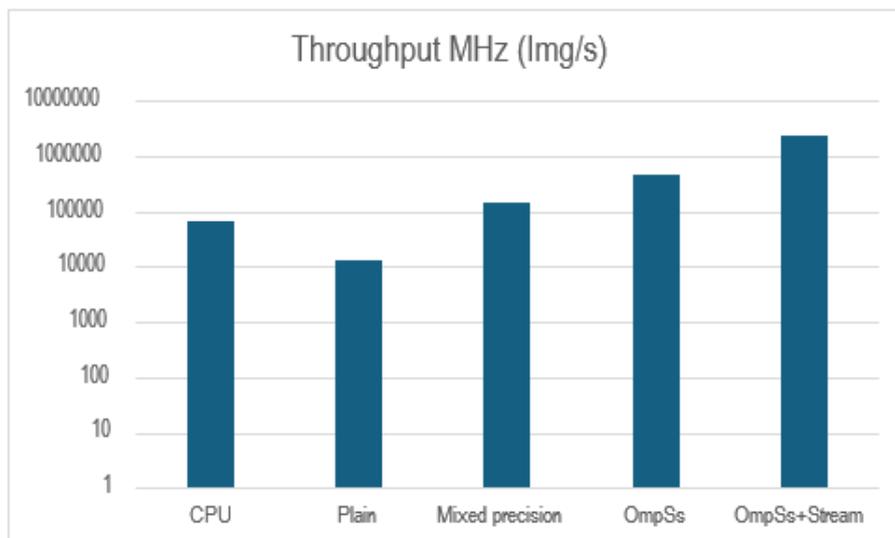


Figure 73: [RAIDER] Throughput with different OmpSs@FPGA implementations

As shown in Figure 73 there were several versions of RAIDER implemented in the OmpSs@FPGA framework. Figure 74 shows how the different versions change the work organization. Also, a description of the features of the main versions follows:

- CPU: This measures the performance of the RAIDER kernel optimized and executed in an Intel Xeon Silver 4208 CPU @ 2.1GHz. It has a throughput of 0.067MHz.
- Plain: This initial version implementation is a direct translation of the code executed in the CPU to the FPGA. In this initial version the OmpSs@FPGA framework is only used to facilitate the communication between the CPU driver program and the FPGA and to measure the FPGA performance. This plain version had a throughput of 0.013MHz so roughly a 5x slowdown over the CPU mainly due to the sequentialization of the image processing and the CPU slow sending of images to process to the FPGA.
- Mixed precision: In order to improve performance, an initial optimization taking advantage of the initial image pixel size was performed to use only the minimum number of resources in the FPGA. This optimization could not be applied to the CPU version as it can only exploit fixed size numbers (i.e byte multiples) but it was as simple as changing some type definitions in the FPGA code. The resulting code not only saved resources in the FPGA but also improved performance due to the lower latency of shorter bit operations and the minimized CPU-FPGA data transfer times reaching a throughput of 0.149 MHz. This version already has a 2x speedup over the CPU version.
- OmpSs: The third reported implementation used the FTS IP developed in Task 2.5 and OmpSs@FPGA to improve the kernel task latency increasing throughput significantly to 0.470MHz. This version managed tasks inside the FPGA but already in a sequential fashion.
- OmpSs+Stream: The final proof-of-concept version used the FTS to send tasks that pipeline the different images among the different processes in dataflow. This significantly improves performance as different images are in process in the same hardware kernel in the different processing steps. This implementation was able to reach 2.44MHz. This version is close to reaching the objective of 10MHz considering that under the OmpSs Framework it is trivial to divide the work into 4 or more kernels (as showed in section 3.12).

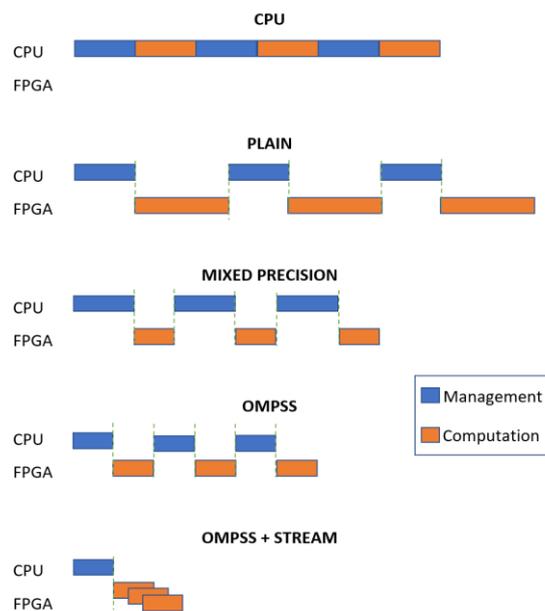


Figure 74: [RAIDER] Different OmpSs@FPGA implementations organization

Taking into consideration that the objective of the work presented in this section was limited to explore the feasibility of joining the task and stream models we consider that the results presented are highly successful as they demonstrate not only the aforementioned feasibility but show a promising venue to leverage it to achieve some significant performance goals. As a consequence of the results shown, more work should be devoted in this direction to integrate the changes into the automatic tools to facilitate development and to evaluate ways to further improve performance and even better mix both models in more intuitive ways.

3.7 TNM – INFN

The tensor network methods (TNM) application combines multiple simulation frontends for simulation quantum systems. For TEXTAROSSA, we consider the following sub-applications:

- Quantum matcha TEA: gate-based quantum circuit emulator for digitized quantum circuits (analyzed in deliverable D6.2).
- Quantum green TEA: solver for the Schrödinger equation or Lindblad equation; within this analysis, we restrict ourselves to finding the ground state of a system.

3.7.1 Towards mixed precision

In the current development goals, mixed precision algorithms are not on the agenda for the tensor network methods. Nonetheless, we take TEXTAROSSA as motivation to compare different precisions for the calculations of the TNM. This data helps us to decide in the future if mixed precision can be beneficial for our application. In deliverable D6.2, we already presented data for Quantum matcha TEA and single-precision versus double precision complex numbers for a quantum circuit. Now, we turn to quantum green TEA, where we added the possibility to run simulations with real data types and switch to higher precision towards the end of search; real data types are suitable as long as the Hamiltonian is also real. The ground

state energy of a system with 128 qubits can be calculated via the Jordan-Wigner transformation ($E=-162.6123001775688(7)$, $J=1.0$, $g=1.0$).

For the comparison, we use the following configuration

- TEXTAROSSA-node Dibona (Atos)
- Quantum green TEA (Version pre_v0.3.17)
- single-cpu simulation
- likwid 5.2.2 and per-cpu energy measurement

Figure 75 and Figure 76 show the same trend for runtime and energy consumption where the double complex data point (black) serves as reference for simulations before the implementations done within TEXTAROSSA. The colormap shows how many of the last sweep are in double precision; not all sweeps have to be executed based on an exit criterion. We observe that using only single-precision sweeps leads to an increased error. Already two double-precision sweeps are sufficient to bring the error on par with the best simulations. The optimal setting here is five single-precision sweeps followed by potentially five double precision sweeps, while actually exiting after the ninth sweep.

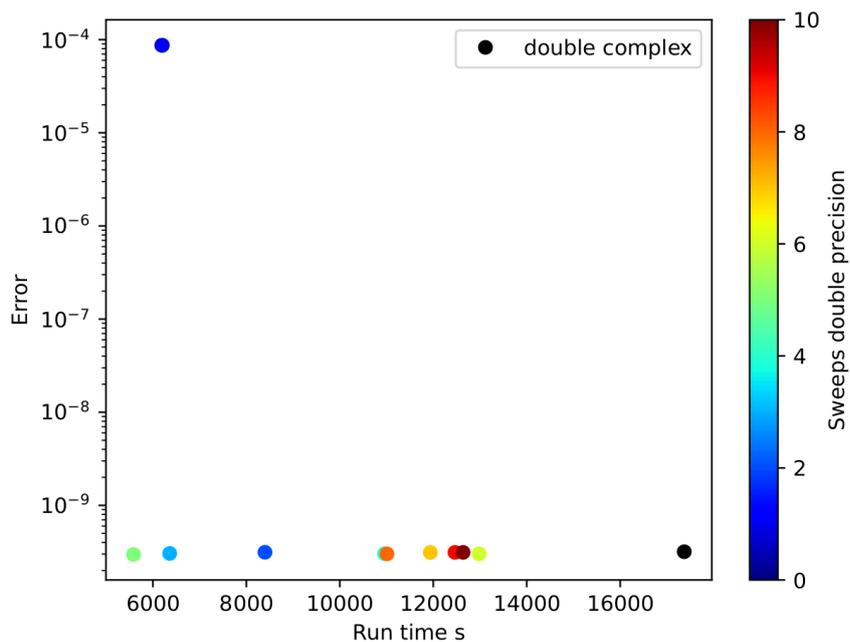


Figure 75: [TNM] Execution time vs error varying the number of final sweeps in double precision

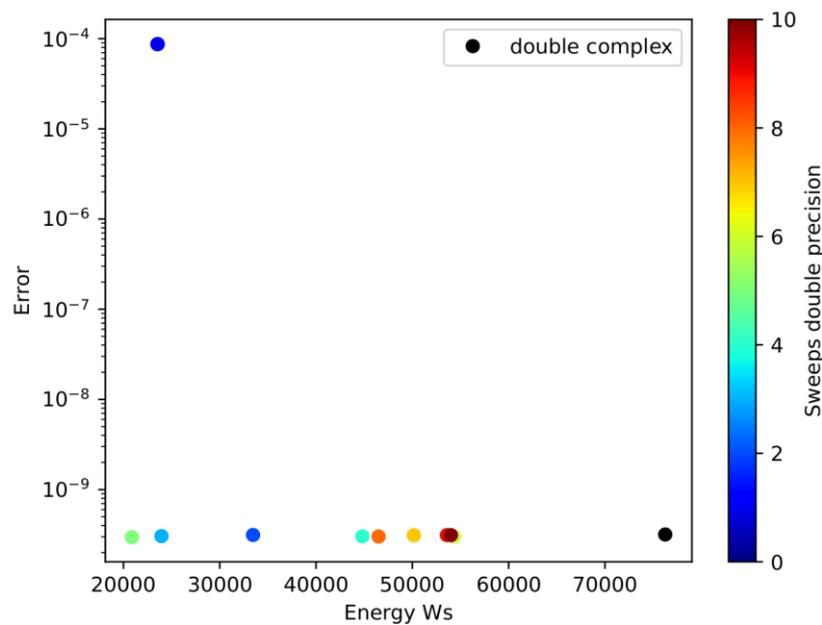


Figure 76: [TNM] Energy consumption vs error varying the number of final sweeps in double precision

In summary, the ability to run with different precisions and increase during runtime the precision is a useful path for quantum green TEA and future implementation, e.g., an automatic increase of the precision independent of the user input.

These changes have been implemented on the Fortran level of the quantum green TEA library.

3.7.2 Hybrid simulation: data parallelism via MPI versus parallelizing a single simulation

On a single node as targeted by TEXTAROSSA, one can find the optimal path between running a batch of simulation in a hybrid MPI-openmp/threading scenario. For the quantum green tea application, we focus on threading via the MKL library. We use the following setup:

- Justus2 cluster (bwHPC): 2xIntel Xeon 6252 Gold with 2x24 cores; we run quantum green tea v0.3.29 on Justus2.
- Batch of 96 identical simulations, overloading node by one MPI thread as one MPI thread is not executing any computational workload in the master-worker approach.
- We have N workers, a total time T for all 96 simulations (wall-time), the time t for a single simulation, and the speedup S for a single simulation. The number of MPI threads is $48 / N$.
- The data points are for the quantum Ising model in one dimension at the quantum critical point.

TNM hybrid simulation MKL threads	Time total T [mm:ss]	Time t for single simulation [mm:ss]	Speedup S for single simulation
1	6:16	3:08	Reference
2	6:17	1:34	1.99
3	7:26	1:14	2.5
4	7:35	0:57	3.3
6	7:43	0:39	4.9
8	18:55	1:11	2.7
12	21:38	0:54	3.5
16	33:30	1:03	3.0
24	51:01	1:03	2.9
48	134:44	1:24	2.2

Table 31: [TNM] Hybrid simulation with MPI and threading via MKL library. We use simulation of the quantum green tea library out of the Quantum TEA suite.

We observe two main trends in Table 31. To minimize the total simulation time, MPI is always the best approach at constant resources. The speedup of the MKL library peaks for this setup at 6 threads. With scalable resources, running 96 MPI threads with each 6 MKL threads across multiple nodes minimizes the wall-time. The minimum of time t as a function of MKL threads might be problem dependent and platform dependent, e.g., which is the size of the generalized matrix problems to be solved in the simulation or the NUMA modes implemented on the architecture.

3.7.3 Device comparison CPU versus GPU

Motivated by the mixed precision results on the Fortran level, we implemented the ground state search of quantum green TEA on the Python side of our library supporting switching precision during the sweeps as well as GPU support. We consider a two-dimensional quantum Ising model and its ground state search via a Tree Tensor Network (TTN). We run on CINECA's Leonardo, single-core for CPU simulations (Intel Xeon 8358 CPU, with 32 cores running at 2.6 GHz), single GPU (NVIDIA custom Ampere GPU, 64GB HBM2) and highlight in the following three figures three aspects:

1. For larger system sizes, the GPU gives a benefit in terms of speedup, see Figure 77;
2. both devices show a speedup when moving sweeps, i.e., iterations of our search, to single precision, see Figure 78;
3. except all sweeps being done in single precision, the same precision for the ground state is reached (lower ground state energy is better), here shown for at least two double precision sweeps at the end, both for CPU and GPU, see Figure 79.

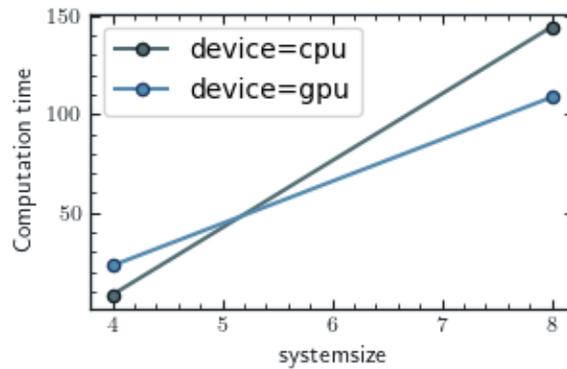


Figure 77: [TNM] Execution time vs system size

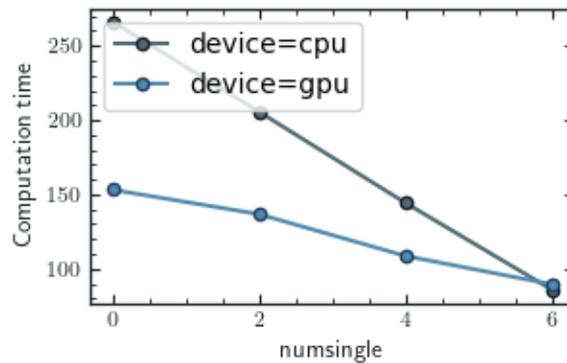


Figure 78: [TNM] Execution time vs number of sweeps in single precision

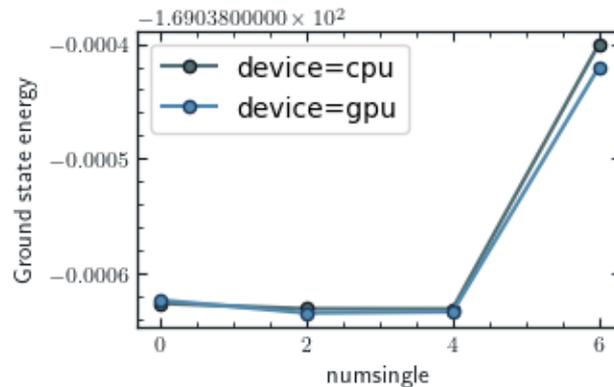


Figure 79: [TNM] Ground energy state (lower is better) vs number of single precision sweeps

3.7.4 Algorithmic improvements for sampling

In addition to the above studies on aspects of different precision and parallelization, we have undertaken some algorithmic improvements which lead to an improvement of the energy consumption of the library via the improved computation time. The following example is on the problem of sampling the classical outcomes of a quantum state represented as a tensor network, which are discussed in detail in M. Ballarin's work [11]. Out of the various examples, we pick the analysis of a tree tensor network representing the ground state of the quantum Ising model while we try to reach a given threshold of probability coverage. Our new algorithm has

shown a speedup, but we verify that the speedup is also reflected in the energy consumption on Dibona node. The key data are:

- TEXTAROSSA Dibona node (Atos)
- qtealeaves python library
- single-core simulation
- likwid 5.2.2 and per-core energy measurement

Figure 80 and Figure 81 the speedup of the new algorithm in computation time and the same ratio in terms of energy consumption.

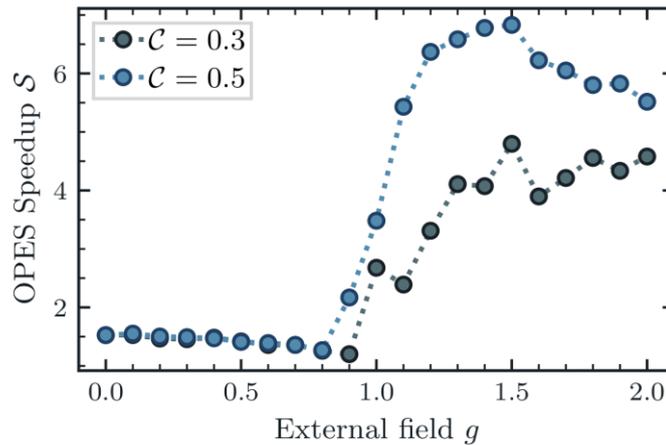


Figure 80: [TNM] Speedup after algorithmic improvements

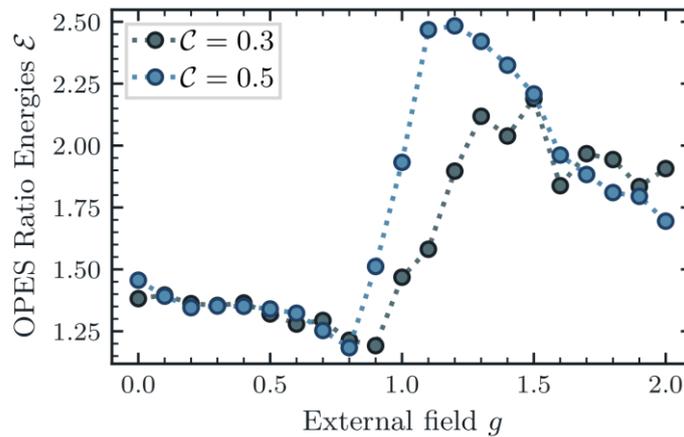


Figure 81: [TNM] Improvement ratio in energy efficiency after algorithmic improvements

3.7.5 Algorithmic improvements for Hamiltonian representation

We were furthermore able to optimize the algorithm by how we represent a Hamiltonian (matrix) and apply it to a state (vector), i.e., a generalized matrix vector multiplication. This application is the foundation of the ground state search and time evolution within quantum green TEA. The results show the relative computation time of a ground state search with respect to the best method, i.e., the icTPO (indexed compressed tensor product operator), see Figure

82. The tensor product operator (TPO) and our starting point before the TEXTAROSSA project are about three times slower in the example system; the indexed tensor product operator (iTPO) and an alternative algorithm running as reference point (sparse-matrix product operators (SPO) and indexed SPO (iSPO)) are about two times slower than the icTPO. The results are obtained for random all-to-all two-body interactions in a quantum Ising type of Hamiltonian with additional randomized local fields. The simulations run on the Leonardo machine at CINECA, as single-CPU simulation. Future steps are to benchmark how algorithmic improvement develops in a parallel algorithm.

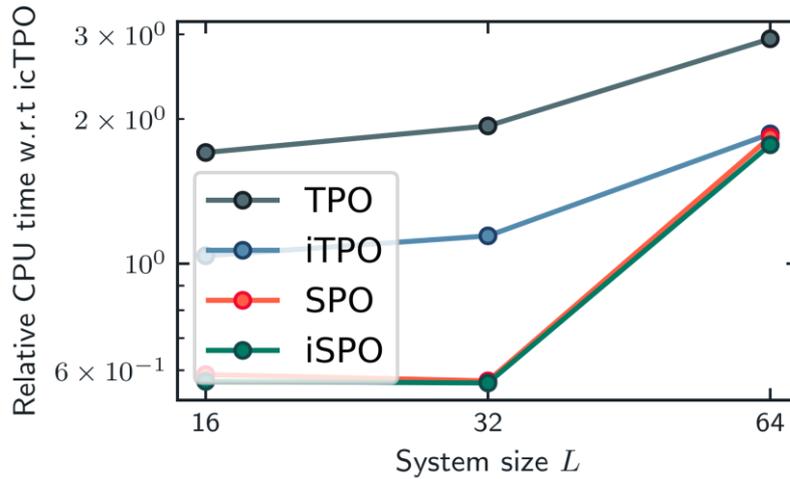


Figure 82: [TNM] Relative CPU time w.r.t icTPO

3.7.6 Key performance indicators on IDV-A

We re-run the KPI measurements for quantum matcha tea on IDV for the quantum Fourier transformation (QFT) on 100 qubits as done in D6.2. Due to the different setup of the two machines, the way of measuring the energy consumption differs. In contrast to D6.2 and the measurement of the energy consumption of a single core, we now use turbostat and measure the energy consumption over the whole machine. We measure for about 15 minutes for each simulation as well as for the idle machine. We use the additional power consumption with respect to the idle machine as input to calculate the KPI. We obtain the results according to the Table 32.

	6-CFF	6-ZFF5	6-ZFF	8-CFF	8-ZFF5	8-ZFF	Idle
Single run bash [s]	9	13	25	56	103s	267	n.a.
Turbostat runs [s]	842	963	947	903	923s	1069	900s
Repetitions turbostat	90	75	38	16	9	4	n.a.
Average turbostat [s]	9.4	12.8	24.9	56.4	102.6	267.3	n.a.
Turbostat Power [W]	406.08	406.26	407.43	407.51	408.62	409.46	404.09
Energy per run [Ws]	18.6	29.1	83.2	193	464.8	1435.0	n.a.
Gates/s [1/s] (KPI)	1266.1	923.1	475.5	194.9	107.2	41.2	n.a.
Gates/Ws [1/Ws] (KPI)	636.2	406.7	142.4	57.0	23.7	7.7	n.a.

Table 32: [TNM] Measurement of key performance indicators for a single-core simulation with Quantum matcha TEA on IDV-A. The simulations repeat the analysis executed for D6.2 on the most recent machine, i.e., IDV-A.

The simulation string N-XFF in the table refers with N to the number of entangled blocks and with X to the data type (C: single complex, Z: double complex, see also D6.2). A postfix 5 refers to simulations run with double complex, but we truncate singular values at the $1e-5$ level equal to the single complex simulations. This way, CFF and ZFF5 are better comparable in terms of the matrix dimensions encountered in the simulation for the linear algebra operations.

We see a performance increase in time-to-solution with respect to the previous machine, which leads to a better KPI for the gates per second. The KPI for the energy consumption indicates an improvement as well, although the measurements with different tools might hide more details.

3.8 ScalFMM (Mathlibs-INRIA)

The Fast Multipole Method (FMM) is a highly efficient algorithm for computing long-range forces in N-body simulations, which are prevalent in fields such as astrophysics, molecular dynamics, and electrostatics. The traditional FMM algorithm reduces the computational complexity of calculating pairwise interactions from $O(N^2)$ to $O(N \log N)$ or even $O(N)$, depending on the implementation and the desired accuracy.

TBFMM - a fork of ScalFMM - is a high-performance task-based FMM implementation where the algorithm is decomposed into a series of discrete tasks, each representing a portion of the computation. These tasks are defined with clear dependencies, forming a Directed Acyclic Graph (DAG) that represents the execution flow of the entire computation. This DAG includes tasks for operations such as particle-to-particle interactions, particle-to-multipole translations, multipole-to-multipole translations, and more, depending on the specific version of the FMM algorithm being used.

In the TEXTAROSSA project, we developed a new scheduler to decide how the tasks should be distributed on the different processing units. Those processing units are usually heterogeneous with a mix of CPUs/GPUs (IDV-A) or CPUs/FPGAs (IDV-E). Our scheduler is detailed in WP4, and WP6 is tied to its use and analysis.

We provide the results in Figure 83. We executed a simulation of 1 million particles on two hardware configurations, Intel CPU + NVIDIA V100 and AMD CPU + NVIDIA A100 (similar to IDV-A). We compared our scheduler against two state-of-the-art alternatives. Our scheduler, MultiPrio, demonstrates very good results. In a more in-depth study, we showed that our scheduler is especially competitive when the tasks are irregular (i.e., of various granularities).

Thus, MultiPrio represents an innovative approach to dynamic task scheduling that significantly benefits the execution of TBFMM by reducing the overall computation time while maximizing hardware utilization.

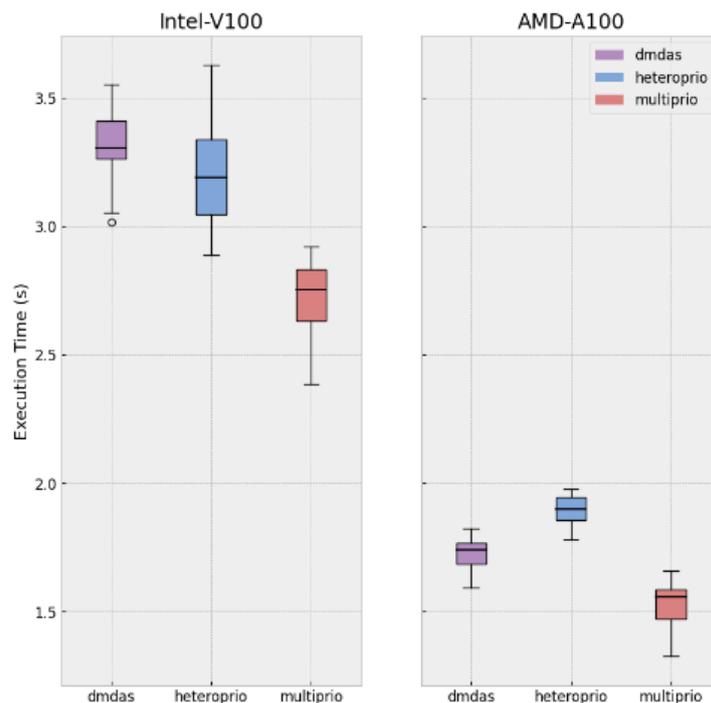


Figure 83: [FMM] Performance results for TB-FMM on two hardware configurations
The execution time is provided (the lower the better)

3.9 Chameleon (Mathlibs-INRIA)

Chameleon is a software library designed to optimize the execution of dense linear algebra computations on heterogeneous computing systems. It provides high-performance implementations of various linear algebra operations, adapting to the underlying hardware to achieve efficient execution. The library leverages task-based programming models to decompose complex operations into smaller tasks, facilitating their scheduling and execution across different types of processing units, such as CPUs and GPUs. This approach allows Chameleon to dynamically balance the workload and utilize available hardware resources effectively.

Within the Chameleon library, key kernels such as GEQRF, POTRF, and GETRF are implemented to support a range of linear algebra computations:

- **GEQRF (QR Factorization):** This kernel computes the QR factorization of a matrix. QR factorization is a process that decomposes a matrix AA into the product of an orthogonal matrix QQ and an upper triangular matrix RR . The GEQRF kernel is crucial for solving systems of linear equations, eigenvalue problems, and many other applications in scientific computing. In the context of Chameleon, the GEQRF operation is broken down into tasks that can be executed in parallel, exploiting data locality and minimizing data movement for improved performance on heterogeneous systems.
- **POTRF (Cholesky Factorization):** The POTRF kernel performs the Cholesky factorization of a symmetric positive-definite matrix AA , decomposing it into the

product of a lower triangular matrix LL and its transpose $LTLT$. Cholesky factorization is particularly useful for solving linear systems where the matrix is symmetric and positive-definite, such as in optimization problems and numerical simulations. Chameleon optimizes the execution of POTRF by scheduling tasks efficiently across available resources, taking advantage of the specific characteristics of the matrix to enhance performance.

- **GETRF (LU Factorization):** This kernel computes the LU factorization of a general matrix. LU factorization decomposes a matrix AA into the product of a lower triangular matrix LL and an upper triangular matrix UU , possibly with a permutation matrix PP due to partial pivoting. GETRF is a fundamental operation in linear algebra, enabling the solution of systems of linear equations, inversion of matrices, and computation of determinants. In Chameleon, the GETRF operation is managed through a task-based approach, allowing for concurrent execution of tasks and efficient use of heterogeneous computing resources.

As for ScalFMM, we used our scheduler on Chameleon and provided the result in Figure 84. We compared MulTiPrio against DMDAS and heteroprio, on Intel CPU + NVIDIA V100 and AMD CPU + NVIDIA A100 (similar to IDV-A). Our new scheduler is competitive with DMDAS but is slightly slower. The main reason is that DMDAS use high-tuned priorities provided by the developers (which for example inform the scheduler about the critical path without the need of analysis). Our scheduler is able to use them too, but in the given results we use a version where our scheduler tries to find the critical task autonomously.

We provide two execution traces, see Figure 85 and Figure 86, to illustrate the difference. We observe that when priorities are utilized, the number of ready tasks is well-controlled, meaning that the most critical tasks are computed first.

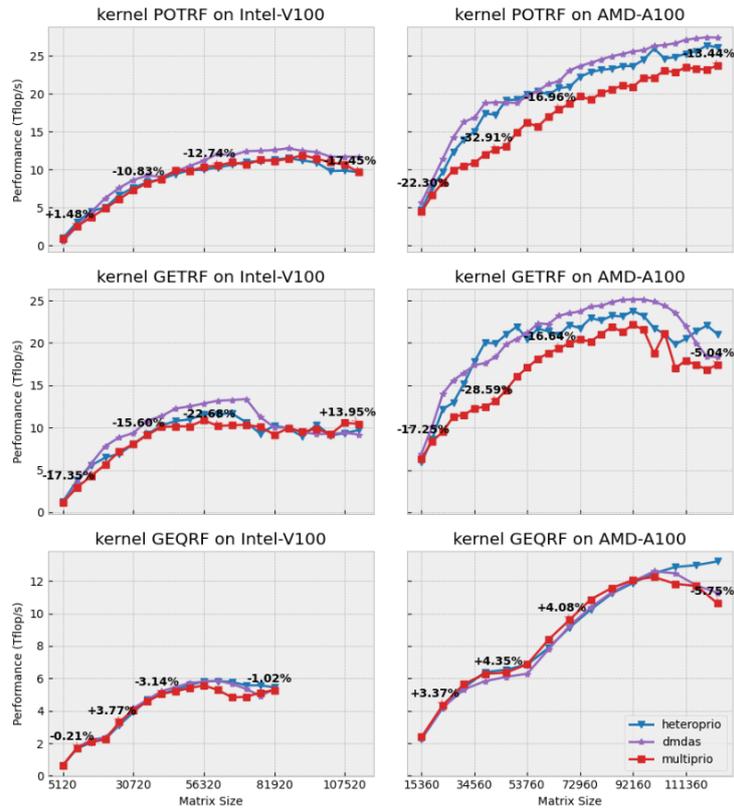


Figure 84: [CHAMELEON] Performance results for Chameleon on two hardware configurations Flops per second are provided (the higher the better)

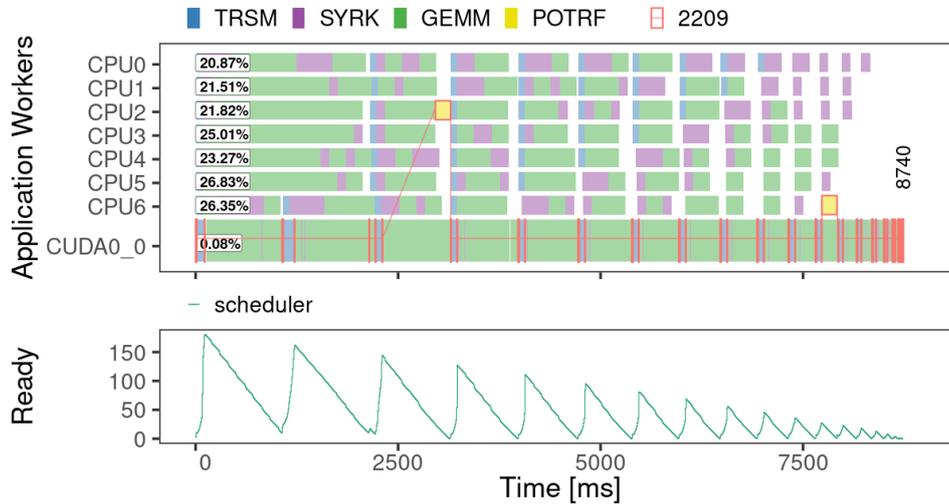


Figure 85: [CHAMELEON] Execution trace (emulation) for a Cholesky factorization without using user's priorities

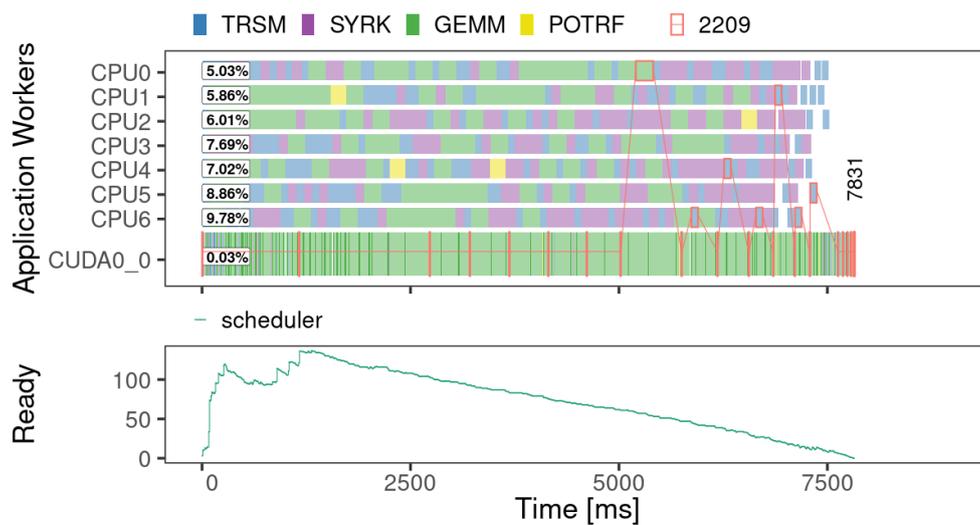


Figure 86: [CHAMELEON] Execution trace (emulation) for a Cholesky with user’s priorities

3.10 StarONNX (INRIA)

StarONNX is a framework for performing inference requests on heterogeneous computing platforms, for any neural network model. StarONNX is based on architecture-specific optimized kernels provided by ONNX Runtime, and on the task-based runtime system provided by StarPU. The asynchronous, task-based approach ensures automatic and efficient data management and overlapping of computations and communications. In addition, this allows for distributing the computation of a single inference request over several heterogeneous processing units: when each processing unit performs the computation for which it is best suited, we obtain a more efficient execution and a better tradeoff between throughput and latency.

We installed StarONNX on the IDV-A node of TextaRossa to evaluate its ability to efficiently manage the 100 CPU cores and 4 GraceHopper GPUs. We use version 1.16.3 of ONNX Runtime, with the CUDA provider based on cuDNN version 8.9.7, and the development version 1.4 of StarPU (commit 0ced6f9a). We perform all tests using the GoogleNet network.

We consider a setting with periodic inference requests being sent to the framework, with a given input rate, or *throughput*, expressed in number of inferences per second. For a given value of throughput, the system may or may not be able to process requests. We do not consider the possibility of dropping queries, so there is a maximum value for the throughput beyond which the system will not be able to process queries. For a fixed throughput, a relevant metric for comparing inference systems is their induced maximum latency, i.e., the maximum interval duration between the arrival of a query and the end of its processing. In practice, for each batch, the maximum latency is obtained for the first data item in the batch.

To efficiently process inference queries and optimize the efficiency of computational resources (especially GPUs), queries are grouped into batches, whose size needs to be determined. The rationale is as follows: increasing the batch size increases the batch processing time and the

average wait time for queries before batch processing, leading to a higher latency. However, increasing the batch size increases the computational efficiency of the GPU, so the ratio of batch processing time to batch size decreases as the batch size increases, which provides a better throughput.

In the plots of both figures, each point corresponds to the average value, for 800 batches, of the (maximum) latency of the first data item in the batch, and the error bars are computed over the 800 batches. The Figure 87 shows latency measurements for a fixed batch size: latency decreases as throughput increases, up to a point where throughput is so high that the next query arrives before the previous one has finished. This leads to congestion, and the maximum latency diverges as the total number of queries increases. This highlights the need to choose a larger batch size to avoid overloading the GPU and to maintain a bounded latency.

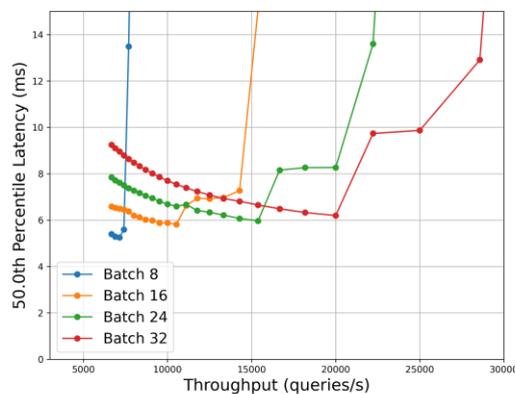


Figure 87: [StarONNX] Performance plots for StarONNX with the GoogLeNet network showing latency as a function of throughput, for different batch sizes

In Figure 88 we show how to minimize the maximum latency for a query. Among all possible batch sizes, we choose the one that minimizes the maximum latency. The dotted plots in grey correspond to the results obtained with different batch sizes, and the blue plot shows the best choice of batch size (the optimal batch size is shown in the blue plot) as a function of throughput.

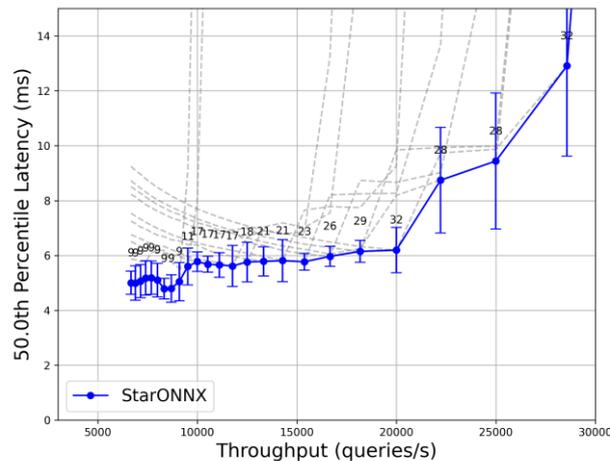


Figure 88: [StarONNX] Performance plots for StarONNX with the GoogLeNet network showing latency as a function of throughput: for each throughput, the batch size that minimizes maximal inference time for a query is chosen

3.11 UrbanAir - PSNC

As described in previous deliverables, selected UrbanAir kernels were adapted to benefit from usage of heterogeneous resources, exploiting CPUs and NVIDIA GPUs. While results are provided for strong scalability, the usual scenarios for assessing air quality run on domain of size exceeding the memory resources of a single GPU, therefore we focus on weak scalability. The toolchain includes C++ compiler, the CUDA toolkit, threads (OpenMP) for shared-memory parallelization (single node) and MPI for data exchange between GPUs, all available at IDV-A platform. To monitor energy consumptions of kernels running on GPUS, GPower project tool was used, also available on IDV-A TEXTAROSSA platform. As previously, the tests were conducted on PSNC Altair machine equipped with NVIDIA V100 and the newly available IDV-A equipped with NVIDIA H100.

The UrbanAir-gcrk starts with solver initialization, followed by a guess of initial wind velocity, estimation of pressure and initialization of boundary condition. Eventually, iterative solver is started where for each iteration reduction and preconditioner (to speedup solver) subroutines are called, and Laplacian operator is solved. Every sub-routine is provided with a separate implementation for CPU and GPUs, so that within a compilation user can decide whether to run on CPUs or GPUs or use a mix of them.

The global domain is divided between the workers, and each receives a 3D block called subdomain, which size is determined by the hardware used. The domain decomposition is done manually before the compilation for additional compiler optimizations. On a single node with CPUs only, OpenMP is used to exploit shared memory. To exchange data between the nodes, and between GPUs, MPI paradigm is used. On order to achieve the best performance on GPUs, the size of subdomain on GPUs is the maximum possible (memory constraint). The implementation allows to execute simultaneously on CPUs and GPUs – on GPUs maximum

possible size of domain is assigned, while the remaining is solved on CPUs – the latter need to be carefully chosen so that parts executing on GPUs do not have to wait for iteration done on

CPUs. However, on currently available fat nodes with multiple GPUs, the most efficient execution is to use accelerators and use CPU only to handle communication between GPUs. Once the domain is divided into subdomains, it is then additionally divided into blocks for optimal cache utilization. As for communication, it is overlapped with computation: computation of boundary conditions is separated from computation of inner domains, so that the data at boundaries is exchanged with neighbors while the inner domain is still computed.

UrbanAir-gcrk iterates over whole domain doing stencil computations. Every subdomain is assigned to exactly one GPU and after each iteration the data between subdomains is exchanged using additional HALO cells (boarders of the subdomain). To increase performance on NVIDIA V100 and newly emerged H100 GPUs cards, and to support wind-flow-specific settings, additional improvements were provided: i) additional boundary conditions /data were placed in shared memory to speedup, ii) further optimization of Laplacian operator implementation for GPUs was introduced, iii) number of communication in *rhsdiv* kernel was limited, iv) further optimal size of the block of threads within subdomain for each sub-kernel is selected based on some auto-tuning (kernel compile – run – compare results). After each iteration, a global sum of variables is calculated – reduction algorithm was improved to benefit from shared memory in multi-GPU environment. With these changes provided, further up to 10% decrease of time-to-solution is observed.

Adaptation to GPUs and optimization is a game changer, not trivial though, in terms of time-to-solution decrease, as reported already in D6.2. For applications where data to be computed can be equally divided between the workers, it is easier to adapt as the workload for each GPU is the same and developer needs to focus on providing implementation for an accelerator and how to just exchange data between subdomains. As the previous results and those in following subsections have described, executions on GPUs are superior to CPUs, though the best efficiency can be achieved only with a GPU fully loaded with data to be computed.

3.11.1 Multi-node environment

To test UrbanAir-gcrk in multi-node environment, PSNC's Altair HPC system was used, with up to 6 nodes available, each equipped with 8 NVIDIA V100 GPU cards. Two tests were conducted based on strong and weak scalability. Figure 89 presents the strong scalability for 5M grid points problem size. With more GPUs added, the number of iterations solved within one second increases. Provided optimizations allow to reach up to 11% increase in performance, which eventually drops when more GPUs are used, because the problem size per GPU becomes too small to efficiently utilize GPU. Figure 90 presents weak scalability for a fixed problem size – 5M grid points per GPU. Although the number of iterations solved initially drops when more GPUs are added, it then oscillates close to 40 for the remaining GPUs, making it possible to solve 240M grid points problem within the same amount of time.

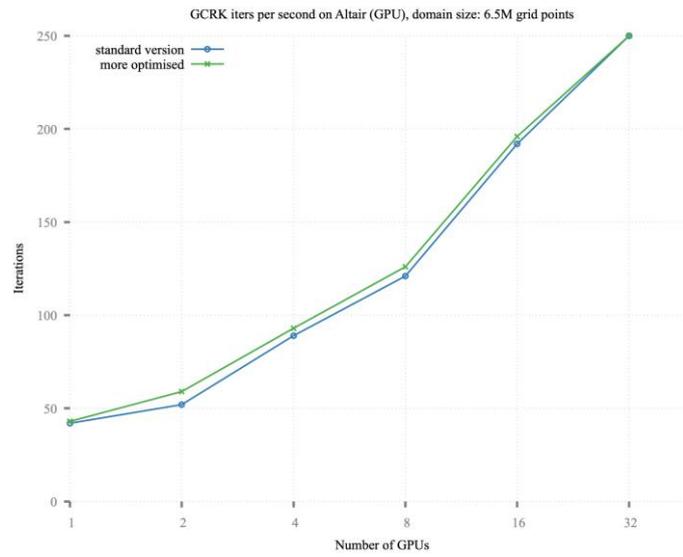


Figure 89: [UrbanAir] Iterations per second (problem size: 6.5M)

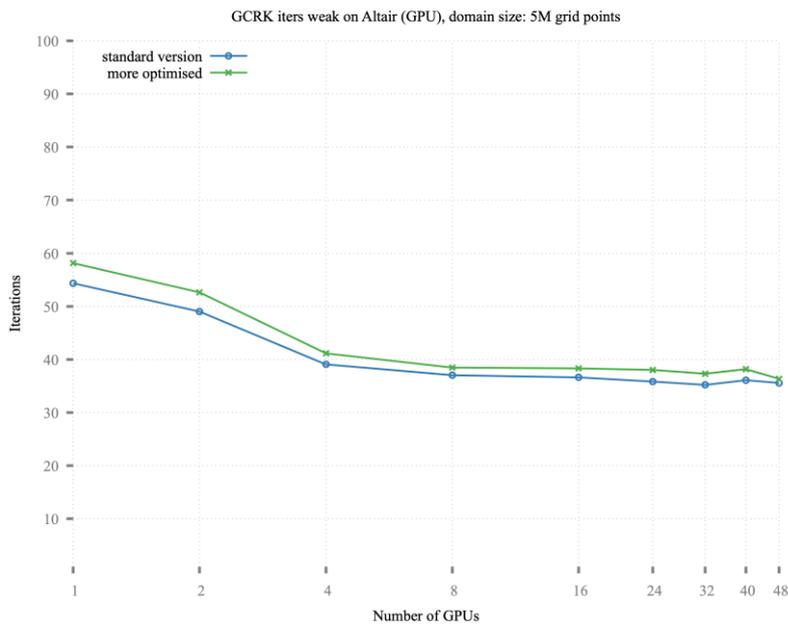


Figure 90: [UrbanAir] Iterations per second (5M per GPU)

The 5M grid points per GPUs is suboptimal because the GPU is not fully loaded with data it computes. Figure 91 presents comparison of achieved speedup between 5M and 50M grid points (maximum size on NVIDIA V100 due to memory constraints). It is clearly visible that the larger the problem size is, the better speedup can be achieved, thus obtaining results in a shorter amount of time. Similarly, weak scalability of UrbanAir-gcrk is better with GPUs fully loaded with data. Figure 92 presents efficiency comparison between 5M and 50M grid points per GPU. While efficiency is better for a maximum executable problem size, it decreases with more GPUs used, and thus the time-to-solution becomes slightly longer. Nevertheless, with GPUs we are able to run UrbanAir-gcrk over much larger areas within a reasonable amount of time.

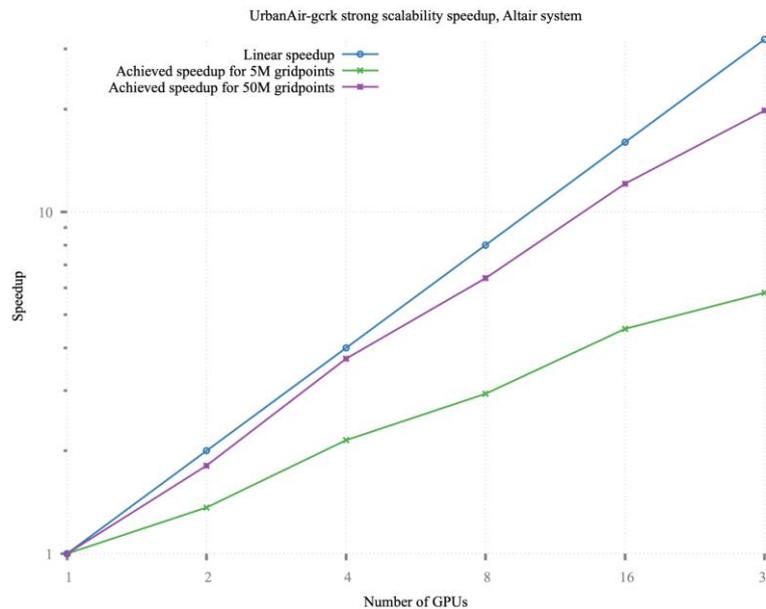


Figure 91: [UrbanAir] Speedup comparison for different domain sizes

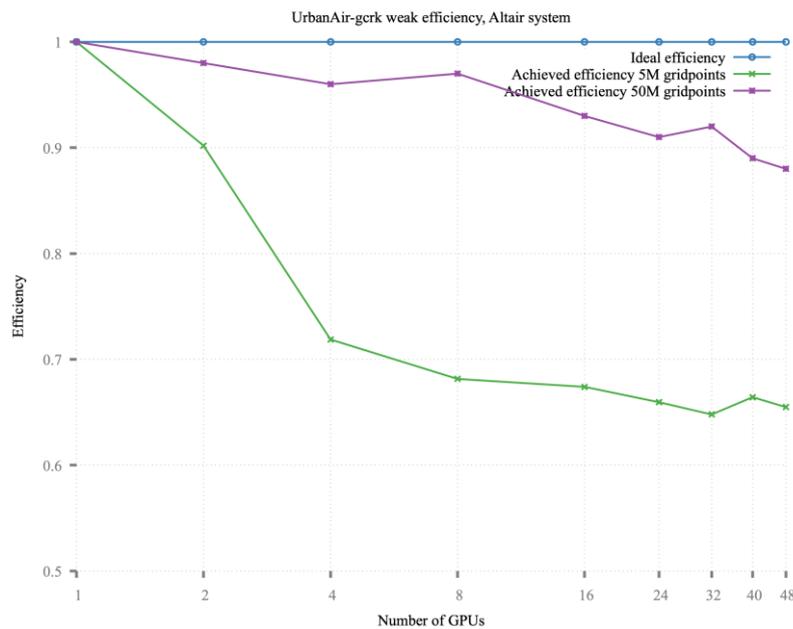


Figure 92: [UrbanAir] Efficiency (weak speedup) comparison for different domain sizes

3.11.2 IDV-A results

Table 33 Table 33: [UrbanAir] KPIs overview presents KPIs for UrbanAir-gcrk, which are then discussed in the following subsections. In previous deliverable, D6.2, initial baseline results were discussed, comparing execution on CPUs vs GPUs. PSNC Altair system and Dibona (predecessor of IDV-A) were used.

KPI for computational efficiency	KPI for energy
<ul style="list-style-type: none"> - Time-to-solution - (strong and weak) speedup - number of iterations / second 	<ul style="list-style-type: none"> Iterations / Watt

Table 33: [UrbanAir] KPIs overview

3.11.2.1 Computational efficiency

NVIDIA H100 GPU cards available at IDV-A node allows to run larger UrbanAir-gcrk problems, comparing to V100 available on Altair PSNC system: 59M grid points instead of 50M. More important difference comes with performance, for the same problem size (50M) the H100 GPU cards offer 2x decrease in execution time compared to V100 GPU, see Figure 93.

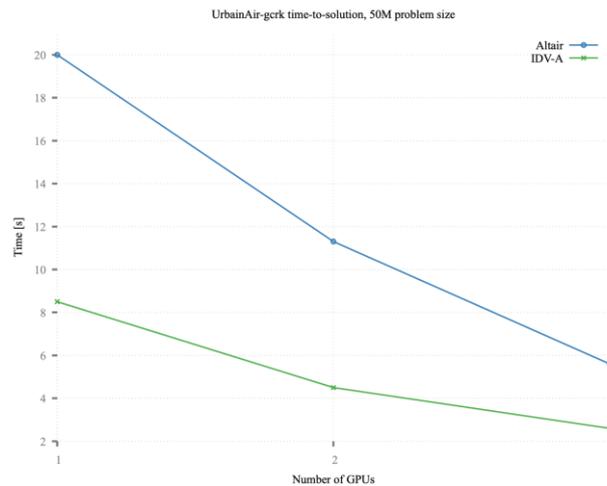


Figure 93: [UrbanAir] Time-to-solution on GPUs: V100 vs. H100

Figure 94 compares speedup of UrbanAir-gcrk across different systems, configurations and code optimizations. Dibona is the former IDV-A system, results for which were presented in D6.2. On final IDV-A, three different configurations are compared:

- IDV-A with traditional cooling system (wo-2phase)
- DV-A with two-phase cooling system installed (w-2phase),
- IDV-A with two-phase cooling system on which optimized version of UrbanAir-gcrk was run (w-phase-opt).

The speedup characteristic is much the same for each configuration, with slight decreases as more GPUs being used. The worse result is observed for w-phase-opt, while the best speedup is achieved on Dibona and wo-2phase. While w-phase-opt provides the worse speedup, it achieves the best iterations/s KPI, see Figure 95. It is worth mentioning that thermal controller of the 2-phase cooling is actively monitoring GPUs for its load to ramp-up their frequency whenever required. Therefore, it provides some overhead, or rather a very slight performance decrease of applications running on GPUs. In case of UrbanAir-gcrk, it impacts time-to-solution the more, the more GPUs are used. Provided optimizations mitigate this impact. Compared to the initial version of IDV-A, Dibona, 10% increase in number of iterations per second is observed.

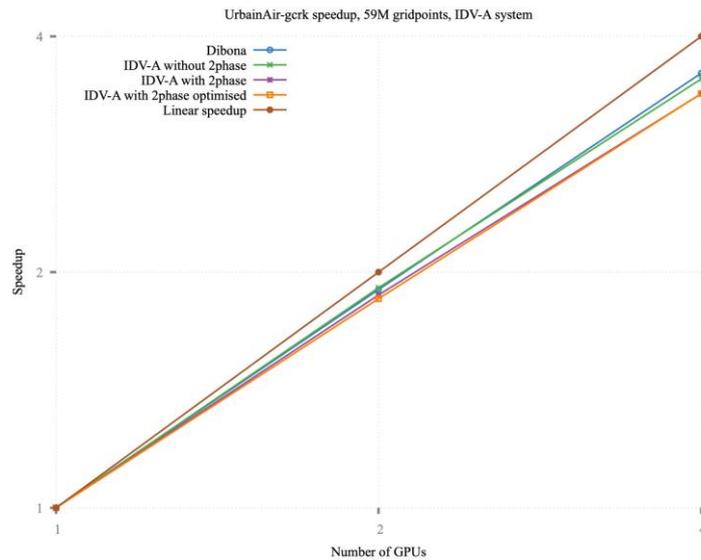


Figure 94: [UrbanAir] Strong speedup comparison for 59M grid points

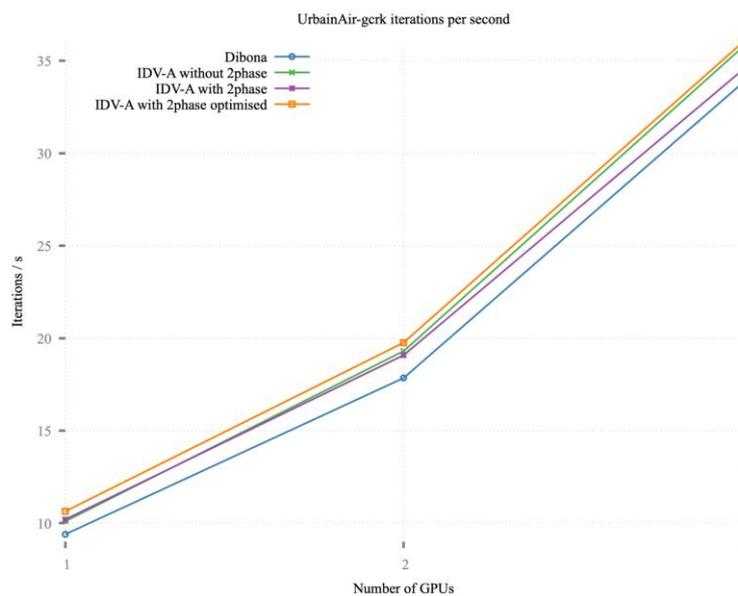


Figure 95: [UrbanAir] Iterations per second for 59M grid points

3.11.3 Energy efficiency

To measure energy efficiency, the methodology proposed in D6.2 is used. Figure 96 details comparison of energy consumption of UrbanAir-gcrk run for 59M grid points domain (strong scalability) for different number of GPUs (rows in Figure) and three different systems (columns in Figure): Dibona (previous version of IDV-A, left column), IDVA-wo-2phase (IDV-A without 2 phase cooling system, left middle column) and IDVA-w-2phase (with two phase cooling system installed, right middle column) and IDVA-w-phase-opt (two phase cooling system installed, additional optimizations of kernels). For each system, adding more GPUs results in less power consumption of each, which is due to fewer grid points to be computed by each GPU. Dibona, IDVA-wo-2phase and IDVA-w-2phase have the same characteristic of

energy consumption – after reaching the highest level of energy utilization, power draw remains more or less constant during the entire execution. It is worth noticing that power draw during execution is 7% higher compared to Dibona. In contrast, additional optimizations introduced (IDVA-w-2phase-opt) provides different power characteristics. It starts with reaching the highest level of energy consumption and almost immediately the power draw drops to an even lower level than observed on Dibona. Taking this into account, and that time-to-solution is shortened when running on IDVA-w-2phase with optimized code, the consumed power should be the smallest among others.

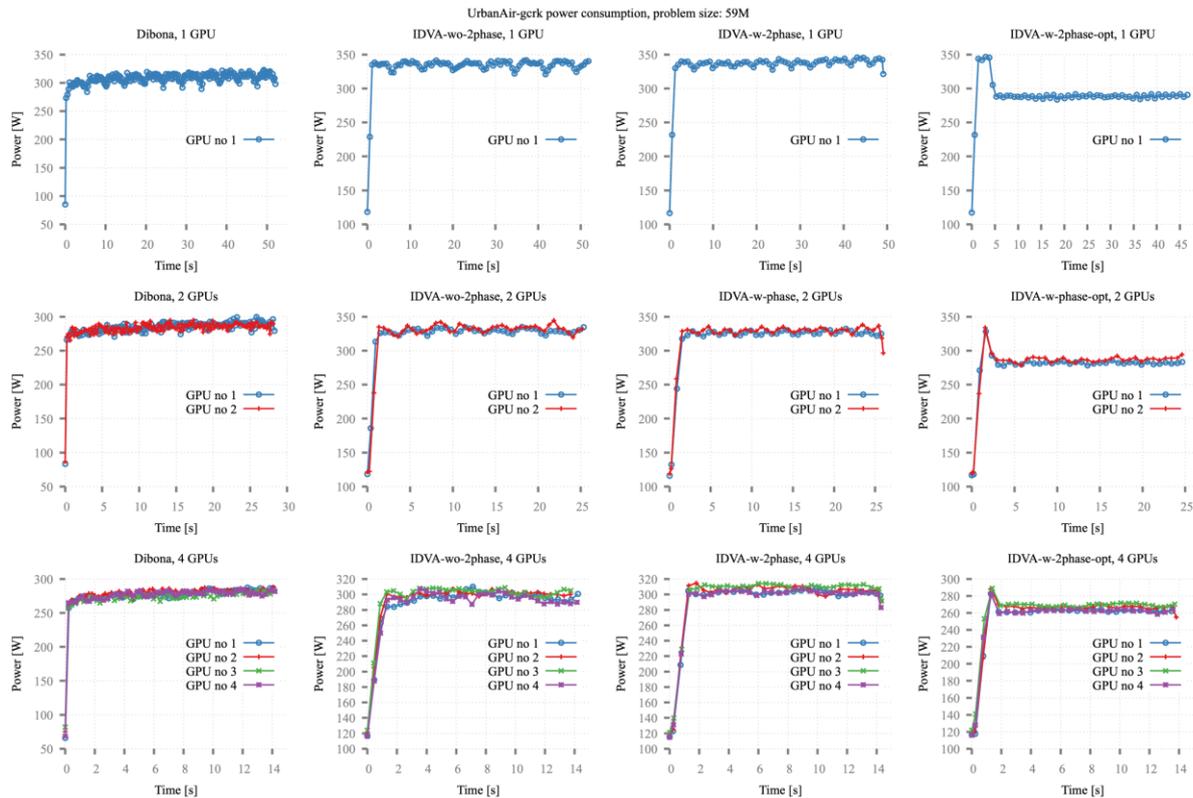


Figure 96: [UrbanAir] Energy usage characteristic for 59M grid points

Figure 97 presents comparison of energy used for running UrbanAir-gcrk on different systems. As discussed above, runs on Dibona were the most energy efficient until more optimizations were introduced to UrbanAir-gcrk. Please note that only GPU power usage being used in computation is taken into account. The CPU power usage is neglected as the CPU is only used to orchestrate data exchange between GPUs. Figure 98 presents energy usage when all GPUs available at the node are taken into consideration, i.e. for a single GPU run, energy usage of all 4 GPUs available are summed. As one can expect, the most energy efficient execution is when all available GPUs are used for computations. It is important to notice that introduced optimizations resulted in the most energy-efficient run on IDV-A .

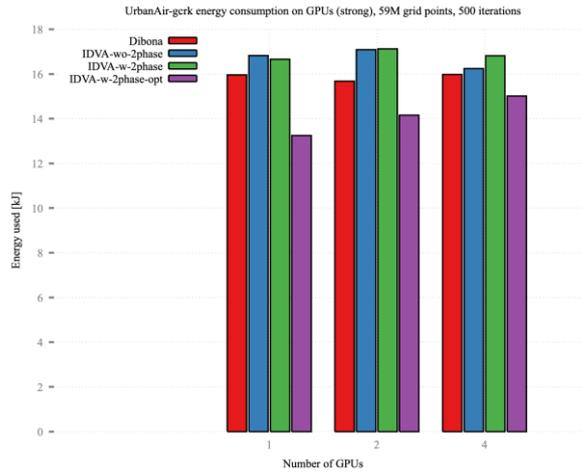


Figure 97: [UrbanAir] Energy consumption for 59M grid points, only used GPUs accounted

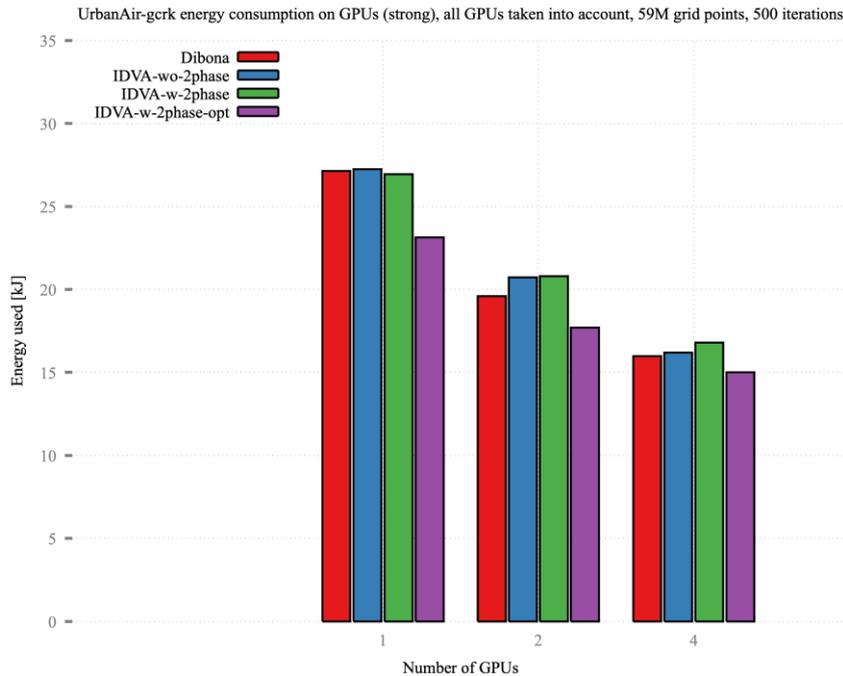


Figure 98: [UrbanAir] Energy consumption for 59M grid points, unused GPUs accounted

UrbanAir-gcrk scales weakly very well when entire GPUs is occupied by data, as discussed in previous subsections. Figure 99 compares IDV-A with and without optimization with respect to time-to-solution (top), and energy consumption (bottom) for a fixed problem size per GPU (59M grid points). With the new version, 4% decrease in time-to-solution is observed, and 22% decrease in energy consumed (when all 4GPUs are used).

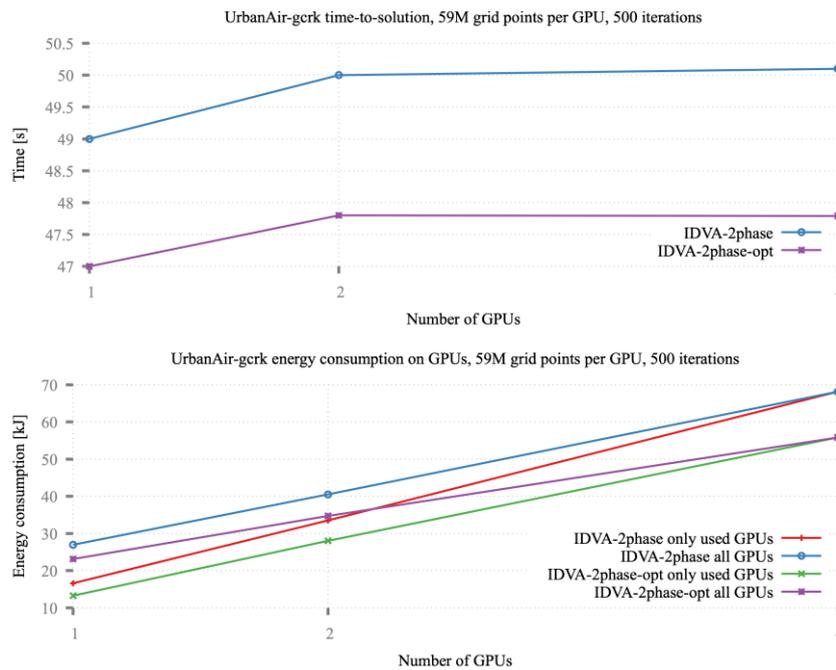


Figure 99: [UrbanAir] Time-to-solution and energy consumption for 59M grid points per GPU

Figure 100 presents the number of iterations per kW for different number of GPUs and systems. When only used GPUs are taken into account, the optimized version running on IDV-A with two phase cooling installed is able to compute the highest number of iterations. The lowest number of iterations per kW is solved with regular version running on IDV-A with or without 2-phase cooling. Similar results are obtained when all GPUs, even unused ones, are taken into account of energy consumption.

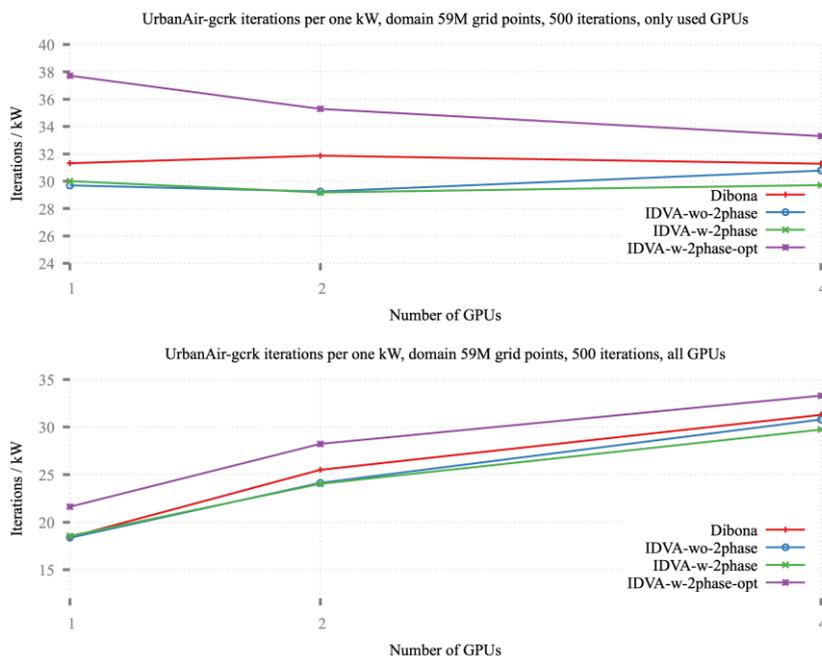


Figure 100: [UrbanAir] Iterations / kW

3.11.4 KPIs

UrbanAir-grck KPIs are discussed in the previous subsections, and are here summarized in tabular form:

- For a fixed problem size (strong scalability):
 - o Table 34 – Time-to-solution
 - o Table 35 – Iterations /s
 - o Table 36 – Iterations / kW
 - o Table 37 – Iterations / kW (unused GPUs accounted)
- For a fixed problem size per GPU (weak scalability):
 - o Table 38 – Time-to-solution
 - o Table 39 – Iterations / s
 - o Table 40 – Iterations / kW
 - o Table 41 – Iterations / kW (unused GPUs accounted)

KPI	Time-to-solution [s]			
GPUs	Dibona	IDVA-wo-2phase	IDVA-w-2phase	IDVA-w-2phase-opt
1	53.2	49.5	49.0	47.0
2	28.01	25.9	26.2	25.30
4	14.8	14.0	14.5	13.90

Table 34: [UrbanAir] KPI Time-to-solution (strong)

KPI	Iterations / s			
GPUs	Dibona	IDVA-wo-2phase	IDVA-w-2phase	IDVA-w-2phase-opt
1	9.4	10.1	10.2	10.64
2	17.85	19.3	19.08	19.76
4	33.78	35.71	34.48	35.97

Table 35: [UrbanAir] KPI Iterations / s (strong)

KPI	Iterations / kW (used GPUs)			
GPUs	Dibona	IDVA-wo-2phase	IDVA-w-2phase	IDVA-w-2phase-opt
1	31.33	29.71	30.01	37.72
2	31.88	29.25	29.18	35.29
4	31.28	30.79	29.73	33.31

Table 36: [UrbanAir] KPI Iterations / kW (strong, used GPUs)

KPI	Iterations / kW (all GPUs)			
GPUs	Dibona	IDVA-wo-2phase	IDVA-w-2phase	IDVA-w-2phase-opt
1	18.43	18.36	18.55	21.62
2	25.50	24.13	24.03	28.23
4	31.28	30.79	29.73	33.31

Table 37: [UrbanAir] KPI iterations /kW (all GPUs, strong)

KPI	Time-to-solution [s]		
	GPUs	IDVA-2phase	IDVA-2phase-opt
1	49	47	
2	50	47.8	
4	50.1	47.79	

Table 38: [UrbanAir] KPI time-to-solution (weak)

KPI	Iterations / s		
	GPUs	IDVA-2phase	IDVA-2phase-opt
1	10.2	10.63	
2	10	10.46	
4	9.98	10.46	

Table 39: [UrbanAir] KPI Iterations / s (weak)

KPI	Iterations / kW		
	GPUs	IDVA-2phase	IDVA-2phase-opt
1	30.01	37.72	
2	14.92	17.85	
4	7.34	8.96	

Table 40: [UrbanAir] KPI Iterations / kW (used GPUs, weak)

KPI	Iterations / kW (all GPUs)		
	GPUs	IDVA-2phase	IDVA-2phase-opt
1	18.55	21.62	
2	12.34	14.4	
4	7.34	8.96	

Table 41: [UrbanAir] KPI Iterations / kW (all GPUs, weak)

3.11.5 Mixed precision

Running UrbanAir to assess air quality in the city requires computational resources and time. To run a 24-h forecast over a city district with 5M grid points, hundreds of CPUs are required with time-to-solution close to six hours. These requirements may be hard to fulfill by individual users or environmental institutions willing to receive prediction in less time and using less resources. Adaptation to GPU accelerators is one of the possible approaches. However, further improvements in time-to-solutions may be achieved by integrating GPU accelerators with approximate computing techniques. This can be done by exploiting TAFFO.

TAFFO is a precision tuning framework based on LLVM, as described in D4.3 and D4.7. While originally developed for CPU-based applications in embedded systems, TAFFO has been extended to GPU accelerators within the TEXTAROSSA project, supporting both OpenCL and CUDA applications.

In order to study the feasibility of using TAFFO for compiling a reduced-precision version of UrbanAir, we started with a simple, yet time consuming kernel: `rhsdiv`. This kernel is written in C++ using templates for stencil computations, using shared memory parallelization on a single node (OpenMP) and MPI for multimode execution (CPU or GPUs).

In order to be able to use TAFFO, several modifications to the code were required, among others to be able to compile using LLVM for both, CPUs and GPUs. Next, improvements in TAFFO were required to support code generation for heavy-templated code. However, at the present time TAFFO still cannot support CUDA compilation flows where the kernel and host code are both present in the same file, and UrbanAir relies on this functionality of the NVidia `nvcc` compiler. Hence, the measurements on GPUs were performed by unbundling the kernel from the rest of the application.

In the experiments, TAFFO was configured for two separate numeric representations: 32 bit fixed point (i.e. scaled integers) and 16 bit floating point. The GPU employed was an NVidia GeForce RTX 4070 with 12GB of RAM. The average absolute error in the results of the `rhsdiv` kernel was $4.4e-8$ for 32 bit fixed point, and $7.1e-5$ for 16 bit floating point, both compared against 32 bit floats - the default data type employed by the UrbanAir application. By exploiting 16 bit floating-point types we obtained a speedup of 1.25 times over the original 32 bit floating-point code, and no speedup from exploiting the fixed-point representation. This is consistent with the experimental results provided in D4.3 and demonstrate the potential benefits of using TAFFO in this context.

More work is required to provide a complete integration, particularly on the TAFFO side. The reliance of TAFFO on the LLVM-IR technology limits its applicability for applications written with language abstractions which disappear at the compiler-intermediate-representation level. To resolve this problem, a different compiler toolchain needs to be adopted. In the current state-of-the-art, the MLIR (Multi-Level IR) approach is favored as a platform on which to build next-generation compilers, and for TAFFO it would provide the ability to transparently analyze and tune applications that exploit complex language features without specific support for them. On the other hand, from the UrbanAir side, the main problems surfaced by the integration work stem from the use of a closed-source GPU programming API, CUDA. The closed-source nature of CUDA makes it difficult to adapt applications in order to exploit open solutions such as TAFFO and LLVM. Recently, the Kronos group has developed an alternative to CUDA which solves this issue: SYCL, a modern open API for GPU-accelerated applications. The adoption of SYCL would make UrbanAir easily scalable to novel compiler solutions such as TAFFO, or other kinds of accelerators such as FPGAs.

3.12 FIPLib: FPGA Image Processing Library – ENEA/INFN

3.12.1 FIPLib on single FPGA

In this section we report the results achieved implementing the FPGA Image Processing Library (FIPLib) through the Vitis HLS flow. FIPLib is implemented as a collection of C++ kernels and utilizes the Vitis HLS flow to translate the C++ library's kernels, interconnected through streams, into bitstreams. Currently, FIPLib encompasses nearly 70 functionalities, each designed with a streaming behavior: data are read from the input streams and are written to the output streams. The width of all streams is controlled by a template parameter that determines both the data path width and the available parallelism. The functionalities within FIPLib can be classified as follows:

1. **Stream Management:** Encompasses stream copy, stream split, stream merge operations. These modules can also adapt the stream protocol (internal stream or AXI stream) and the stream width (increase/reduce stream width by factor 2 and 4)
2. **Data Mover:** Serves as the interface between DDR/HBM and streams.
3. **Parameter Initialization:** Establishes weight values for distinct filters like Sobel, Gaussian, and Box filters and for structuring elements for morphological operators.
4. **Pixel-based Operations:** Orchestrates pixel-based operations on input images, yielding output images through operations such as multiplication, subtraction, addition, complementation, logical operations, minimum, maximum, and linear scaling.
5. **Color Space Conversion:** Undertakes pixel-based transformations, converting images from the RGB color space to the YUV color space and vice versa.
6. **Image Masking:** Executes pixel-based operations that derive output images based on pixel values from either of two input images, based on the binary masking image.
7. **Line-based Operations:** Requires complete reading of one or more lines prior to initiating the output of processed images. This category encompasses diverse filters (2D convolution with kernels of dimensions 3x3, 5x5, 7x7), median filtering (3x3, 5x5), morphological operators like dilation and erosion, horizontal mirroring, and 2x up-sampling through bilinear interpolation.
8. **Store and Forward Transformations:** Demands full image read before generating the output image. This includes operations such as histogram equalization, contrast maximization, and vertical mirroring.

All functions share a consistent structure, and they encompass the subsequent parameters:

- Input and Output streams
- Function specific parameters (such as kernel weights in convolutional filters)
- Image size (including the number of rows and columns)
- Number of images to be processed before the function restarts.

These functions, or tasks, can be interconnected via streams and are invoked within a designated DataFlow section. Tasks within a DataFlow section are executed concurrently, with synchronization achieved through streaming communications. In Figure 101 we provide an excerpt of a basic image-processing kernel, accompanied by its graphical representation within a process network:

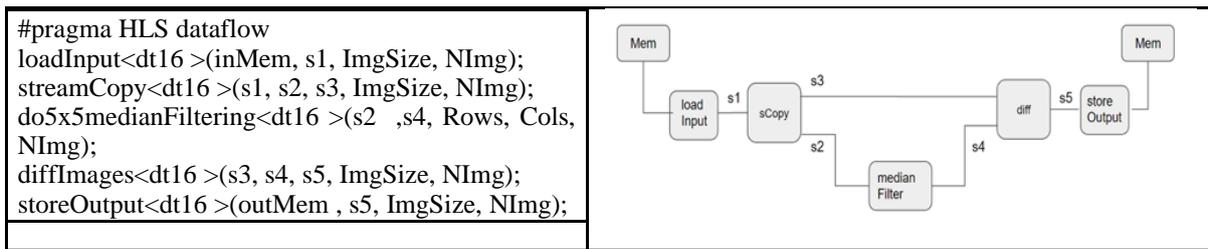


Figure 101:[FIPLib] Codelet and graphics showing a simple image processing kernel

The key point in achieving optimal FPGA processing performance is the maximization of parallelism, both in data access and data processing. Within FIPLib, four distinct types of parallelism are employed:

- **Spatial Parallelism:** This form of parallelism involves the concurrent processing of N image components, optimizing the utilization of parallel resources.
- **Fine-Grain Pipeline Parallelism:** Employed in tasks such as stream management, data movers, and pixel-based operations, this type of parallelism operates at the granularity of an individual pixel component.
- **Medium-Grain Pipeline Parallelism:** Implemented in line-based operations, this type of parallelism is centered around processing an entire line of the image, hence the unit of parallelism is one line of the image
- **Coarse-Grain Pipeline Parallelism:** Observed in store and forward transformations, this parallelism form encompasses the entirety of an image as the unit of parallel processing. It is applicable only when processing a sequence of images.

To demonstrate FIPLib capabilities, let’s refer to an algorithm, ALGO1, which draws a pencil-like sketch of a given input color image. ALGO1 uses the following set of kernels, *ImgProc*, as building block, see Figure 102.

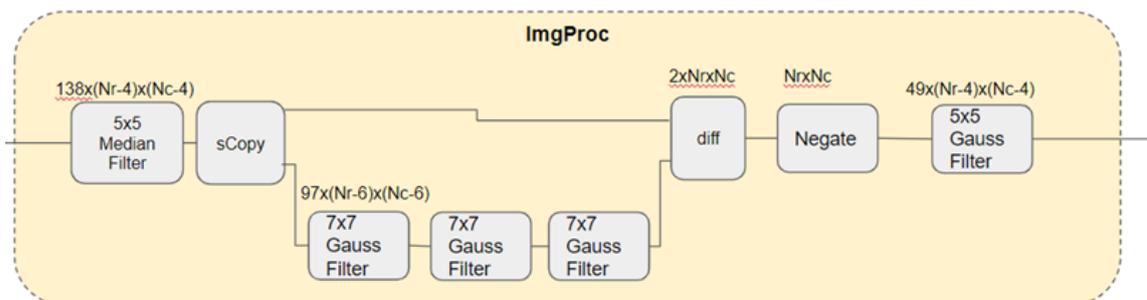


Figure 102: [FIPLib] Basic image processing pipeline used in ALGO1

The *ImgProc* set of kernels processes images (Rows, Cols) with 8 bits per pixel (8bpp, 256 levels of grey). It cleans the input image through a 5x5 median filter and uses a cascade of three 7x7 gaussian filters to remove the high frequency components of the image; the difference with the cleaned input image allows to extract the high-frequency components (smoothed borders) which are reported as a drawing on a whit paper (Negate kernel) and further smoothed through a final 5x5 Gaussian filter. Thanks to the streaming behavior of the kernels, at the steady state *ImgProc* can process N pixel components of the input image per clock cycle, being N the width, in bytes, of the streams.

The basic pipeline is instantiated three times to process the R, G, and B channels of a color image. This image is read/written from/to the DDR memory through data movers kernels, with the stream width being twice the width of the processing kernels.

The input image is separated/merged using the splitRGB and mergeRGB kernels. The structure of ALGO1 is reported in Figure 103.

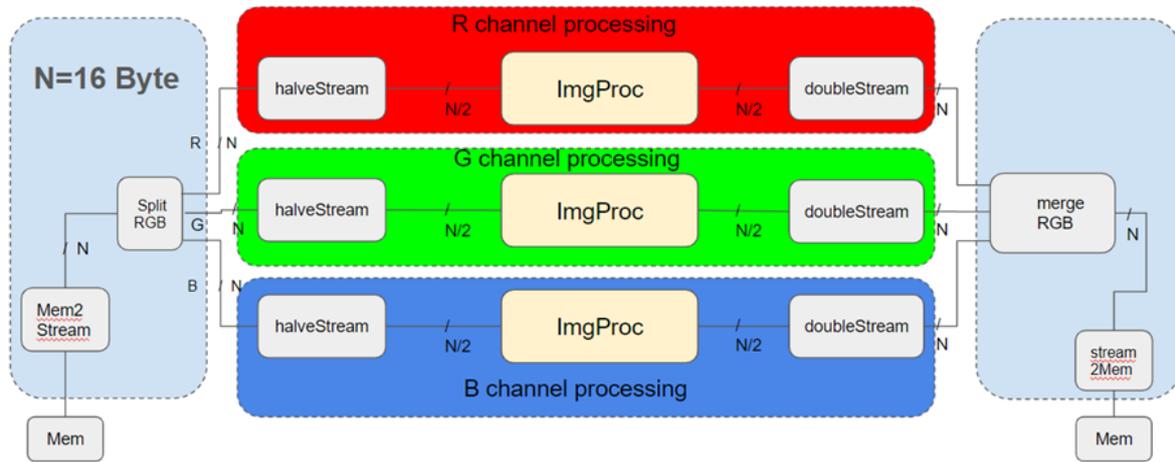


Figure 103: [FIPLib] Kernels in ALGO1

As observed, the size of the streams of the split/merge RGB kernels is twice the size of the in/out streams of the Image Processing Pipelined kernels. This is because data arriving at the R, G, and B channels has 1/3 of the throughput of the memory channels. As a result, new data is presented at the R (or G, or B) channel processing pipeline every 3 clock cycles. Therefore, we halved the size of the streams (and consequently the parallelism) by a factor of two. This approach allows the inner pipeline to process data with an efficiency of 2/3 instead of 1/3. Counting the number of operations performed by ALGO1 to process a (Rows; Cols) RGB image (refer to the expressions reported in the ImgProc figure), we find that the number of operations to compute ALGO1 is given by

$$N_{ALGO1} \approx 1.4 \times 10^3 \times \text{Rows} \times \text{Cols}$$

being operations seized as either 8-bit integer or 16-bit integer (multiply in the Gaussian Filter) or 32-bit (addition in the Gaussian Filter).

After synthesizing the design produced by the HLS flow, we ran ALGO1 to process a sequence of 3000 4096x4096 RGB images. We measured the time from the start of the send of the first image to the FPGA memory to the end of the reception of the last image processed by the FPGA. The time to process all the images is 65.8 s, so one image is processed in $T_{\text{ExeFPGA}} = 21.9$ ms and the sustained speed is

$$S = N_{ALGO1} / T_{\text{ExeFPGA}} = 2.35 \times 10^{10} / 2.19 \times 10^{-2} = 1.07 \times 10^{12} \text{ Op/s}$$

The resources used by the ALGO1 are reported in Table 42.

	#	%
LUT	121127	10.9%
BRAM	328	19.4%

DSP	4376	48.6%
-----	------	-------

Table 42: [FIPLib] Resources used by ALGO1

ALGO1 clock frequency is $f_{ck}=200$ MHz. As the number of clock cycles to process one image is $N_{ck}=T_{Exe} \times f_{ck}=4.3 \times 10^6$, ALGO1 is sustaining

$$N_{ALGO1} / N_{ck} = 5.5 \times 10^3 \text{ Ops/cycle}$$

To evaluate the advantages of FIPLib, we implemented ALGO1 on a multicore Intel Xeon Haswell CPU using the OpenCV library [see <https://opencv.org>]. We left to the OpenCV the management of the parallelism and we saw that 9 cores were fully used during the processing. We processed the same sequence of 3000 RGB images used in the FPGA tests. The processing time for the whole sequence is 508.5 s, so the time needed to process one image is $T_{ExeopenCV}=170$ ms and the speed-up achieved through the FPGA implementation w.r.t. openCV is $S = 7.8$

FPGAs are commonly regarded as computing devices capable of reducing power consumption. Let's assess and quantify this aspect through measurements on FPGA kernels implemented using the FIPLib and compare their energy usage with the corresponding OpenCV implementation on the CPU.

To quantify the power consumption during a computation, let's consider the energy used by the application to be run, being the energy defined as

$$E_A = \int_0^D P_A(t) dt$$

where D is the running time of the application, and

$$P_A(t) = P(t) - P_{Idle}$$

the power used by the application at time t , computed as difference between the node instant power and the power absorbed by the node when no user processing is running.

To perform the power measurements, we accessed the Baseboard Management Controllers (BMCs) of the computing node through the Intelligent Platform Management Interface (IPMI): in this way, with a sampling period of 1 s, we can read the instant power absorption of the node. As $P(t)$ is the power erogated by power supply unit, it includes the power used both by the CPU and by the FPGA.

To compute E_A , we first estimate P_{Idle} measuring the instant power absorbed when there are no user processes running, then we take the instant power measurements, along the application run, to compute E_A .

To express the global performance of the execution of an algorithm implementation on a platform, we use the Energy Delay Product (EDP) [12].

$$EDP = E_A \times D$$

EDP [Js] considers both the energy consumed by the application and its execution time. Using EDP, we can distinguish between two designs that consume the same amount of energy by identifying the faster one. Smaller EDP indicates more efficient design implementation.

Figure 104 reports $P_A(t)$ when running ALGO1 to process the 3000 RGB images on CPU, through OpenCV, and on FPGA, through the FIPLib.

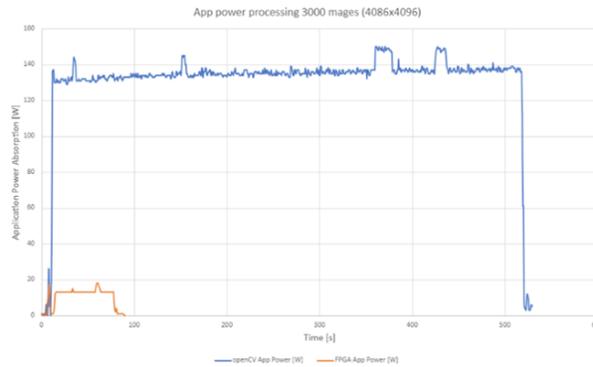


Figure 104: [FIPLib] Power absorption of ALGO1

The values for E_A (i.e. the areas below the two plots), the computing times, and EDP are reported in Table 43 (computing time includes FPGA programming time, filesystem accesses, and all the needed initializations).

	Energy used by ALGO1 run (kJ)	Time used by ALGO1 run (s)	EDP for ALGO1 run [kJ*s]
FPGA / FIPLib	0.92	91	60
CPU / openCV	69.22	530.6	35198

Table 43: [FIPLib] Energy and time used for ALGO1

FPGA implementations outperform CPU implementations in both energy efficiency and speed, using 75 times less energy and 7.8 times less time.

3.12.2 Multi-FPGA Implementation

One of the possible improvements for the application is to perform image processing by working on streams of pixel of 256 bit width. However, this possible update translates in a huge requirement in terms of resources to implement the RGB kernels described in the previous section (in particular, since each `ImgProc` hw function requires ~32% of the available DSPs of the board). This made it impossible to implement the total application on a single FPGA. To cope with this problem, we chose to deploy FIPLib on a multi-FPGA setup exploiting the APEIRON framework: in this way, we were able to split the overall image processing by implementing a single RGB kernel on each node, each of them dedicated to a single-color stream processing. Implementing FIPLib HLS kernels as APEIRON tasks means changing the interface of each of them (to cope with the standard required by the framework to compile the entire project. An example of this type of change is reported in Listing 1.

<u>SINGLE FPGA FPLib IMPLEMENTATION</u>	<u>MULTI-FPGA FPLib IMPLEMENTATION (APEIRON)</u>
<pre>extern "C" { void ImgProc(hls::stream<io_stream_16B> &s_in, hls::stream<io_stream_16B> &s_out, unsigned int ImgSize, unsigned int NbImages, unsigned short int ImgRows, unsigned short int ImgCols, unsigned int channel) { ... } }</pre>	<pre>#include "ape_hls/hapecom.hpp" extern "C" { void ImgProc(message_stream_t message_data_in[N_INPUT_CHANNELS], message_stream_t message_data_out[N_OUTPUT_CHANNELS], unsigned int ImgSize, unsigned int NbImages, unsigned short int ImgRows, unsigned short int ImgCols, unsigned int channel_id) { ... } }</pre>

Listing 1: [FIPLib] Single vs multi FPGA

Another requirement for the implementation via APEIRON was to substitute the basic Vitis HLS stream connection approach (stream.write, stream.read) with the HAPECOM APIs send() and receive() functions in order to allow the communication between them by exchanging Apelink packets through the network. A snippet of the code is reported in Listing 2.

<u>SINGLE FPGA FPLib IMPLEMENTATION</u>	<u>MULTI-FPGA FPLib IMPLEMENTATION (APEIRON)</u>
<pre>while (NbWordToTransfer > BUFFER_SIZE) { if (phase){ buffer2Stream(outStream, Buff1, BUFFER_SIZE); stream2Buffer(inStream, Buff2, BUFFER_SIZE); } else{ buffer2Stream(outStream, Buff2, BUFFER_SIZE); stream2Buffer(inStream, Buff1, BUFFER_SIZE); } phase = !phase; NbWordToTransfer -= BUFFER_SIZE; } void buffer2Stream(hls::stream<io_stream_16B>& outStream, dt16 Buff[BUFFER_SIZE], unsigned int size) { #pragma HLS inline off io_stream_16B tmp; tmp.keep = 0xFFFF; tmp.last = false; // copy Buff to stream for (unsigned int i = 0; i<size; i++){ #pragma HLS pipeline tmp.data = Buff[i]; outStream.write(tmp); } }</pre>	<pre>#include "ape_hls/hapecom.hpp" while (NbWordToTransfer > BUFFER_SIZE) { if (phase){ send(Buff1, BUFFER_SIZE*sizeof(word_t), coord, task_id, ch_id, message_data_out); stream2Buffer(inStream, Buff2, BUFFER_SIZE); } else{ send(Buff2, BUFFER_SIZE*sizeof(word_t), coord, task_id, ch_id, message_data_out); stream2Buffer(inStream, Buff1, BUFFER_SIZE); } phase = !phase; NbWordToTransfer -= BUFFER_SIZE; }</pre>

Listing 2: [FIPLib] Single vs multi FPGA implementation

The system was composed by 4 interconnected Xilinx® Alveo U200 (installed in the INFN Roma1 APE Lab) in a ring topology. The bitstream flashed on each board implemented a project depicted in Figure 105 and composed by 3 different HLS kernels:

- *RGB2Mem*: connected to INFN Communication IP intranode port 0, it is built combining the overcited *mem2stream* and *splitRGB* kernels
- *Mem2RGB*: connected to INFN Communication IP intranode port 0, it is built combining the overcited *stream2mem* and *mergeRGB* kernels
- *ImgProc*: connected to INFN Communication IP intranode port 1, its functionality is described in the previous setup.

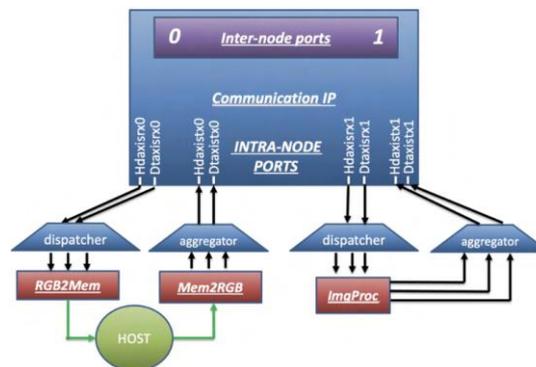


Figure 105: [FIPLib] APEIRON project for single FPGA bistream

In the final setup to be deployed on the four-FPGA system, each board (and so each hardware project flashed) is enabled to work on a specific task. In fact, the FPGA 0 is used to load the RGB images from the memory, to split them into color streams of data and to send them through the network in which the FPGA 1, 2 and 3 are used to process respectively a single-color stream of data and to send back results to the FPGA 0 memory.

The FIPLib multi-FPGA deployment and execution schemes are depicted in Figure 106, while a detail on which kind of processing HLS kernels (previously described) are used on each board is represented in Figure 107.

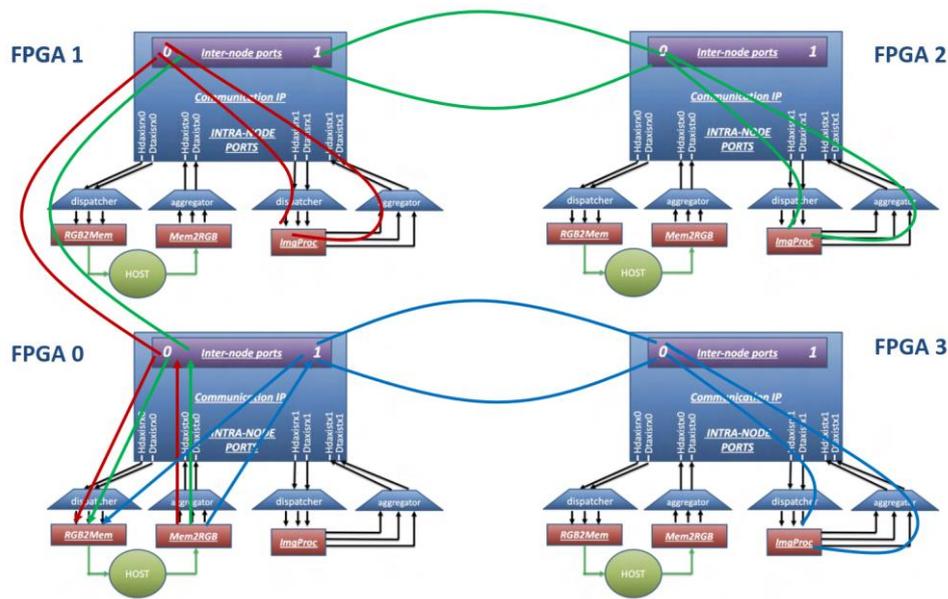


Figure 106: [FIPLib] multi-FPGA deployment and execution scheme

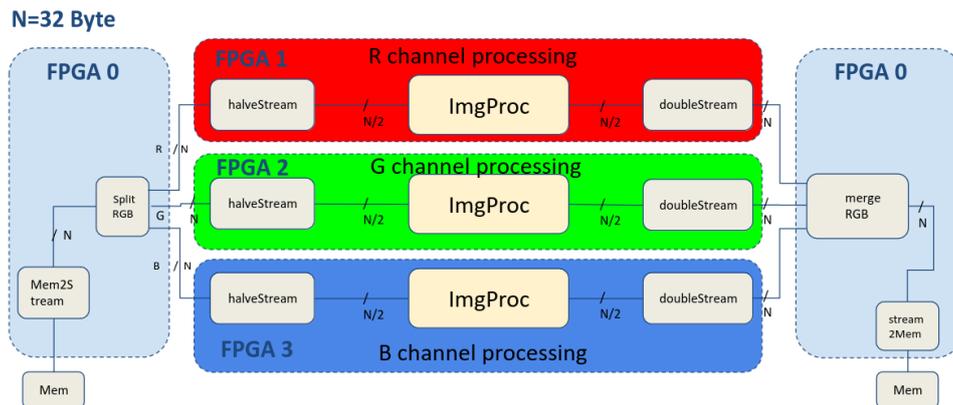


Figure 107: [FIPLib] Kernel displacement in multiple FPGA setup. Internal datapath set at 32 bytes

To test the application performances, we choose FIPLib on images of difference sizes, measuring the execution time and the power consumption in each of the cases.

In detail, FIPLib execution time and power consumption has been measured exploiting the XRT runtime library, which allows the user to measure FPGA internal voltages and currents values during the application run and to insert time checkpoints to analyze the timeline of running HLS kernels.

The performances result for processing images of sizes 512x512 and 4096x406 are respectively depicted in Figure 108 and Figure 109 and reported in Table 44.

Image Processing APEIRON (256bit datapath@100 MHz) [9000 images 512x512]

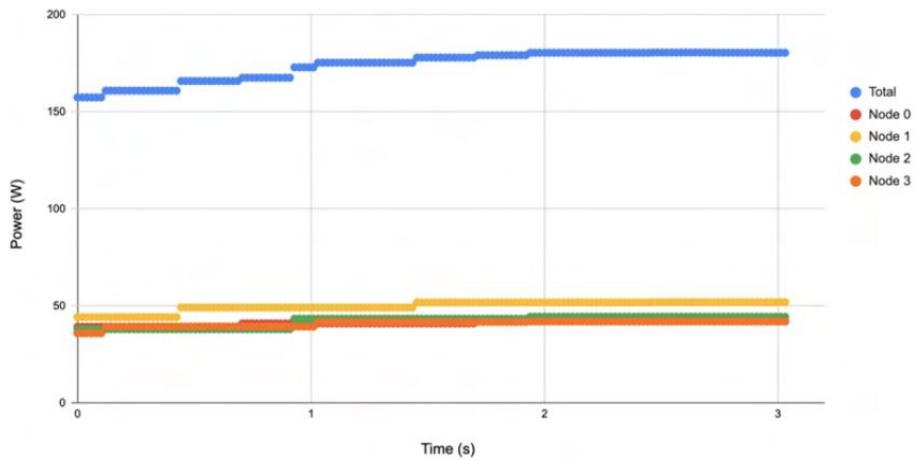


Figure 108: [FPLIB] Multi FPGA APEIRON Image Processing on 512x512 images execution power profile

Image Processing APEIRON (256bit datapath@100 MHz) [200 images 4096x4096]

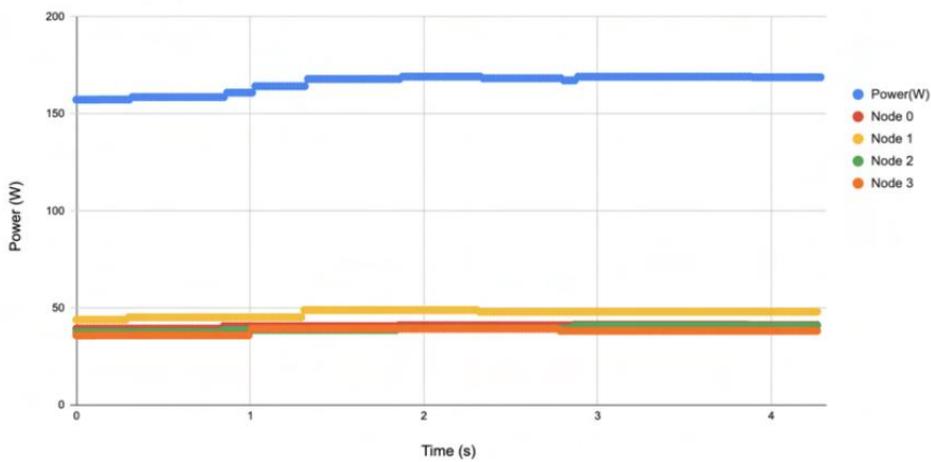


Figure 109: [FPLIB] Multi FPGA APEIRON Image Processing on 4096x4096 images execution power profile

	Image Size	
	512x512	4096x4096
Processing Time	3.05 s	4.28 s
Throughput	2950.62 fps	46.70 fps
Energy consumption	0.53 kJ	0.71 kJ
Processed Images per Joule	16.98 fpJ	0.28 fpJ

Table 44: [FIPLib] Multi-FPGA APEIRON Implementation performances

3.12.3 Single- and Multi-FPGA implementation comparison



Figure 110: [FIPLib] CPU, Single- and Multi-FPGA implementation performances comparison

In Figure 110, performances obtained from the previous described FIPLib implementations (CPU, Single-FPGA, Multi-FPGA) are compared focusing on throughput and processing images per joule values obtained from the processing of different sizes of images.

In terms of throughput, it can be seen a large improvement while working on 512x512 images on a multi-FPGA setup: this is coherent since the possibility of increasing the system datapath to 32 byte wrt to the limit of 16 bytes in the single-FPGA system. A slightly improvement can be seen also while working on 4096x4096 images: in this case the throughput is limited from the bandwidth of the INFN Communication IP used in the APERION framework (since the size of packets to be sent on the network scales with the numbers of columns of the image).

In terms of energy, the scaling in number of used FPGA boards is coupled to a huge loss in processing images per joule (higher loss in energy when lower improvement in throughput).

FPLIB kernels have also been used to demonstrate the possibilities offered by the FastFlow FPGA library extension developed in WP4 and discussed in deliverable D4.7 and whose final results have been outlined in D4.9. As an example, the kernel(s) outlined in Figure 103 have been integrated in a FastFlow pipeline. The FastFlow code executing the kernel(s) on a U50 FPGA board achieved performances comparable to the ones achieved running the very same kernel(s) using OpenCL only host code, thus showing there is no additional overhead introduced by the FastFlow kernel execution orchestration. Further experiments have been run on the IDV_E node demonstrated that the FastFlow FPGA library perfectly manages the execution of kernels on the available U280 FPGAs and achieves results comparable to the ones achieved on the U50 as well as to the ones achieved on IDV_E node using OpenCL only host code to manage the FPGA kernel executions.

3.13 NBody - BSC

This section shows the results of the N-Body simulation application. This application allows evaluation and comparison of different use cases of the IDV-E computing node. It also provides a way to verify the successfulness of the different software stack integrations developed in the project.

This section presents results when running the application in a single FPGA, multiple FPGAs without FPGA-to-FPGA communication, multiple FPGA with direct FPGA-to-FPGA communication using APEIRON communication framework and finally, multi-node multi-FPGA using OMPIF over ethernet.

N-body simulation computes the interaction of a set of particles due to gravitational forces. The input is a set of particles with initial positions, velocities, and masses. Both positions and velocities are represented by 3-dimensional vectors. Data is represented with single precision. During the simulation, two steps are repeated iteratively. The first one is the computation and accumulation of forces that each particle exerts over all other ones. This part is the most computationally expensive, because the number of forces grows proportional to n^2 where n is the number of particles. The second one is the update of the particle velocity and position for a given time step. Therefore, most of the resources are dedicated to the force computation step.

3.13.1 OmpSs@FPGA on IDV-E single FPGA

This implementation uses one of the two available FPGA devices, it's used as a baseline and by a means to evaluate thermal and power profile of the IDV-E cooling system described in deliverable D3.2.

The design is created using the OmpSs@FPGA toolchain [13] which is described in deliverable D4.6. Figure 111 shows a high-level diagram of the resulting design.

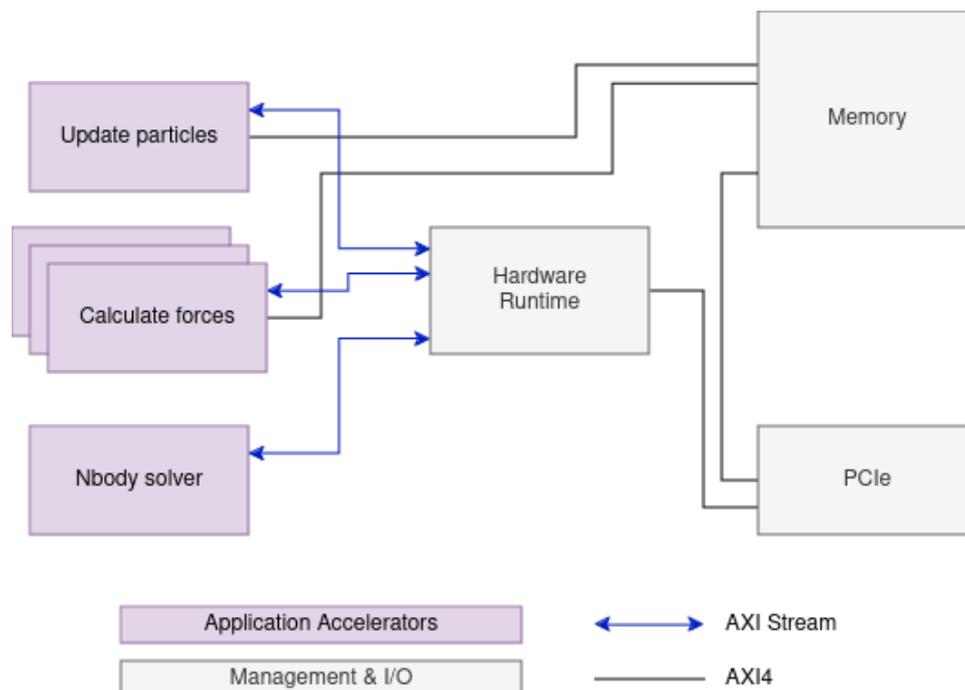


Figure 111: [NBody] Single node n-body FPGA design

Each block represents an IP core or module, and lines represent data buses. Purple boxes represent each of the application accelerators in the design. The stacked boxes represent multiple instances of the same accelerator. Gray boxes represent management modules and I/O infrastructure.

The n-body design implemented contains one instance of the *Nbody solver* and *update particles* kernels and 8 instances of the *calculate forces* accelerators. Each of the *calculate forces* kernels calculates forces for 16 particles in parallel. This distribution of resources is done because force calculation is the most compute intensive part of the application.

Table 45 shows the amount of resources used by this design.

Resource	Used	Available	Used %
LUT	600625	1303680	46.071507
LUTRAM	99094	600960	16.489285
FF	800548	2607360	30.703392
BRAM	614	2016	30.456348
URAM	8	960	0.8333334
DSP	5131	9024	56.859486
IO	6	624	0.9615385
GT	16	24	66.66667
BUFG	12	1008	1.1904762
MMCM	1	12	8.333334
PCie	1	6	16.666668

Table 45: [NBody] Resource usage of the single-node n-body for the IDV-E alveo U280

In this implementation, the host copies data to FPGA memory and then submits a single FPGA task that will solve the problem for a given number of particles and timesteps. This is shown in Figure 112.

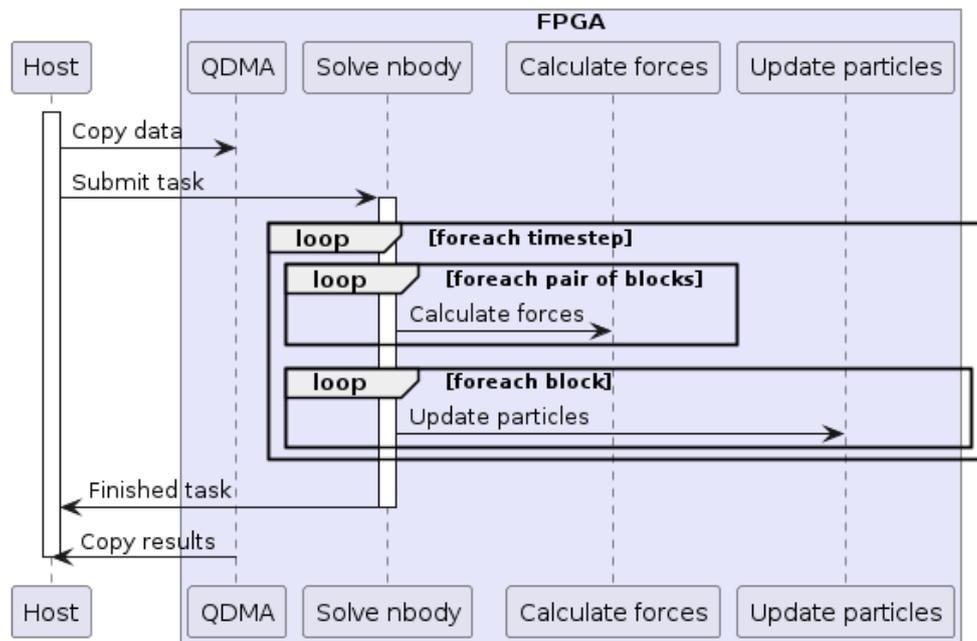


Figure 112: [NBody] Diagram sequence of the n-body computation for single FPGA

The host first copies data to the FPGA memory by submitting a DMA operation to the QDMA module, which implements the PCIe communication between the host and the FPGA. Then, a single *solve nbody* task is submitted. This task then communicates with the hardware runtime to create and submit tasks to the *calculate forces* and *update particles* accelerators in order to solve the n-body problem for a given number of particles and time steps. Finally, the host copies the results back to main memory after the *solve nbody* task is finished.

Single FPGA implementation reaches a performance of 37.43 Gpps (Giga pairs per second) while consuming 94.84W. This yields a power efficiency of 0.395 Gpairs/Watt. More details as well as comparison with other platforms are described in deliverables D1.4 and D4.6.

Regarding thermal performance, IDV-E shows better behavior when comparing to other platforms. We compared against actively air-cooled Alveo U200, a passively cooled Alveo U55c and the U280 using in quattro two-phase cooling technology from IDV-E as described in deliverable D3.2.

Regarding memory configurations, Alveo U200 uses 64GB DDR memory, Alveo U55c integrates 16GB HBM memory instead of DDR and IDV-E Alveo U280 has 8GB HBM and 32GB DDR, however, in our tests, only HBM is used.

All three devices make use of the same design, using roughly the same amount of resources. Therefore, power draw from design logic is expected not to change across the devices as they are the same device family (AMD Virtex UltraScale+) using the same manufacturing process.

Figure 113, Figure 114, Figure 115 and Figure 116 show temperature and power usage for during an n-body execution of the same problem size in U200, U55, and IDV-E's U280 air cooled and 2 phase cooled devices respectively.

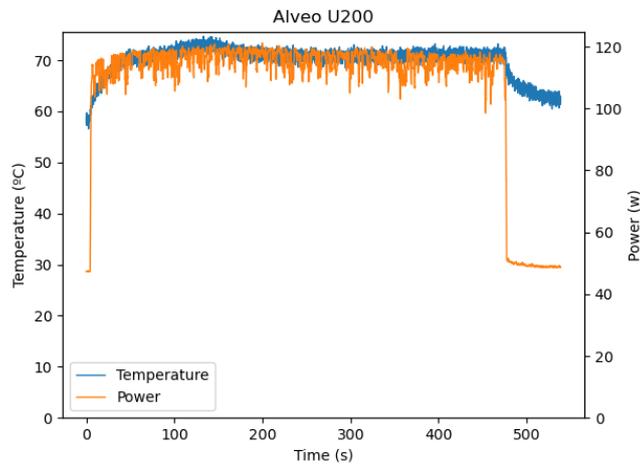


Figure 113: [NBody] Temperature and power for an actively air-cooled Alveo U200

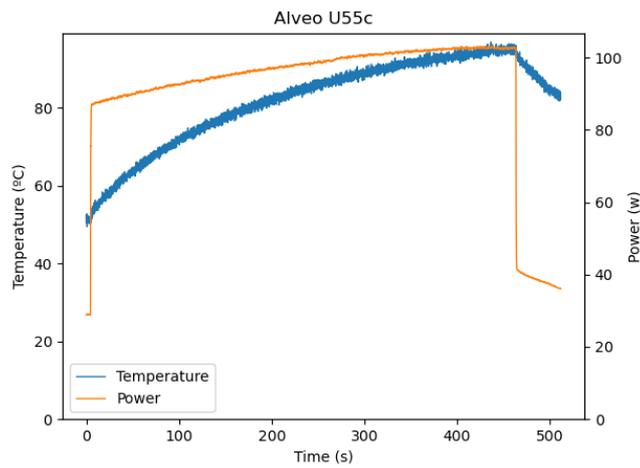


Figure 114: [NBody] Temperature and power for a passively air-cooled Alveo U55c

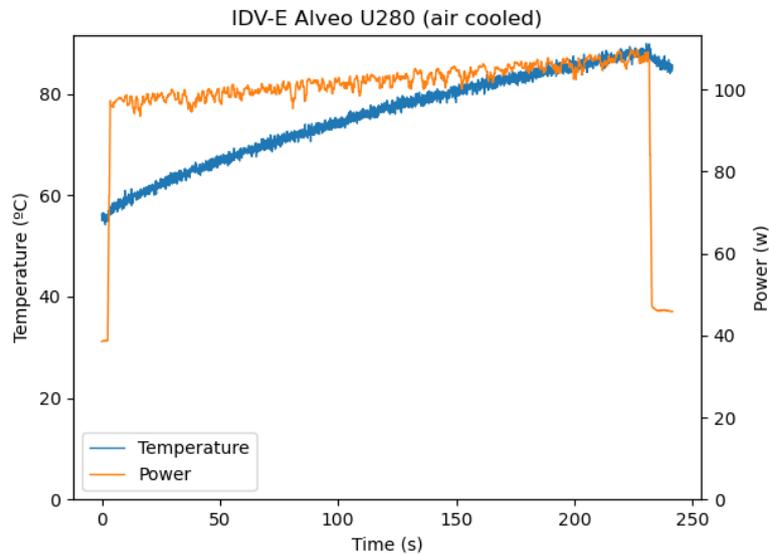


Figure 115: [NBody] Temperature and power air cooled IDV-E's Alveo U280

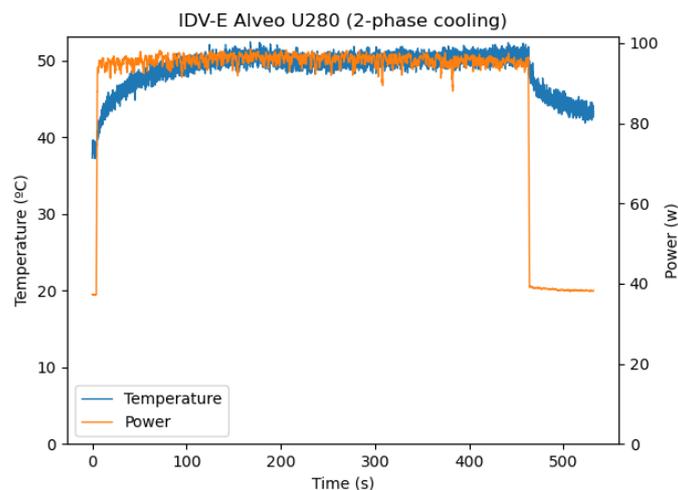


Figure 116: [NBody] Temperature and power 2-phase cooled IDV-E's Alveo U280

Those figures show many differences regarding temperature highlighting differences in the cooling solutions used in each case.

The temperature for the Alveo U200 seems to remain stable throughout the execution. This is the expected outcome as the fan controller should adjust fan speed depending on the thermal load. This can be seen in Figure 113 at around 150s into the execution, the temperature slightly drops as the fan speed increases. Also, idle temperature is higher than the other two platforms as the fan will run at less speed at lower thermal loads.

Alveo U55c shows lower idle temperature. This is due to the case airflow being relatively high disregarding FPGA's thermal load, which is the expected behavior in a passive cooling solution. However, when computation starts, the amount of dissipated power increases, but the airflow does not significantly increase, causing the temperature to increase throughout execution reaching 96°C in this execution. Insufficient airflow when thermal demand is high could lead to thermal protection shutting down the device to prevent catastrophic failure.

Finally, IDV-E cooling system shows bigger thermal mass as temperature increases at a lower rate than other devices at the start of the execution. A larger thermal mass allows to absorb peaks in thermal load. Temperature drops slightly at around 160 seconds; this suggests that the cooling system is reacting to the change in temperature as was the case for active cooler in the Alveo U200.

Those charts show also large differences in power usage across the evaluated devices and temperatures.

On Alveo U200, average power usage is 110W, which is higher than the 95W that the U280 uses. This difference can be attributed to the different memory used. This board uses 64GB DDR memory instead of the 16GB DDR + 8GB HBM of the U280. Also, active cooling has an impact, since the integrated fan has a non-negligible power draw that also is included in measurements.

On Alveo U55c, even though power draw at the start of the execution is 80W, this quickly increases due to the temperature change in the FPGA chip. This is a known effect in MOSFET-based integrated circuits [14]. At the end of the execution power draw is just over 102W. On longer executions, power seems to stabilize at around 103W

On the IDV-E, when using air cooling, behavior is similar to the U55c. However, cooling performance is worse. In fact, executions are shorter because on longer executions, overtemperature protection is tripped to avoid catastrophic failure. This can be seen as the x axis on Figure 115 shows a shorter span of time. This is caused by the airflow in the node being insufficient to properly cool the boards. On longer executions, the power draw reaches around 115W just before the device powers down itself. Reported power consumption reported is slightly lower as we run shorter executions in order not to get so close to the thermal limit.

In the 2-phase cooled IDV-E power usage remains stable throughout the execution, staying close to the 94W average. On top of running cooler overall, this cooling method allows us to run workloads for extended periods of time without triggering thermal protection.

All in all, power efficiency of the FPGA in 2 phase cooled IDV-E is highest of the devices evaluated, as shown on Table 46. Cooling solution plays an important role in power efficiency, as the device running hotter results in higher power draw. Even in the case of passive solutions the influence of temperature on power consumption can negate power savings of the active components. Furthermore, temperatures can reach dangerous levels in particular workloads if server chassis air flow is not adjusted accordingly.

Device	Performance (Gpps)	Power (W)	Energy efficiency (Gpps/W)
U200	36.4214	112	0.3265
U55c	37.4103	103	0.3629
U280 IDV-E (air)	37.4343	107	0.3487
U280 IDV-E (2-phase)	37.4348	94.8	0.3947

Table 46: [NBody] Comparison of performance, power and power efficiency across different devices

These power measurements only consider the FPGA device's power. In all cases except for the Alveo U200, cooling is passive and therefore power used by the cooling infrastructure is not included in measurements. Table 47 shows power consumption of the full node including

cooling infrastructure for the air-cooled and 2-phase-cooled versions of the IDV-E when using 1 and 2 FPGAs.

Power (W)	1 FPGA	2 FPGA
Air cooling	354.3846154	424.8186813
2 phase cooling	384.1043956	447.6978022

Table 47: [NBody] Comparison of power consumption on air cooled and 2-phase cooled IDV-E

As table shows, when considering the power consumed by the cooling infrastructure, it consumes more power than the air-cooled alternative for this configuration. However, the difference in consumed power when using one or two FPGAs is lower in the two-phase cooling system (+16.5%) than in the air-cooled version (+19.8%). Figure 117 shows a comparison between both versions of the IDV-E. For the cooling version, a breakdown of the used power between the computing node and the cooling system. For the air-cooled version, measuring power consumed by cooling fans is not possible as only full node power usage is reported.

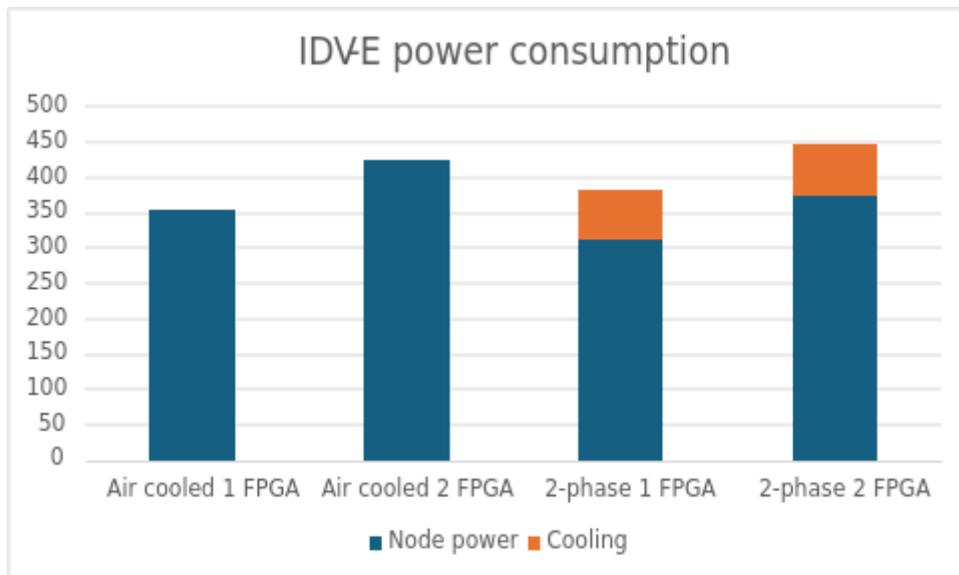


Figure 117: [NBody] Sequence diagram of multi-fpga n-body execution

Power consumption of the 2-phase cooling system remained constant at 73W throughout all experiments independently of the power used by the IDV-E node. Therefore, if the system used more FPGAs and the trend continues, a 2-phase cooled IDV-E node with more FPGAs would be more efficient than the air-cooled alternative. This trend is shown in Figure 118.

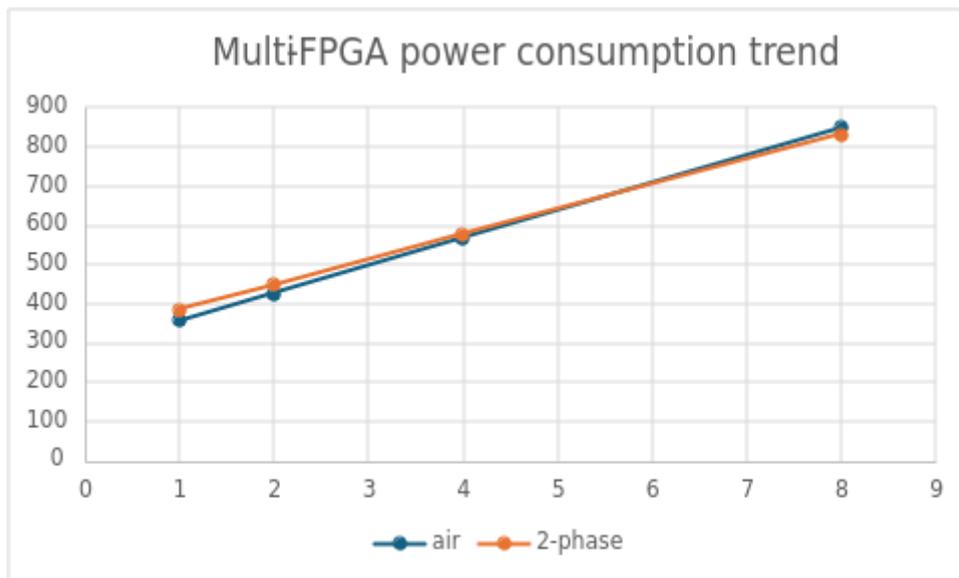


Figure 118: [NBody] Sequence diagram of multi-fpga n-body execution

This figure shows the power consumption trend for air-cooled and 2-phase-cooled IDV-E for different number of FPGAs.

For nodes with 8 FPGAs as the ones used in MEEP project [15], two-phase cooling would be more efficient than air-cooling the FPGAs. Also, in our experiments, CPUs are only used for FPGA task submission, no workload is being executed in the CPU, leaving almost all the cores in idle state. If workloads actively use CPUs, we expect 2-phase cooling to be even more efficient as cooling load increases.

It is also worth noting that the current air-cooled configuration is insufficient and cannot be used for workloads that run in the FPGAs for an extended time. In order to use the machine in production-ready environments more energy consumption would be expected from air-cooling, making an even better case for the use of 2-phase cooling.

3.13.2 OmpSs@FPGA on IDV-E multiple FPGA

In this use case, we use both FPGAs in the IDV-E to run the n-body simulation. In this case, the host processor manages execution across the FPGAs without relying in the FPGAs directly communicating between them.

To distribute work among multiple FPGAs, it's not possible to submit a single task that computes the full n-body problem. The host process must create tasks for each particle block and move data from one FPGA to another. This workflow is shown in Figure 119.

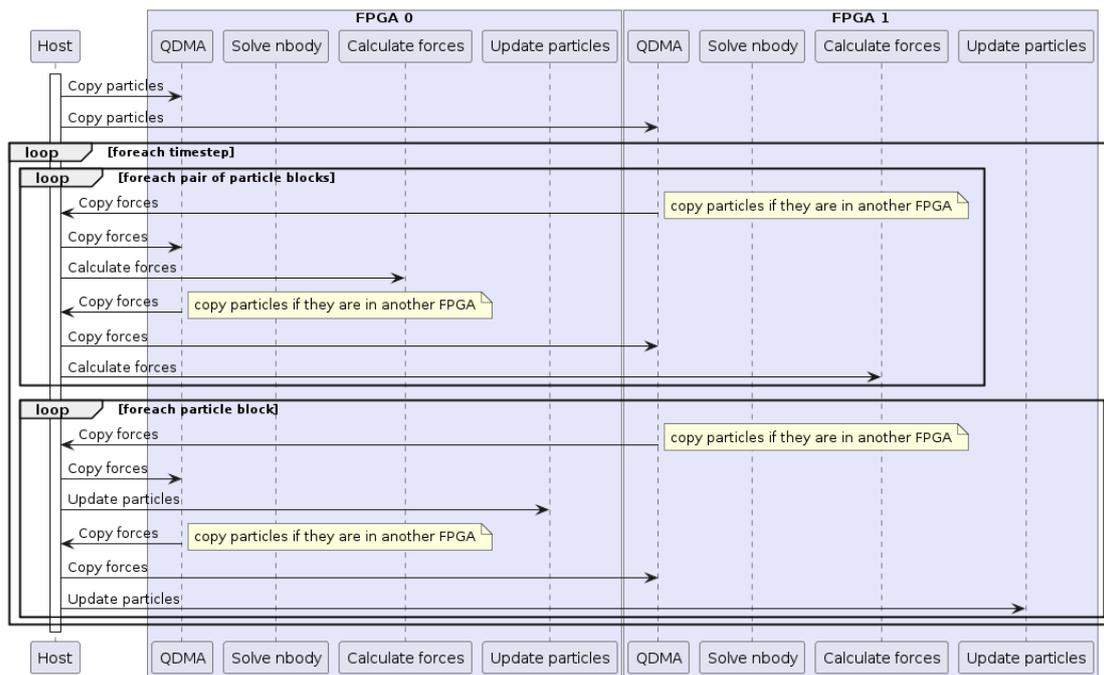


Figure 119: [NBody] Sequence diagram of multi-fpga n-body execution

The sequence diagram in Figure 119 clearly shows that the amount of communication between the host system and the FPGA is much larger than the single device implementation described in Figure 111. This is true for task execution commands since the host must send an individual command for each task to process every block of particles. Also, the amount of data copies quickly grows. Data that needs to be moved from one device to another needs to be copied from the source device memory to host main memory, and then copied again to the destination device memory. If P2P transfers were supported, data could be copied from one device to another without going through the host. This can decrease the total amount of data movements, but keeping track of where each individual block of particles is stored in a given point in time has a non-negligible overhead that will limit performance as the problem size and number of devices grow.

Another disadvantage of this approach is that task scheduling quickly becomes a bottleneck. The host may not be fast enough creating tasks, checking that their dependences are ready and sending them to the appropriate device to keep accelerators busy.

All in all, there's a performance degradation when moving to a 2-device implementation. This is shown in Figure 120.

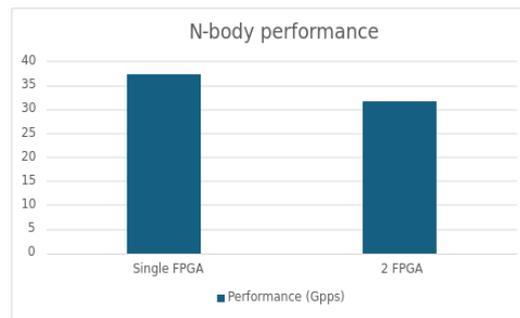


Figure 120: [NBody] Performance comparison between using 1 and 2 FPGA

This chart shows performance in Giga pairs per second (Gpps) for single and multi FPGA implementation. The design used for the multi-FPGA is the same as in the single node implementation shown in Figure 111, but the *solve nbody* accelerator is not used. Tasks are submitted directly to *calculate forces* and *update particles* accelerators.

Even though this approach may work for problems that can be partitioned in larger blocks that do not need communication between them, the all-to-all pattern of the n-body simulation makes it unfeasible to efficiently implement this application following this approach.

3.13.3 OmpSs@FPGA + APEIRON

To overcome issues regarding data copies and task scheduling of the host directly managed multiple devices, we implemented direct FPGA-to-FPGA data transfers. Communications between FPGA devices are implemented using APEIRON framework from INFN (described in D2.9). We use the apeiron API to send and receive data from the application accelerated kernels to the switch. The switch then routes packets to the destination FPGA in the network. Figure 121 shows a diagram representing the design of the of the n-body using apeiron to implement inter-FPGA communications.

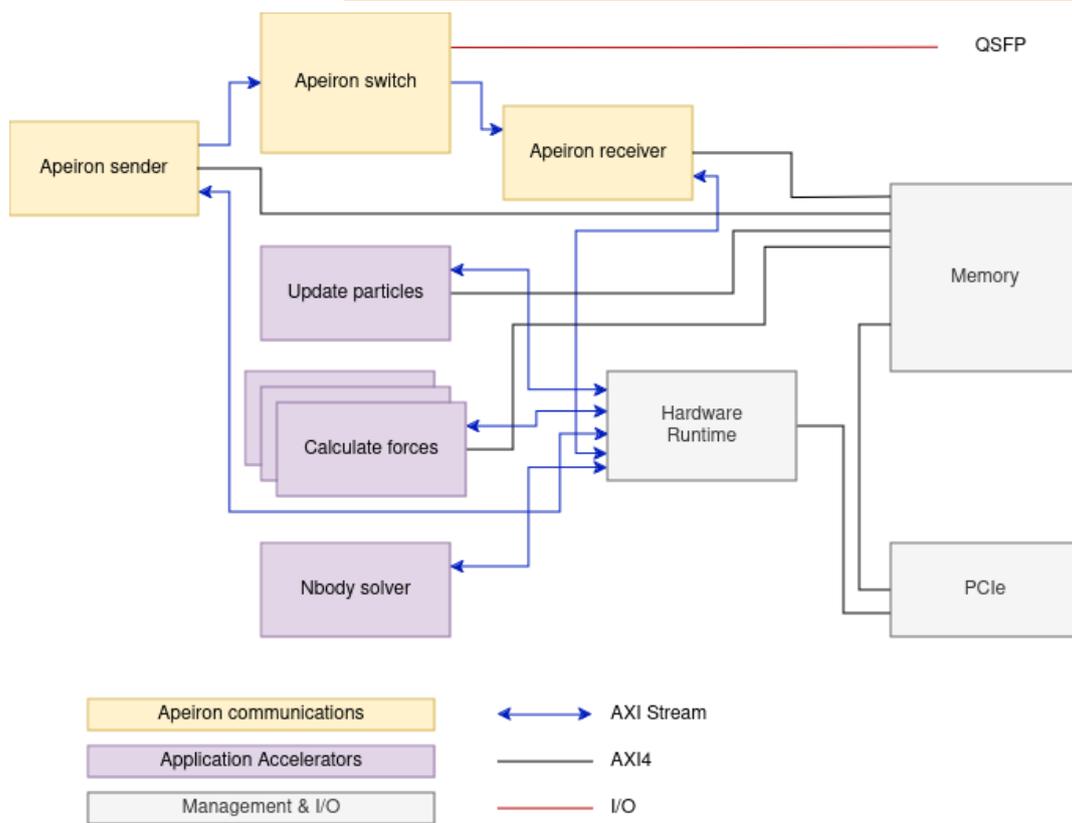


Figure 121: [NBody] body design diagram showing ompss accelerators and apeiron communication modules

Each of the blocks shown in the figure represents different IP cores in the design. The diagram shows application accelerators (purple blocks) connected to the hardware, which schedules tasks and to device memory. Hardware runtime sends commands to accelerators using an AXI stream interface. Accelerators read and write data to memory using AXI-4 interfaces. Yellow blocks are the apeiron communication modules. Apeiron sender and Apeiron receiver modules are connected like other application accelerators.

When an accelerator needs to send or receive data, a command is sent to the apeiron sender or receiver. Then it reads data from memory, builds apeiron packets and sends them to the switch through a stream interface. Then the switch sends the data through the appropriate QSFP port on the board. The receiver works in the opposite direction, receives data from the switch and writes it to memory. This process is shown along with all n-body execution flow in Figure 122.

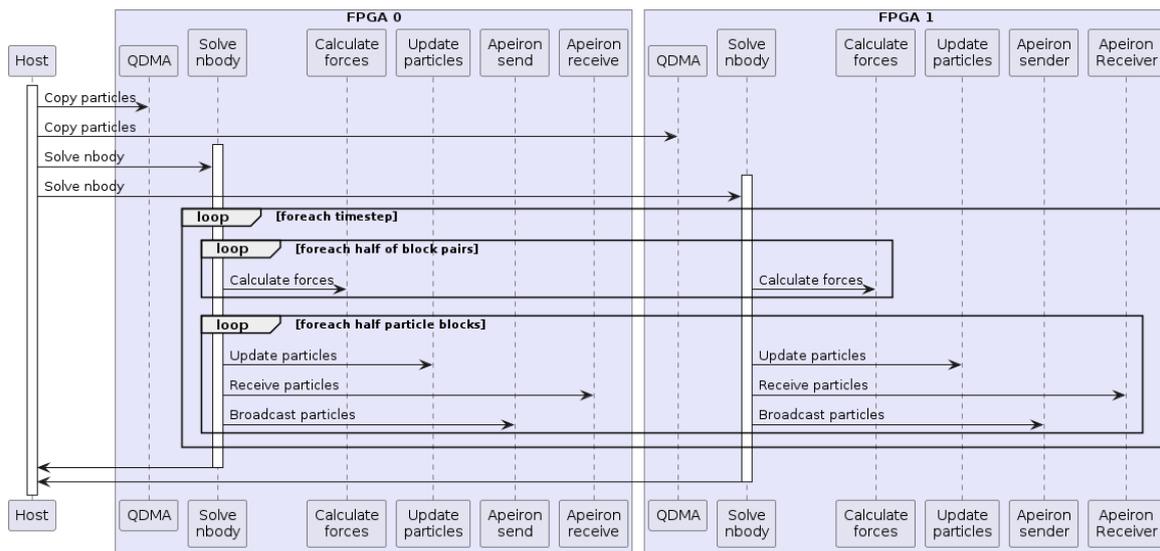


Figure 122: [NBody] Ompss + apeiron sequence diagram

This figure shows a sequence diagram showing a high-level view of the interaction of application accelerators and apeiron modules. This diagram shows the case of two FPGA devices connected to the host via PCIe using the QDMA module and between devices via QSFP using apeiron modules.

First, data is copied to all devices, then a single *solve nbody* task is submitted to each device. Then, for each timestep, forces between each particle pair are computed. Each of the two nodes compute forces for half the particles. Then, velocity and position of the particles assigned to each node are updated. Finally, the updated particles are sent to all other nodes.

Data send and receive operations are usually done in parallel and data movements can be overlapped with task execution. Also, multiple *calculate forces* tasks are run in parallel in multiple accelerators, although not shown in the diagram for clarity.

Comparing with flow described in Figure 119 for the host directly managing multiple devices it is obvious that the number of data movements are greatly reduced when data can be transferred between devices without host intervention. Also, the number of communication operations related to task submission and synchronization involving individual task execution for each particle block is greatly reduced.

This approach allows performance to scale when using multiple devices, as shown in Figure 123.

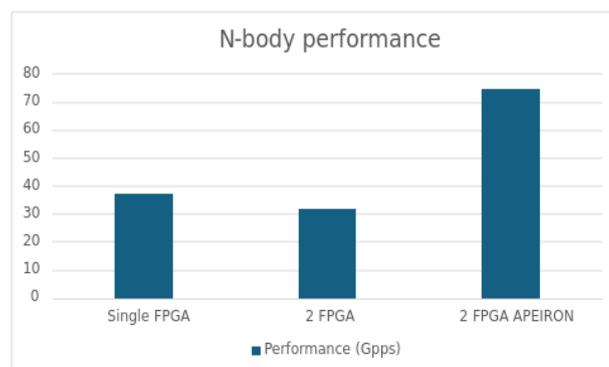


Figure 123: [NBody] Performance for different n-body FPGA implementations

The chart clearly shows that performance using apeiron for communications between both FPGAs performs much better than the previous host-centric approach. In fact, speedup achieved using this approach provides a speed up of 1.99x when compared to the single FPGA implementation. This is very close to the ideal speedup of 2x the results from doubling the available resources.

Also, power efficiency is maintained when moving from one to two FPGA devices. This is thanks to the low overhead caused by communications between FPGAs and between the host and FPGA devices. Energy efficiency results for different implementations are shown in Figure 124.

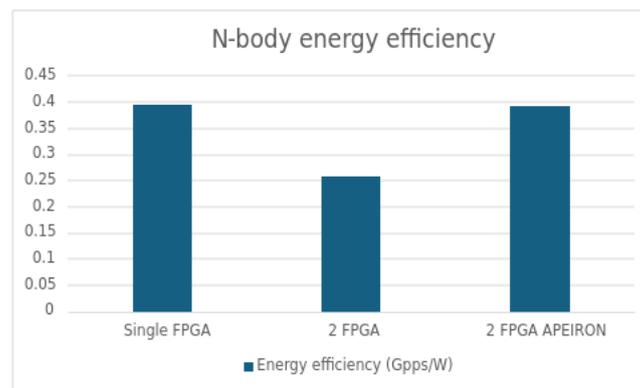


Figure 124: [NBody] Energy efficiency of n-body implementations

3.13.4 OmpSs@FPGA multinode applications

We evaluated the distributed n-body implementation across a larger number of devices. The design follows the same approach as the implementation using APEIRON described in section 2.12.3. Scalability tests are run in the MEEP cluster [14]. This cluster contains 96 Alveo U55c accelerator cards connected using 100G Ethernet installed across 12 different nodes. These cards use the same FPGA part as the Alveo U280, but without DDR. Also, the same application accelerators are used and therefore, per-FPGA performance should be the same as the IDV-E. In the MEEP system however, an external 100G ethernet switch is used instead of using the apeiron switch to route data packets to their destination. Nevertheless, setup should provide a good approximation of the performance achievable by a cluster of IDV-E nodes and allows us to evaluate the distributed approach across multiple nodes with multiple FPGA devices each.

Tests have been carried using up to 64 FPGAs across 8 different nodes. Scalability results are shown in Figure 125.

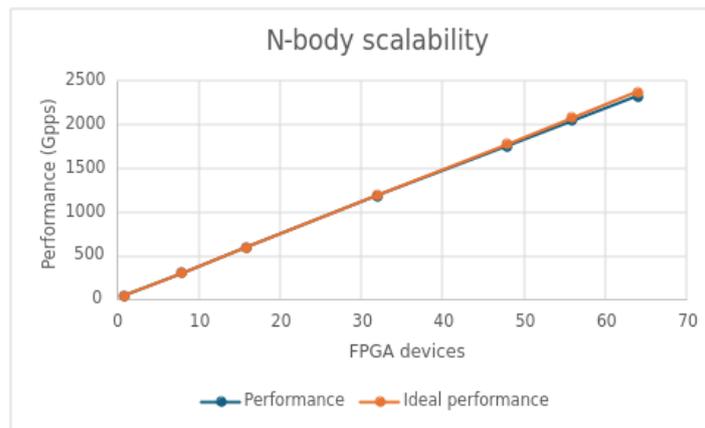


Figure 125: [NBody] Measured performance and ideal across different number of devices

This figure shows scalability results for the n-body application. For each number of nodes, it shows the actual, measured performance and the ideal performance. Ideal performance is the result of multiplying the performance of one FPGA by the number of devices.

Data presented in Figure 125 shows that measured performance very closely matches the expected performance. In the case of 64 FPGAs performance is 98% of the ideal performance. This demonstrates that the proposed distributed solution could still be efficient in a larger system comprised of several IDV-E nodes.

We compared scalability with a classical CPU-based system. The chosen system is MareNostrum 4, which as 3456 nodes with two Intel Xeon Platinum 8160 CPUs, 94GB of DDR4 RAM each, and 14nm Intel technology. It also features a 100Gb Intel Omni-Path Full-Fat Tree network. In our test, we measure performance for up to 56 nodes. Results are shown in Figure 126.

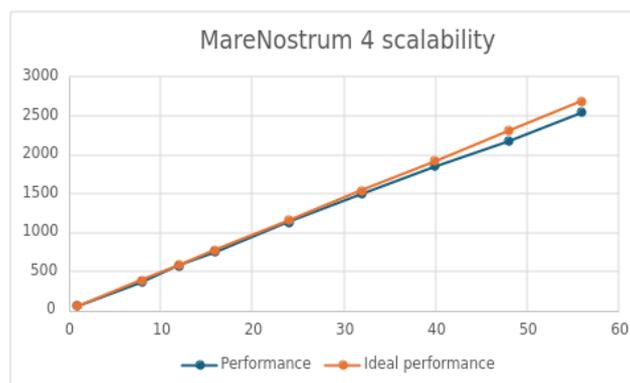


Figure 126: [NBody] MareNostrum 4 n-body scalability

This figure shows performance across different number of nodes. The chart shows two lines, one with the actual measured performance and another line with the expected ideal performance. Ideal performance is the product of the number of nodes by the performance on single node. While scalability is good, reaching 94.6% of the expected ideal performance, it's slightly lower than the 98% of the FPGA system.

Table 48 shows a summary of performance and power efficiency for IDV-E, MEEP multi-node FPGA system and MareNostrum4. The table reports the main KPIs (performance and power

efficiency) obtained by the multinode system compared against the same numbers obtained in the MareNostrum4 supercomputer to show that the experimental system developed is able to compete with a current supercomputer in terms of performance while delivering more than 2x times better power efficiency even if the FPGAs technology is 16nm against the next generation 14nm MareNostrum4 processors.

Cluster (size)	Performance (Gpairs/s)	Power (W)	Power efficiency (Gpair/W)
IDV-E (1)	37.43	94.84	0.394
IDV-E (2)	74.69	189.64	0.393
MN4 (32)	1488.08	9459.0	0.157
MN4 (56)	2536.60	15504.2	0.163
MEEP(32)	1175.26	3227.9	0.364
MEEP(64)	2322.98	6472.7	0.359

Table 48: [NBody] Summary of performance, power consumption and power efficiency across different clusters

4 Concluding Remarks

In this deliverable we presented application benchmarks and results with respect to the KPIs defined in the D6.1 deliverable and as they were advanced throughout D6.2 up to the state that they reached on the final platforms. Next subsections detail the relevant results and for each application.

4.1 Smart Cities - CINI-UNIFI

The key lessons learned thanks to the work at application levels are:

- Edge server applications for surveillance of smart cities can be set-up but to sustain in real-time combination of multiple complex algorithms like YoloV5+ Deep Sort a compute blade with a powerful GPU (Tesla A100) and processor (Xeon) should be adopted
- RISC-V computational capabilities, particularly in scalar version, and considering commercially available solutions are still far in performance from other platforms like those based on Intel and/or ARM.
- For ARM the SVE version in Fujitsu performed worse than the ARM Neoverse N1 in Ampere Altra Max (the one in IDV-E) that is not using SVE; this can be justified by the fact that the algorithms were not optimized for a scalar vector extension version of the processor.
- The availability of FPGA Xilinx Alveo in IDV-E can be exploited by porting the Yolo calculation on it while keeping the DeepSort on the ARM processor and hence a gain in speed can be achieved roughly by a factor of 2. However, the performance is still below those that can be achieved using Intel Xeon plus a GPU Like Tesla A100.
- Posits alone, implemented via SW library emulated by the processor, but without a hardware support, cannot bring an advantage vs. architectures like Tesla A100 already supporting natively mixed-precision (FP64, FP32, BF16, INT8, INT4)

These results are important output for the European community working on defining an EU-based alternative to the monopoly of Intel/AMD plus NVIDIA GPU solution.

In future projects and benchmarking activity of CINI-UNIFI will be extended to consider also RISC-V architectures with Vectorized Instruction Set and compare them to ARM-based (Neoverse N1, A64FX) and Intel-based GPP architectures.

4.2 MathLib - CNR

CNR developed some multi-GPU kernels for efficient use of heterogeneous architectures embedding last generation of Nvidia GPUs, as in the IDV-A platform. Activities were focused on the design of new algorithms and implementation patterns which can exploit at the best the combination of distributed-memory/shared-memory programming models, leveraging multiple GPU accelerators. Benefits of the library have been demonstrated both on IDV-A and the Italian Leonardo supercomputer, accessed by an Early Access Grant obtained by the CNR

team. The lesson learned is that for effective use of the new architectures it is not sufficient to re-factor legacy libraries and codes, but it is needed to rethink basic algorithms. The library was used as benchmark for the IDV-A platform and some project toolchain for energy consumption measurements before and after the installation of the new 2-phase cooling system, showing that this system has small impact on the performance results of the main kernel although can reduce the average power and then its total energy consumption. Concerning the aspect of energy consumption of the library modules, main concluding remarks are that using more GPUs for solving problems with fixed size is not efficient because generally, also if execution time is reduced, energy consumption increases in a significant way. When the problem size increases linearly with the number of GPUs, energy consumption of the main solver increases in an almost linear way, leading to good weak scalability not only in performance but also in energy consumption. Some more investigations are instead needed for the behavior of energy consumption of the SpMV, MWM and SpMMM kernels.

4.3 RTM - Fraunhofer

Both Floating Point formats Posit 16 bit, with one exponent bit, as well as Float 16 bit demonstrated the potential to be successfully used for RTM calculations based on the Marmousi reference data set. Different scenarios are available. Common is the storage of the wavefield in reduced precision. This can be extended by also calculating the kernel in reduced precision as well as stacking the image in reduced precision. If the Kahan summation is used for the stacking procedure all the numerical operations can be executed very well in reduced precision. The deviation towards the reference solution can be improved by calculating the kernel in Float32 or stacking the image in Float32. This presumes that the core provides the capability to perform Float32 as well as reduced precision operations.

4.4 HEP - INFN

The comparative analysis of the different available architectures yields valuable insights on the KPIs (throughput and energy efficiency) for the applications of interest. For the CPU only version of both the applications we observed a better scaling in energy efficiency running on the Intel(R) Xeon(R) Platinum 8470 based platform (IDV-A) with respect to the Dibona AMD-based cluster. Obviously, we have to take into account that the IDV-A platform shows a baseline power consumption of around 400 W. So in absolute terms, in the current state of the IDV-A, the Dibona shows better results regarding energy efficiency. We expect that these results will be reverted after some optimization activities on the thermal control daemon for the two-phases cooling system.

On the other hand, the ARM-based IDV-E node proved more effective on both throughput and energy efficiency compared to the Dibona and IDV-A nodes. The IDV-E architecture reaches the same throughput as the more performant of the two X86_64 ones (IDV-A) while surpassing by roughly a 50% the more energy efficient of them (Dibona) in the energy efficiency KPI.

When considering the performance of the IDV-A platform in comparison with the Dibona one in the CPU+GPU configuration, the improvement in throughput KPI for Pixeltrack is noticeable and expected in the IDV-A node. On the other hand, we measured a slightly worse performance of CLUE on IDV-A for the throughput KPI compared to that measured on Dibona node. Similarly, considering the energy efficiency KPI, there is a slight advantage in using the

Dibona node. These results hint for further investigation and specific code tuning for the Nvidia H100 GPU architecture.

4.5 NEST-GPU - INFN

The comparison of the CPU and GPU versions of the *hpc_benchmark* test is smashing in favor of the latter; in order to have a more complete reference against the D6.2 baseline this same multi-node benchmark is to be performed on the Dibona ARM cluster – the IDV-E used here is a single node so that inter-node communications were not gauged – but we do not expect a reversal. The same test was also performed on the state-of-the-art Intel CPUs of the IDV-A; on those, as in D6.2 we observed somewhat shorter runtimes at the cost of significantly increased energy-to-solution compared with those obtained on the ARM CPUs of the IDV-E mentioned in section 3.5 above but nothing close to the more than 17 times faster and 10 times smaller energy-to-solution a single IDV-A GPU can muster.

Even if the no-communications benchmark is faster than the multi-GPU ones the weak scaling that can be assessed with up to 4 GPUs of a single node IDV-A is quite good; for the accessible problem sizes that can be accommodated in the IDV-A in its current deployment the communications do not seem to be pose a significant hindrance yet – it remains to be seen for larger number of GPUs, possibly on more interconnected nodes – and there seems to be some performance to be squeezed out of optimizing the inter-GPU communications, as mentioned in section 3.5.

It is however worth noting that, at least for the present version of the neural network simulator, the NVIDIA H100 GPUs equipping the IDV-A are undoubtedly faster than all platforms they were compared to but as regards more general figures of merit like energy-to-solution and synaptic updates per Joule, it was found that the previous generation of GPUs (as the A100 ones on the Dibona cluster) are still quite competitive – and in some cases even winning – against the much more expensive H100.

4.6 RAIDER - INFN

The further co-development of the application, of the APEIRON scalable HLS framework, and of the INFN Inter/Intra-FPGA Communication IP led to significant improvements wrt the baseline reported in Deliverable D6.2.

In particular, in both testbeds used, we showed the scalability of the application in terms of throughput KPI, leveraging on the capabilities of the APEIRON framework in terms of multi-FPGA deployment and execution.

As can be extrapolated from the energy efficiency KPI values reported for the two testbeds used, for this class of applications a good scaling in term of energy efficiency can be reached only if the added nodes have enough processing capabilities that allow to increase the integrated processing throughput at a higher rate than that of the additional energy consumption they introduce in the system. This is the case of the IDV-E testbed (1 IDV-E node with 2 Alveo U280 FPGA boards), while in the INFN APE-Lab testbed we found the opposite situation, where the increased processing throughput capabilities when scaling the number of nodes (1 Intel node with 1 Alveo U200 FPGA board) were not enough to compensate for the increased energy consumption in the system. This is due to the higher density of FPGAs per node in the

IDV-E system (2 per node vs 1 per node), to the higher energy efficiency of the Ampere Altra processors compared to the Intel Xeon Silver counterparts in the APE Lab small cluster, and to the higher energy efficiency of the U280 FPGAs compared to the U200.

Regarding RAIDER integration with OmpSs, mixing task-based and streaming models can lead to increased performance with good programmability. However, tools and models need to be developed to exploit the advantages of both models. We also plan to further evaluate the integration of stream kernels into the task-based model. We plan to port the application to the MEEP machine with multiple FPGAs (64) and take advantage of the FPGA-to-FPGA direct link communication with the APEIRON communication framework (to reach a set of 32 FPGA pairs) if the kernel requires more than one FPGA to be deployed with enough performance as the preliminary tests suggest.

4.7 TNM - INFN

Regarding the Quantum matcha TEA (gate-based quantum circuit emulator for digitized quantum circuit, the relevant KPIs were reported in deliverable D6.2 as a baseline for the IDV-A node performance. For Quantum green TEA, i.e. solver for the Schrödinger equation or Lindblad equation, this is a new application developed after Deliverable D6.2. We reported a strong scaling study with multi-threading and MPI on single node towards finding the optimal split in hybrid multi-threading/MPI. The study was conducted using a node of the Justus2 cluster (bwHPC), with 2xIntel Xeon 6252 Gold with 2x24 cores, showing a ~5 speedup wrt the single core configuration. We pointed out the bandwidth towards the main memory as the main limiting factor to scalability. Furthermore, algorithmic improvements and comparison between performance on CPU and GPU on a CINECA Leonardo Booster node have been reported. We re-run the KPI measurements for quantum matcha tea on IDV-A for the quantum Fourier transformation (QFT) on 100 qubits as done in D6.2 on the Dibona node. We measured a performance increase in time-to-solution with respect to the Dibona machine, which leads to a better KPI for the gates per second, a the KPI for the energy consumption indicates an improvement as well. The outcomes and knowledge gained during the TEXTAROSSA project for the TNM will be applied via the collaboration of the Quantum TEA team with HPC or via its open-source code. In 2023, we have been working together with CINECA and Quantum Matcha TEA was installed on Leonardo. Moreover, the software is open-source and will be made public step-by-step in the future such that features as the data type switching are available to the public. Moreover, the code is used for European projects as Pasquans2 or EuRyQa, Marie-Curie fellowships like the UniPhD program at the University of Padova, or the QRydDemo project funded by the German Ministry of Education and Research. The applications within these projects have and will profit from the intellectual property gained during TEXTAROSSA.

4.8 MathLib - INRIA

Our work focuses on task scheduling for heterogeneous architectures. Our novel scheduler can support any combination of hardware types, including CPU/GPU and CPU/FPGA configurations, for instance. We are now actively developing an upgrade of our scheduler to reduce energy consumption, and our short-term objective is to demonstrate its benefits.

Accordingly, our next milestone will be to evaluate its performance on the Chameleon and ScalFMM applications.

During the project, we focused on developing applications that use FPGAs or GPUs, and CPUs. For FPGAs, we considered that a significant challenge is related to the software stack, which makes development costly both in terms of the learning phase and the programming/profiling efforts. Moreover, there are several outdated documentations (even on the vendor website) or installations on the available nodes. Besides, there are usually several possible technologies that can be used to produce the same result, which leads to ambiguity and complicates the development process. In addition, we did not achieve the expected performance, as the granularity required to be competitive against CPUs or GPUs was too large for the task-based parallelization, where the problem must be divided into tasks. This observation has been corroborated by the results of our partners, that achieved high-performance with large test cases.

In contrast, our work with GPUs, particularly on task scheduling, highlighted the critical role of CPUs when aiming to maximize performance. We found that CPUs are indispensable, especially with irregular applications where the performance difference between GPUs and CPUs varies significantly for different task types. This highlights the importance of creating advanced schedulers.

As a result, our guidelines have evolved modestly from the beginning of the project: we advocate for parallelizing applications using a task-based approach, as it allows getting high performance and productivity. This method facilitates the management of several critical aspects such as data transfers and scheduling decisions, allowing for near-optimal utilization of available hardware. However, the challenge of porting existing applications to a task-based parallelization framework appears tedious, and few applications were successfully ported during the project. Therefore, it seems that using task-based parallelization early in a project seems much more appropriate than porting an existing HPC application from MPI+X to a task-based runtime system.

4.9 UrbanAir - PSNC

The work carried out focused on exploiting heterogeneous resources, namely GPU accelerators, by UrbanAir-gcrk. Comparing to execution on CPUs, the time-to-solution was decreased 9x (multinode environment), while energy efficiency was increased by 2x (single node). With further work, for the same accelerators up to 10% increase in iterations per second was achieved, and up to 22% increase in energy efficiency. It clearly demonstrates that with proper implementation GPU accelerators are bringing possibility to run more efficiently in terms of execution time and energy efficiency compared to execution on CPUs. While there is memory constraint on GPU limiting the maximum size of a problem to be solved, there is no such constraint in multiple node environment, where each accelerator can compute its portion of data, so that a significantly larger problem can be solved. The eventual gains in speedup and energy efficiency depend on the problem and algorithms used. For problems similar to UrbanAir-gcrk, where there are grid points in 3D space and the value of each is calculated at each iteration/timestep by the same scheme (e.g. exchanging data with its neighbors), the adaption to heterogeneous may be simpler and benefits easier to achieve. The reason for this is

that in such cases the problem can be divided equally between accelerators, where each accelerator is computing their data independently, and just exchange the data with when necessary at each timestep. It is important to remember that only fully utilized GPUs are bringing the maximum gain in computational performance (and thus time-to-solution).

While TAFFO enables to introduce mixed-precision to the existing code, it is quite complicated to adapt large and complex ones using CUDA API, such as in case of UrbanAir. Developers of such codes should consider switching to SYCL, which makes using TAFFO easier and more efficient, and opens also possibility to use other accelerators such as FPGAs.

4.10 FIPLib - ENEA

The aim of the study performed over this image processing application was to design a library, this is FIPLib, written in C++ and useful for a stream-based hardware design to be implemented on an FPGA board via the Vitis HLS tool. It has been proved that, taking into account the limited computing resources on the accelerator card, it is possible to synthesize a bitstream to be flashed on the FPGA capable of overperforming (in terms of throughput and energy consumption, as reported in the previous section) its CPU model counterpart.

In addition to this, exploiting the capabilities of the APEIRON framework built over the INFN Communication IP, it has been proved that it is possible to implement even larger (in terms of resources) image processing application thanks to the possibility of multi-FPGA deployment. The results presented in the previous section show, in fact, how it is possible to implement FIPLib processing kernels over multiple boards maintaining the stream-based execution via the usage of the APEIRON HAPECOM Communication API. However, even if throughput performances seem to scale well with the number of boards in the setup, the energy consumption remains a parameter to take into account since the improvement in terms of processed image per second could result into a loss in terms of processed images per Joule with respect to the single FPGA implementation.

4.11 NBody - BSC

Evaluation of OmpSs@FPGA in different scenarios has highlighted that in task-based programming models, task management overhead can be critical for performance. It quickly can become a bottleneck, especially when using multiple accelerator devices as shown in section 3.13.2. Also, efficiently moving data between accelerator devices is critical when managing multiple devices, as shown in section 3.13.3.

Furthermore, proper cooling, which is often overlooked, can cause a non-negligible impact on energy efficiency and even on system usability if it is not properly addressed, as shown in section 3.13.1.

The work in the multinode capabilities of OmpSs@FPGA as shown in section 3.13.4, is expected to continue. Next steps include the inclusion of a Virtual Memory system in the OmpSs@FPGA framework to provide better scalability. Furthermore, we expect to leverage the Implicit Message Passing programming model, already under review, in a RISC-V manycore system with the help of the Fast Task Scheduler developed in Task 2.5. We plan to continue using NBody as a part of the test benchmarks used to evaluate these new developments along with other different applications (as maybe RAIDER, HPCG, etc.).

5 Achieved project objectives, lesson learned and guidelines

In this section focus is given to the evaluation of project results with emphasis on overall project objectives applications directly contributed to. In the next subsection, lesson learned and guidelines are summarised.

5.1 Project objectives

The following objectives of the project addressed by applications and use cases are discussed:

- **Energy efficiency** (Increase the energy efficiency by 2-5x compared with current available HPC solutions based on CPU or CPU+GPU thanks to: i) new energy-efficient processors and accelerators, ii) integrated use of FPGAs, iii) bi-phase cooling, iv) mixed-precision and data compression techniques, v) direct communication between FPGAs.)
- **Sustained application performance** (Achieve an increase in performance by 2-10x compared with current solutions based on CPU or CPU+GPU thanks to new AI accelerators and high-performance GPUs and FPGAs, HW fast schedulers, low-latency inter-/intra-node links, better cooling allowing higher frequencies, fast flow enabled at programming mode and run time levels.)
- **Seamless integration of reconfigurable accelerators** (TEXTAROSSA vendors will adopt tools and solutions of the project.)
- **Integrated Development Platforms (IDVs)**
- **Development of new IPs**
- **Fine-tuned thermal policies integrated with an innovative cooling technology**

5.1.1 Energy efficiency

The goal of a 2-5x increase in energy efficiency is reached by orchestrating the design and developments done in WP1, WP2, WP3, WP4, WP5 and WP6 in a co-design approach. The assessment of the objective involved the measurement of KPIs for energy efficiency identified in D6.1 for each application, both on IDV-A (CPU+GPU) and IDV-E (CPU only and CPU+FPGA) and the comparison against the baselines obtained on the AMD EPYC 7402 x86-64 CPU with and without the NVIDIA A100 GPU acceleration provided by the ATOS partner in Dibona OpenSequana blade system and reported in D6.2, and then on NVIDIA H100 GPU accelerators provided in IDV-A.

CNR Mathlibs – IDV-A – very good energy efficiency scalability

CNR designed and implemented some kernels for sparse matrix computations and iterative linear solvers, which efficiently use heterogeneous architectures embedding last generation of Nvidia GPUs, as in the IDV-A platform to increase the energy efficiency. For an execution on 512 GPUs on Leonardo EuroHPC-JU supercomputer, the time-to-solution is 2.6x shorter comparing to NVIDIA AmgX solver, which results in a similar increase in energy efficiency.

RAIDER – IDV-E – APEIRON, Communication IP, mixed-precision, OmpSs – 1.6x increase

The RAIDER High Performance Data Analytics on FPGA application addressed the energy efficiency goal along two directions: i) energy efficient dataflow processing using dedicated pipelines in the FPGA fabric and using mixed precision techniques to implement the Convolutional Neural Networks performing the inference task, and ii) energy efficient intra/inter-FPGA communication mechanism provided by the Communication IP integrated in the APEIRON framework to enable scalability when deploying the application on a multiple FPGA system with a stream processing farm scheme. In D6.2 we reported a O(10) improvement factor in energy efficiency – measured as processed events per Joule – for the preliminary single-FPGA version of the RAIDER application when comparing it to those obtained on a CPU only and CPU+GPU platforms performing the same computational task. Performance of this single-FPGA version, implemented on a Xilinx Alveo U200 board before IDV-E became available, has been taken as baseline to assess the scalability performance of the final multi-FPGA version of the application. Using two IDV-E nodes available at E4 premises, each hosting two Xilinx Alveo U280 boards, we increased the system size to include up to four directly interconnected FPGAs, measuring a linear scaling in energy efficiency, with a +60% improvement w.r.t. the single U200 FPGA baseline in the four U280 FPGA configuration.

N-Body – IDV-E – OmpSs, IMP, CommunicationIP, FTS – 3x increase

N-Body application executed using OmpSs@FPGA (integrating the Fast Task Scheduler IP developed in Task 2.5) over a cluster of IDV-E FPGAs achieves a nearly 3x improvement in terms of energy efficiency when compared against Marenostrum 4 supercomputer.

NEST – IDV-A – 22x increase

Achieving the efficiency goal by NEST applications was considered by exploiting state-of-the-art GPU accelerators available on IDV-A and low-power CPU available on IDV-E. Execution of the neural simulation with NEST on the IDV-E Ampera Altra CPUs showed an improvement by roughly a factor of two w.r.t. AMD EPYC 7402 x86-64 CPU, while the same simulation with NEST-GPU using a single NVIDIA A100 yielded a more than twenty-fold increase in energy efficiency w.r.t. the same X86-64 CPU. The measured energy efficiency KPI between the IDV-A node and the Dibona OpenSequana one yields a ~20% decrease running the simulation with NEST-GPU on the four available GPUs, despite the measured increase in the throughput KPI. Considering pure CPU performance, running the simulation with NEST on the two IDV-E Ampera Altra CPUs showed an increase of ~50% in energy efficiency compared to the two X86-64 Intel Sapphire Rapids IDV-A CPUs, reducing the efficiency gap measured between the IDV-E ARMv8 and X86-64 measured with the Dibona node.

HEP – IDV-A, IDV-E – 4-11x increase

High Energy Physics aimed to reach energy efficiency goal by exploiting GPU accelerators available on IDV-A and low-energy CPU cores available on IDV-E. On the new IDV-A platform, energy efficiency yielded increases of x11 for CLUE application and x4 for Pixeltrack application. Moreover, the energy efficiency remains constant while scaling the number of used GPUs. When compared with the baseline results using four GPUs, the energy efficiency on IDV-A shows a decrease of about 10%, iii) When run on CPUs only, the ARMv8 Ampera Altra based IDV-E's energy efficiency is three times higher than that measured on the X86-64 Intel Sapphire Rapids based IDV-A; comparing IDV-E energy efficiency with the

Dibona CPU baseline, we see that for Pixeltrack, IDV-E surpasses it roughly by a 50% while being on-par with it on CLUE. We investigated the rather poor performance of the IDV-A CPUs and identified the high-power consumption of the CPUs node when idle as the main reason for it. This was a clear indication for a possible tuning for the platform.

TNM – IDV-A – mixed-precision techniques – 25-150% increase

The TNM application received several optimizations on the algorithmic side (regarding sampling the quantum state represented as a tensor network and the Hamiltonian matrix representation and its application to a state vector) with gains from 25% up to 150% glimpsed from the Dibona node and on the Leonardo HPC cluster at CINECA. Another line of investigation pursued using a mix of single and double precision steps in the search for the ground state energy of a 128-qbit system, which showed that with the optimal mix the same final precision can be reached with less than one third energy-to-solution compared to a double precision-only computation.

Smart Cities – GPUs and FPGAs – mixed-precision

Smart cities achieved increase in computation speed, i.e. time to solution KPI, at the node level by exploiting hybrid programming models for heterogeneous architectures and use of mixed-precision available in accelerators like H100 GPU that is the same GPU of IDV-A. Quantitative evaluations are reported in the journal publication [16].

UrbanAir – IDV-A – 2x increase

UrbanAir increase in energy efficiency was expected by exploiting GPU accelerators available on IDV-A. Kernels demonstrated that by using GPU accelerators we can increase energy efficiency of the execution. Moreover, for a fixed problem size, with the increasing number of GPUs used, we can keep the same level of energy consumption. With the increasing number of GPUs and fixed problem size per GPU (weak scaling), the energy consumption increases linearly. Obtained results report up to 2x improvement in energy efficiency when using GPU accelerators comparing to CPU-only execution. Further improvements for IDV-A resulted in 22% increase in energy efficiency when executing on GPU accelerators.

FIPLib – IDV-E – APEIRON, Communication IP – 75x less energy

FIPLib application expected energy efficiency increase by exploiting FPGA accelerators available on IDV-E. The results are that implementation for FPGA uses 75x less energy compared to CPU implementation, while the Energy Delay Product is 582x smaller compared to CPUs.

5.1.2 Sustained application performance

The goal of increasing the performance by 2-10x is achieved by developments done in WP1, WP2, WP3, WP4, WP5 and WP6 thanks to new AI accelerators and high-performance GPUs and FPGAs and tools developed within TEXTAROSSA.

CNR-Mathlib – IDV-A – 2.6x better than NVIDIA AMGX

CNR-MathLib benefitted from high-throughput GPU accelerators and showcased very good speedup and scalability potential in exploiting multi-GPU nodes for computations on memory/communication bound sparse matrices/graphs, which also outperform the state-of-the-art Nvidia library for similar computations by 2.6x.

Smart Cities – GPUs and FPGAs

Smart cities achieved increase in sustained performance, i.e. number of processed frames , at the node level by exploiting heterogeneous CPU-GPU architectures and use of mixed-precision available in accelerators like H100 GPU that is the same GPU of IDV-A. Quantitative evaluations are reported in the journal publication [16].

RAIDER – IDV-E – APEIRON, Communication IP, mixed-precision, OmpSs – 8x increase

The RAIDER multi-FPGA application benefitted from the Communication IP services, offered by the APEIRON software stack, showing a close-to-perfect linear scaling of processing throughput – expressed in terms of processed events per second – with the number of used FPGAs in the two IDV-E nodes setup (four U280 FPGA boards), reaching an x8 increase over the single U200 baseline (and hence a O(100) increase over the CPU+GPU baseline used in D6.2).

IDV-E – OmpSs, IMP, CommunicationIP, FTS – 100x increase

BSC Fast Task Scheduler showed impressive performance increases, reaching over 100,000 tasks per second in FPGA systems and achieving a 100x speedup over the software alternative in RISC-V manycore systems when executing a set of HPC benchmarks.

TNM – IDV-A – mixed-precision techniques – 50% increase

The TNM application was motivated by the gains obtained investigating the single-double precision mix in testing a GPU-accelerated ground state search implementing such mix; results are variable, with up to a 50% better runtime on the GPU for large enough system sizes but with a smaller and smaller margin as soon as the mix leans more and more towards single precision.

NEST-GPU – IDV-A – 8x increase

NEST-GPU reported significant improvements when running on the Dibona OpenSequana GPU accelerated blade, as it was reported in D6.2, where we measured a 8x gain in performance KPI moving the execution from the two EPYC CPUs to a single A100 GPU. In the second reporting period we performed a weak scaling study on the four H100 GPUs on IDV-A and compared results with the four A100 GPUs Dibona baseline, observing only a marginal increase in performance KPI (+8%), but still x15 higher than what can be obtained running on the IDV-A two X86-64 Intel Sapphire Rapids CPUs and x12 than what is achieved on the IDV-E ARMv8 CPUs.

HEP – IDV-A, IDV-E – 5-14x increase

HEP applications, CLUE and Pixeltrack, were used to assess the attainable processing throughput in terms of reconstructed Physics events per second both in the IDV-A (CPUs, CPUs plus one to four GPUs) and in the IDV-E (CPUs). Comparing results with the Dibona OpenSequana blade baseline we observe that: i) on IDV-A, processing throughput scales linearly with the number of GPUs for both applications as it was the case for the baseline, showing an increase of +11% (Pixeltrack) and a decrease of -24% (CLUE) of the maximum attainable value, ii) on IDV-A, throughput scales linearly with the number of used cores on both CPUs, without saturating before reaching the max available cores as it was in the case of the baseline, and it shows increases of +28% (Pixeltrack) and of +6% (CLUE), and iii) on

IDV-E ARMv8 CPUs, processing throughput are marginally higher (+17% for CLUE and +7% for Pixeltrack) than those obtained on IDV-A X86-64 CPUs.

INRIA Mathlibs – IDV-A – StarPu, scheduler – 30% increase

INRIA's Mathlibs focused on contributing to performance objective by scheduling of task-based applications (scheduler that can improve the performance up to 30% on computing nodes with multiple GPUs, compared to state-of-the-art schedulers) and the GPU accelerators - the gain can be 60% on extreme test cases.

UrbanAir – IDV-A – TAFFO – 9x increase

UrbanAir benefitted from high throughput GPU accelerators to achieve 9x speedup by leveraging multi-GPU accelerators comparing to execution on CPUs. Further 4% increase in computational performance is observed with additional optimisations for IDV-A. 1.25x speedup with TAFFO.

FIPLib – IDV-E – APEIRON, Communication IP – 7.7x increase

FIPLib implemented wide parallelism and deep pipelining for FPGA available on IDV-Eto increase speed by 7.7x comparing to execution on CPUs.

RTM – FPGA EPI SGA2– mixed-precision – up to 2x increase

In RTM, FPGA platform which was developed within the EPI SGA2 project was used. Different scenarios were analysed where half precision is used just for storage, for calculating parts of the main propagation kernel or to use everywhere. The speedup depends on the chosen algorithmic setting, especially the storage of the calculated image (F32, P16 or P16 with Kahan summation) other than the storage of the simulation domain. 2x increase is an upper boundary for the total calculation.

5.1.3 Seamless integration of reconfigurable accelerators

The availability of reconfigurable accelerators integrated in a seamless way allows applications to achieve energy efficiency and sustainable performance goals.

APEIRON is an extension of the Xilinx Vitis HLS framework able to support a network of FPGA devices interconnected with a low-latency direct network implemented by the INFN Communication IP as the reference execution platform. Developers can define scalable applications using a dataflow programming model that can be efficiently deployed on a multi-FPGAs system: the APEIRON communication IPs allow low-latency communication between processing tasks deployed on FPGAs, even if hosted on different computing nodes. Thanks to the use of High-Level Synthesis tools in the workflow, namely Xilinx Vitis, tasks are described in a high-level language (C++) while communication between tasks is expressed through a lightweight C++ API based on non-blocking *send()* and blocking *receive()* operations. The APEIRON framework has been validated on the two IDV-E nodes system – each of them hosting two Alveo U280 FPGA boards – installed at E4 premises through the execution of the HPDA/AI RAIDER and of the HPC N-body applications, and on a four X86-64 nodes system – each hosting one Alveo U200 FPGA board – installed at INFN APE Lab through the execution of the ENEA FIPLib image processing library and of the RAIDER application.

OmpSs@FPGA for IDV-E supports seamless integration of reconfigurable accelerators. As explained in D4.1, D1.4 and D6.3, the OmpSs@FPGA framework can submit C/C++ programs directly to CPU+FPGA heterogeneous systems only by using the pragma “target device(fpga)”. The accelerator is compiled directly through our LLVM+AIT system (using Clang and Vitis HLS as backends) and data transfers and synchronization are managed directly by Nanos runtime and our Fast Task Scheduler IP. The system has been extended using the IMP (Implicit Message Passing) and OMPIF (OmpSs + MPI for FPGAs) programming models to support multi-FPGA multi-node systems. Consequently, currently reconfigurable accelerators can use the OmpSs@FPGA framework to be seamlessly integrated in multi-node heterogeneous FPGA systems using the same source codes used to program single-node SMPs as reported in D6.3 section 3.13. The fact that only one pragma needs to be added to an original SMP source code to obtain an accelerator working in IDV-E FPGA is the best key indicator of the successfulness of the task. As an example, the N-Body code uses three pragma lines (because it has three user accelerators) to execute in FPGA a code that is hundreds of lines long.

FastFlow is a structured parallel programming environment targeting shared and distributed memory (COW) architectures. In TEXTAROSSA FastFlow has been extended to allow the seamless integration of FPGA kernels, developed through the Vitis HLS flow in such a way the kernels can be executed in the places where currently only CPU computations (threads) may be executed (e.g. pipeline states, map workers, etc). Once the kernel has been compiled through the HLS flow, generating the bitstream, it is sufficient to pass to FASTFLOW the bitstream name and the top-level kernel parameters to have the kernel automatically instantiated in a FASTFLOW computation. All the low-level details concerning the kernel execution and data movement are transparently managed by FASTFLOW, including synchronization with other FASTFLOW nodes and double buffering to hide communication.

5.1.4 Integrated Development Platform

Two integrated development platforms (IDV) have been designed and built during the project to support both existing and emerging High-Performance Computing (HPC) applications. Their architecture is heterogeneous, comprising CPUs, GPUs and FPGAs. Also, frameworks (OmpSs@FPGA, APEIRON, StarPU, Fastflow or Taffo) have been developed targeting these specific platforms to support efficient execution of applications on both IDVs. The applications development contributes to the development of IDVs by confirming the correctness of the developed frameworks and obtaining the desired KPIs.

5.1.5 Fine-tuned thermal policies integrated with an innovative cooling technology

Two-phase cooling system demonstrated its capacity to efficiently cool CPUs, Memory, FPGAs, and GPUs even in the worse scenario with water-glycol temperature of 45°C. The tests revealed that the selected condenser for the IDV-A showed its limitations and for future developments it must be optimized, while pumping power remained low (1%) in comparison to the total thermal power rejected.

The tests carried out during the TEXTAROSSA project confirmed the technical feasibility of rejecting waste heat effectively in hot climates without energy-intensive measures of chillers (HVAC) or without high water consumption of evaporative towers. This confirms that two-phase cooling is a new sustainable solution for effective thermal management of exascale systems.

Moreover, it is important to underline that for HPC platforms the critical sources from the power standpoint it is not only the CPU but encompass all the main element of the platform (e.g., CPU, GPU, Memory, etc) that must be cooled by adding ad-hoc evaporators, with possibly ad-hoc designed global control strategies. In fact, computer architectures for HPC have high power density with spatial nonuniformity (hot spots) and temporal nonuniformity (thermal transients). Facing such problems with solely package-level cooling could not be enough and an on-chip dynamic thermal management is a viable solution explored by both industry and academia. To this purpose a dedicated control strategy must be put in place.

The solution successfully explored in TEXTAROSSA operates with multiple actuators each with their specific characteristics. A fast on-chip actuator, such as DVFS, is used to guarantee fast enough operation to counteract thermal transients, although at a computational performance penalty. To limit said performance penalty, the control algorithm operates also on the parameters of evaporative cooling cycle, such as controlling the coolant flow rate. Given the different timescales at play, a hierarchical control strategy appears to be the most promising solution, with a fast inner control loop acting on DVFS and a slower outer loop controlling the coolant flow rate. This control architecture also has the added benefit to allow to control the evaporative coolant vapor quality, thus improving the cooling performance of the evaporative cooling while preventing it from reaching a critical heat flux condition, which lead to computing hardware overheating and failure.

During the project we showed that it is possible to keep under control the operating temperature of both GPU and CPU with limited performance penalty. For the GPU the accuracy of the control was in the order of 1° (w.r.t. the single sensor per GPU that is present) while for up to 104 X86 cores, the accuracy was of 2° , considering that we have one sensor per core, so the control is capable to operate at a very fine grain. Details are reported in WP3 deliverables.

5.1.6 Development of new IPs

The list of the exploitable results and generated IPs is reported in D7.5. Specifically, 17 out of 18 innovation products reached the target TRL or better, including both open source and commercial innovation products.

It is worth noticing that University of Torino purchased the first commercial E4 prototype of a GPU-enabled server with evaporative cooling, based on the TEXTAROSSA technology. It is installed in the HPC4AI green data center for wider testing with the existing commercial users and to foster further scientific cooperation with the project partners.

Concerning HW IPs developed in WP2, Posit Processing Unit (PPU) in T2.1, Shake and NTT (Number Theoretic Transform) for advanced cryptography in T2.2, Light PPU for compression in T2.3, Communication IP in T2.4, Fast Task scheduling in T2.5, for all of them a data base is available as open source. All of them have been synthesized on target Alveo Xilinx FPGA, the same available in IDV-E, and all of them have a standard interface that allowed an integration with open-source RISC-V cores (32 bit and 64 bit versions, single and multi cores). The IPs have been exchanged among partners and used with WP4 tools and WP6 applications to assess their performance, their usability and their functionality. The value of the IPs has been also verified with the scientific community with journal publications and presentations to conferences/workshops.

Details of applications using new IPs are presented within Sections 5.1.3 and 3.

5.2 Lesson learned and guidelines

This Section summarises lesson learned and guidelines for Integrated Development Platform, project toolchain and applications

IDV-A:

- Speed-up for all-kind of applications, including increase energy efficiency
- It is favoured over FPGA in particular for small granularity problems or mixed-precision due to built-in support for these data types
- Software stack is more mature on GPUs
- Many sources of noise (system daemons, different power governors and policies) are present while measuring power – it is not always straightforward how to take them into account and sensibly remove them

IDV-E:

- Speed-up can be obtained for problem whose processing nature fits FPGA
- Gain is observed in particular for energy efficiency
- The IDV-E platform offers suitable FPGA resources to enable the scaling of performance and energy efficiency in Multi-FPGA HLS streaming applications leveraging the APEIRON framework and the low-latency direct network implemented by the Communication IP.
- IDV-E offers the basis for multi-node FPGA execution of OmpSs@FPGA based applications.
- IDV-E, in addition, offers the platform necessary to evaluate large nodes of many-core RISC-V SMPs that incorporate the FTS IP. FPGAs are a key enabler for prototyping.

APEIRON (including the Communication IP as part of the framework):

- The co-design of the Communication IP and of the software stack and tools, allowed to offer a framework for the straightforward development of scalable multi-FPGA High Level Synthesis applications. This has been demonstrated in the project with three applications belonging to different domains and making use of different parallelization paradigms: RAIDER (HPDA/AI, stream processing farm), FIPLib(Image Processing, data parallel), and N-body (HPC, distributed memory/message passing).
- Multi-FPGA systems, such as the one implemented with two IDV-E nodes (four FPGAs) and the APEIRON framework, can offer unmatched performance for high throughput/low latency data analytics tasks both in terms of sustained performance and of energy efficiency.
- Algorithms requiring many kernels with a high degree of parallelism may not fit in the FPGA. In such cases, application can use the APEIRON communication HW/SW stack, which allows enhancing the implementation by spreading the kernels among different FPGAs

FastFlow:

- Seamless management of Vitis kernel offloading in streaming applications can be achieved at no additional overhead with respect to OpenCL-only offloading management

- Fast prototyping of alternative solutions for stream parallel orchestration of Vitis kernel offloading does not require specific FPGA/Vitis toolchain knowledge to the application programmer
- Runs on IDV-e but also on different kind of Linux Intel/AMD nodes (experimented with single or double U50 or U280 Alveo boards).

TAFFO:

- The use of mixed-precision can improve latency in GPGPU at the price of some accuracy loss. This is suitable for “urgent computing” scenarios, where the need for a timely result justifies the loss in accuracy. Furthermore, it must be noted that performance gains are strictly bound by the availability of vectorized units. On GPGPU, this is usually limited to replacing IEEE754 32-bit floats with 16-bit floats or bfloat16, with a maximum speedup of 2x. Furthermore, on heterogeneous platforms such as GPGPUs, it is necessary to co-optimize the use of data types between host and device code. Without this optimization, no benefits can be reaped, since data type conversions needed at the interface are likely to offset the benefits of vectorized computation. Finally, employing reduced size data types in data transfer between host and device can speed up data intensive tasks.
- Posit data types can improve precision over IEEE754 floating point types. Thus, they are a good candidate for use in mixed precision computing. In particular, a multiprecision Posit unit including 32-bit Posit and 2x packed 16-bit Posit arithmetics could be effective, by significantly reducing error while allowing performance improvements. Since Posit units are also much smaller than comparable FPUs, the trade-off would be more favorable also in terms of density on custom hardware. TAFFO proved effective at performing the tuning on such number systems, also in mixed precision scenarios.

OmpSs@FPGA:

- In large FPGA systems the SMP can be the bottleneck hindering overall performance (or energy efficiency) unless part of the control code is moved to the FPGA fabric.
- Dense computations with simple control codes are best suited to be executed in GPUs
- Applications should drive the improvements, but compilers are necessary to make such improvements reach computer systems users.
- Programmability and portability are key to move to next large heterogeneous platforms.

FTS IP:

- Moving from hardware task-aware runtimes to task-aware architectures (with hardware support) is key to unlock the next level of performance in many-core architectures
- Accelerator-aware architectures (where hardware IPs manage the hardware resources as communication IPs) are also a convenient step to achieve the necessary scalability in large systems

Applications:

- Effective use of complex architectures, such as heterogeneous GPU-accelerated nodes, requires rethinking of basic algorithms. E.g., in the case of CNR Mathlib we demonstrated that algorithms which are not optimal from the convergence point of view are better suited for SIMD computations and then are able to obtain better time to solution than the optimal-convergence alternatives. Therefore, re-factoring of legacy

libraries can be not sufficient to fully exploit the potential of heterogeneous computing nodes. On the other hand, suitable data structures, adaptivity of the computation with respect to the sparsity pattern, overlapping data communication among GPUs/CPU's to computation and reduction of global synchronizations are key approaches to obtain high efficiency. We also demonstrated that using more GPUs in a weak scaling approach is more energy efficient than using more resources to accelerate fixed-size problem solution. Indeed, good scalability at performance level produces good scalability in energy consumption, resulting in an almost constant Dofs (Degrees of freedom) to Joule ratio.

- GPUs can be competitive with actual supercomputers with the same technology node for the right applications, especially from the point of view of energy efficiency
- RISC-V computational capabilities, particularly in scalar version and considering commercially available solutions, are still far in performance from other platforms like those based on Intel and/or ARM plus GPU
- FPGAs can be competitive with actual supercomputers with the same technology node for the right applications, especially from the point of view of energy efficiency
- Programmability and portability are key to enhance the spectrum of computing platforms
- Posit16 demonstrates advantages vs FP32 when compared to FP32 in RTM case, with storage size saving by a factor roughly 2
- Single source code (Alpaka) makes difficult to reach optimal performance on new GPU versions without recasting the application.
- GPUs are unquestionably the platform of choice for the neural network simulation, either regarding raw performance or power efficiency

6 References

- [1] Zenker, Erik & Worpitz, Benjamin & Widera, Rene & Huebl, Axel & Juckeland, Guido & Knüpfer, Andreas & Nagel, Wolfgang & Bussmann, Michael. (2016). Alpaka - An Abstraction Library for Parallel Kernel Acceleration. 10.1109/IPDPSW.2016.50.
- [2] Emanuele Ruffaldi, Federico Rossi, Marco Cococcioni, "cppPosit: a C++ template-based library implementing Posit arithmetic ready for machine learning applications", <https://www.techrxiv.org/doi/full/10.36227/techrxiv.24139605.v1>
- [3] [2021-05-06-NVIDIA Ampere Architecture In Depth NVIDIA Developer Blog.pdf \(yyrcd.com\)](#)
- [4] Heroux MA, Dongarra JJ. Toward a New Metric for Ranking High-Performance Computing Systems. Sandia National Lab. Tech. Rep., SAND2013-4744, 2013.
- [5] Bernaschi M, Celestini A, D'Ambra P, Vella F. Multi-GPU aggregation-based AMG preconditioner for iterative linear solvers, 2023. IEEE Transactions on Parallel and Distributed Systems, vol. 34, 8, 2023. <https://doi.org/10.1109/TPDS.2023.3287238>
- [6] NVIDIA, Algebraic multigrid solver (AmgX) library, rel.2.1, 2020. <https://github.com/NVIDIA/AMGX>.
- [7] Naumov M, Arsaev M, Castonguay P, Cohen J, Demouth J, Eaton J, Layton S, Markovskiy N, Reguly I, Sakharnykh N, Sellappan V, Strzodka R. AmgX: a library for GPU accelerated algebraic multigrid and preconditioned iterative methods, SIAM Journal on Scientific Computing, 2015, 37, S602–S626.
- [8] A Modular Workflow for Performance Benchmarking of Neuronal Network Simulations – Neuroinform., 11 May 2022 Volume 16 – <https://doi.org/10.3389/fninf.2022.837549>
- [9] Extremely Scalable Spiking Neuronal Network Simulation Code: From Laptops to Exascale Computers – Front. Neuroinform., 16 February 2018 Volume 12 – <https://doi.org/10.3389/fninf.2018.00002>
- [10] Runtime Construction of Large-Scale Spiking Neuronal Network Models on GPU Devices – Appl. Sci. 2023, 13(17), 9598 – <https://doi.org/10.3390/app13179598>
- [11] Ballarin, M., Silvi, P., Montangero, S., & Jaschke, D. (2024). Optimal sampling of tensor networks targeting wave function's fast decaying tails. arXiv preprint arXiv:2401.10330.
- [12] Xizhou Feng, Rong Ge, K.W. Cameron: Power and Energy Profiling of Scientific Applications on Distributed Systems. Proceedings. 19th IEEE International Parallel and Distributed Processing Symposium, 2005.
- [13] J. M. de Haro *et al.*, "OmpSs@FPGA Framework for High Performance FPGA Computing," in *IEEE Transactions on Computers*, vol. 70, no. 12, pp. 2029-2042, 1 Dec. 2021, doi: 10.1109/TC.2021.3086106.
- [14] K. DeVogeleer, G. Memmi, P. Jouvelot and F. Coelho, "Modeling the temperature bias of power consumption for nanometer-scale CPUs in application processors," *2014 International Conference on Embedded Computer Systems: Architectures, Modeling, and Simulation (SAMOS XIV)*, Agios Konstantinos, Greece, 2014, pp. 172-180, doi: 10.1109/SAMOS.2014.6893209.
- [15] MEEP: Marenstrum Exascaleme Entry Project, D6.4 -Full Emulation prototype release
- [16] <https://www.mdpi.com/2079-9292/13/9/1757>