



A Parameterizable CGRA-based Engine for Machine Learning Acceleration On-board Satellites



Matteo Monopoli matteo.monopoli@phd.unipi.it PhD Student University of Pisa Dept. Information Eng.







L. Zulberti, M. Monopoli, G. Mystkowska, S. Moranti, P. Nannipieri, L. Fanucci

6th SEFUW: SpacE FPGA Users Workshop

eesa

Gabriela Mystkowska gabriela.mystkowska@phd.unipi.it PhD Student University of Pisa Dept. Information Eng.



26th March 2025



Outline

- 1. Background
 - a. Motivation
 - b. Hardware for Al Onboard
- 2. State-of-the-Art
- 3. CGRA-based HW Accelerator
 - a. ESA OSIP Context
 - b. CGR-Al Engine
 - c. Ongoing Work
- 4. Std. Cell Synthesis & FPGA Prototype
- 5. Conclusion

SEFUW 2025









Outline

- Background 1.
 - a. Motivation
 - b. Hardware for Al Onboard
- 2. State-of-the-Art
- 3. CGRA-based HW Accelerator
 - a. ESA OSIP Context
 - **CGR-Al Engine** b.
 - Ongoing Work
- 4. Std. Cell Synthesis & FPGA Prototype

Conclusion

SEFUW 2025











Background – Motivation

- Al Applications in orbit:
- Remote Sensing
 - Object detection
 - Weather forecast
 - Earth observation
- Autonomous spacecrafts
 - Vehicle docking
 - Probes
 - Landers, rovers
 - Deep space missions
 - Fault detection and isolation, recovery
- Data privacy increase
 - Downlink requirements reduction

SEFUW 2025









Background – Hardware for AI Acceleration On-board

- High Computational Efficiency
 - Inference
 - Real-time processing in resource-constrained environments
 - Tailored for parallel processing, and NN operations
- Energy efficient
- **Compact and Integrated Design:**
 - Small form factors
- Customizability and Flexibility:
 - Reconfigurable to meet specific AI requirements
 - Adaptable to evolving AI models and algorithms without requiring new hardware
- Reliability
 - Must be able to withstand exposure to radiation
 - Protection against Single Event Effects (SEEs) through radiation-hardening or fault mitigation

SEFUW 2025









Outline

- Background a. Motivation b. Hardware for Al Onboard
- 2. State-of-the-Art
- 3. CGRA-based HW Accelerator
 - a. ESA OSIP Context
 - **CGR-Al Engine** b.
 - Ongoing Work
- 4. Std. Cell Synthesis & FPGA Prototype

Conclusion

SEFUW 2025











State-of-the-Art

	Architecture	Chip	Pros	Cons	
NOT RAD-HARD	GPU	NVIDIA Xavier, NVIDIA Jetson Orion, AMD Ryzen	COTS, Not rad-hard,		
	VPU	Intel Myriad X	TOPS/W)	only LEO	
	TPU	Google Coral			
RAD-HARD	FPGA	Xilinx Zynq UltraScale+, Microchip PolarFire	Flight heritage, Rad-hard,	General purpose,	
	CPU	Gaisler GR740	Flexible, Reprogrammable	performance	
	Spatial Architecture	Xilinx Versal	Rad-tolerant, High performance (up to 430 TOPS)	Low interpretability ("black box") Low predictability in time critical tasks, Relatively high power consumption (35W)	
	Systolic Array	HPDP	High flexibility High reconfigurability Low power consumption	Relatively old technology (STM 65 nm), Proprietary tools	

SEFUW 2025









SotA – Technology Available over the Years



SEFUW 2025









SotA – Resource Optimised CGRA Hardware Accelerator

- 2D array of PE interconnected through NoC
- Speeding up compute-intensive inner loops
 - Digital signal processing
 - AI/ML workloads
- Heterogeneous computing, high parallelism
- Highly predictable (timing model)
- Highly parametrized architecture
 - Optimised resource utilisation
 - Scalable and flexible
- Low reconfiguration time

SEFUW 2025





	 Power efficient 		FPGA
CGRA is	Power efficientHigher performance	compared to	CPU
	 More flexible 		Systolic Array
	More flexible		ASIC







Outline

- Background a. Motivation b. Hardware for Al Onboard
- 2. State-of-the-Art
- 3. CGRA-based HW Accelerator
 - a. ESA OSIP Context
 - CGR-Al Engine b.
 - Ongoing Work **C.**

4. Std. Cell Synthesis & FPGA Prototype

Conclusion

SEFUW 2025







UniPi – CGR-Al Engine

- CGRA-based accelerator for Al on-edge applications
 - Time-constrained applications (e.g., autonomous operations)
 - High-reliability applications
- Technology:
 - RHBD DARE65T std. cell
 - TSMC40LP std. cell
 - 7 nm radiation-hardened (in progress)
- Heritage of:
 - Three ongoing ESA supported projects
 - **OPERAND** project supported by Italian Ministry of Education and Research

OSIP I-2021-03237 Innovative Coarse-Grained Reconfigurable Array Platform for Computing Artificial Intelligence On-Board OSIP I-2022-04765 Risc-V Based SoC Featuring A Soft-GPU Hardware Accelerator for Artificial Intelligence On-Board OSIP I-2023-09415 UDSM Al-engine For Reliable, Energy-efficient Next-generation Satellites

SEFUW 2025







ESA OSIP - UDSM AI-Engine for Next-Gen. Satellites

- AI Engine IP
 - CGRA core
 - Simplified programming environment
 - CGRA configurations
 - Scheduling firmware
- What we have
 - CGRA core
 - Software-controllable SmartDMAs
 - Tiny RISC-V with DMA extension for data management
- What we need to carry out
 - Reliability analysis
 - SoC integration in UDSM 7nm
 - Programming environment
 - Benchmarking with multiple CNNs

SEFUW 2025







HW Development – CGRA Architecture



L. Zulberti, M. Monopoli, P. Nannipieri and L. Fanucci, "Architectural Implications for Inference of Graph Neural Networks on CGRA-based Accelerators", doi: <u>10.1109/PRIME55000.2022.9816810</u>. L. Zulberti, M. Monopoli, P. Nannipieri, L. Fanucci and S. Moranti, "Highly Parameterised CGRA Architecture for Design Space Exploration of Machine Learning Applications Onboard Satellites", doi: <u>10.23919/EDHPC59100.2023.10396632</u>.

L. Zulberti, M. Monopoli, P. Nannipieri, S. Moranti, G. Thys and L. Fanucci, "Efficient Coarse-Grained Reconfigurable Array architecture for machine learning applications in space using DARE65T library platform", doi: 10.1016/j.micpro.2025.105142

SEFUW 2025



Four hierarchy levels:

- Functional Units: ADD, MUL, SHF, BTW, LUT, MUX
- Processing Elements
 - Tiles
 - CGRA

High parametrization and scalability

- Number of rows and columns
- Implemented FUs and connection matrix
- FIFOs depth and pipeline stages
- Design optimized for ML execution
 - MACC operations
 - Non-linear activation functions

Matteo Monopoli and Gabriela Mystkowska



13

HW Development – SmartDMA Architecture

- RISC-V CPU
 - NeoRV32 CPU
 - Xdma custom extension
- Data Mover Engines
 - Support for 2D/3D data movement
 - Unaligned access
- Event Switch
 - Configurable DME dependencies
 - Data is loaded to the CGRA after transfer to local memory is finished

Nolting, S., & All the Awesome Contributors. (2024). The NEORV32 RISC-V Processor (v1.10.2). Zenodo. doi: <u>10.5281/zenodo.13127811</u>. L. Zulberti, A. Monorchio, M. Monopoli, G. Mystkowska, P. Nannipieri and L. Fanucci, "SmartDMA: Adaptable Memory Access Controller for CGRA-based Processing Systems", doi: <u>10.1109/DSD64264.2024.00048.</u>

SEFUW 2025





Matteo Monopoli and Gabriela Mystkowska 14



14

HW Development – CGR-AI Engine

- Scalable and parametrizable architecture:
 - CGRA size (rows, columns)
 - Size of banked local memories (for data and firmware)
 - Number of reading and writing DMEs
- External host loads firmware and kernel parameters
- **RISC-V** loads CGRA configuration from banked TCM
- RISC-V executes only data movement operations to instruct DMEs
- Event Switch module synchronizes high-speed DMEs communicating with the host with CGRA low-speed DMEs

SEFUW 2025











Ongoing Work – Reliability Enhancement Approach

- hardened technology
- Limited area and power consumption budget
- Evaluation of possible fault detection and mitigation strategies (TMR, EDAC)
 - Trade-offs between PPA metrics and robustness
 - Different level of granularity to consider (from register level to engine level) -
- activity

Criticality	Class 1	Class 2	Class 3	Class 4	Class	Operational Element	Memory
PA			Class J	C1055 +	C1	Partial Duplication	Dority Dit
Class 1	C2	C3	C4	C4	CI	w/ Self Check	Failty Dit
Class 2	C2	C3	C3	C4	C2	Partial Distributed TMR	DED Hamming Code
Class 3	C1	C2	C3	C3	C3	Block TMR w/ 1 Voter	SEC-DED Hamming Code
Class 4	C1	C1	C2	C2	C4	Block TMR w/ 3 Voters	TMR

M. Monopoli, M. Biondi, P. Nannipieri, S. Moranti, C. Bernardeschi and L. Fanucci, "Enhancing a Soft GPU IP Reliability Against SEUs in Space: Modelling Approach and Criticality Analysis on a Radiation-Tolerant FPGA", doi: 10.36227/techrxiv.173273657.76155928/v1

SEFUW 2025



Strict reliability requirements for space applications, especially when implementing on non-radiation-

Tabular-based classification approach already developed by University of Pisa as result of OSIP







Ongoing Work – Reliability Validation Approach

- Use of GUI fault injection application develope internally
- Netlist-based approach to emulate different ty of fault on the FPGA implemented netlist
 - Stuck-at
 - Transient -
- Addition of LUT-based fault injection cells
- Integration in a SystemVerilog simulation enviro



M. Monopoli, M. Biondi, P. Nannipieri, S. Moranti and L. Fanucci, "RADSAFiE: A Netlist-Level Fault Injection User Interface Application for FPGA-Based Digital Systems", doi: 10.1109/ACCESS.2025.3539932

SEFUW 2025



d	Fault Injection Application - RADSAFiE Project OOO							
	File Edit View Window Help							
	$\Box \textcircled{\bar{eq}} \land $							
	Import Netlist c:/Users/user/Desktop/ Import	Netlist Hierarchy Window						
pes	Fault Injection Type Simple Random	✓ Alu mult add sub						
		regFile_0 regFile_1 regFile_2 ✓ Cu						
	Fault Injection Information:							
	Tool Vendor Xilinx ~							
	Netlist File Extension Verilog ~							
	Module Name Alu ~	cu_scheduler						
	Module Instance Alu ~							
	Injection Level Module Cell	Fault Injection Report Window						
	Module Output [22:0]D; ~	FaultID Module Value						
	Fault Type Transient Stuck-At	Fault_0 Alu 0x00						
nment	Injection Method O Truth Table Injection Method	Fault_1 Cu 0xff						
	Fault Value0x00							
	Clear All Inject Export	Remove Clear All Save						
	Log Window							
	[2024-07-05 21:51:54] - [INFO] - New fault injection requested [2024-07-05 21:51:55] - [INFO] - New fault injected successfully							

Matteo Monopoli and Gabriela Mystkowska



17

Outline

- Background a. Motivation b. Hardware for Al Onboard 2. State-of-the-Art 3. CGRA-based HW Accelerator a. ESA OSIP Context
 - **CGR-Al Engine** b.
 - Ongoing work

4. Std. Cell Synthesis & FPGA Prototype

Conclusion

SEFUW 2025











CGR-AI on DARE65T RHBD Library Platform

- Optimized for std. cell implementation
- Synthesized on the 65nm RHBD DARE library platform
- Performed DSE activity with different:
 - Data width
 - CGRA rows
 - Local memory size
 - High-speed ports
- Tested against classic CNN kernels (e.g., conv2D, ReLU)
- Collected area and power consumption data across increasing clock frequencies
 - Great scalability
 - Less than 0.5 W @ 300 MHz on average

SEFUW 2025



	Config	Element Size	Vector Size	Rows	Memory (kB)	HS Ports
m	А	8-bit	32-bit	8	512	2R/2W
	В	8-bit	64-bit	3	64	1R/1W
	С	8-bit	64-bit	6	256	2R/2W
	D	8-bit	128-bit	1	64	1R/1W
	E	16-bit	64-bit	3	64	1R/1W











CGR-Al on Xilinx Zynq Ultrascale+ MPSoC

- System functionality confirmed in simulation
- FPGA overlay as intermediate milestone to demonstrate the readiness of the system
- Scale up directly based on the numbers obtained through std. cell synthesis
- Xilinx Zynq Ultrascale+ ZCU104 MPSoC
 - PS ARM core as external host
 - 600 MHz maximum operating frequency for the CGR-AI Engine
 - Contained resource utilization and power consumption, can still be optimized for FPGA implementation
- IP approach used for ease of portability

SEFUW 2025





3x2 CGRA with 9MB of local memory







Outline

- Background a. Motivation b. Hardware for Al Onboard 2. State-of-the-Art 3. CGRA-based HW Accelerator a. ESA OSIP Context **CGR-Al Engine** b.
 - Ongoing work

4. Std. Cell Synthesis & FPGA Prototype

Conclusion 5.

SEFUW 2025















dall'Unione europea



A Parameterizable CGRA-based Engine for Machine Learning Acceleration On-board Satellites

Matteo Monopoli and Gabriela Mystkowska matteo.monopoli@phd.unipi.it, gabriela.mystkowska@phd.unipi.it

- **CGRA-core for calculation acceleration**
 - Heterogeneous computing
 - Scalable and flexible
 - Low reconfiguration time -
- Developed systematic approach to increase reliability and fault injection application to validate results
- CGR-AI Engine synthesized on DARE65T RHBD
- FPGA overlay as intermediate milestone to demonstrate the readiness of the system

This study received funding from the European Union - Next-GenerationEU - National Recovery and Resilience Plan (NRRP) - MISSION 4 COMPONENT 2, INVESTMENT N. 1.1, CALL PRIN 2022 PNRR D.D. 1409 14-09-2022 – (OPERAND - a recOnfigurable Platform & framEwork for Ai inference on the eDge) CUP N.153D23006070001, from the Italian Ministry of Education and Research (MUR) within the framework of the FoReLab project (Departments of Excellence), from the European Space Agency (ESA) within the framework of the Open Space Innovation Platform (OSIP) co-funded research under the ESA contract n. 4000144254 and from the Space It Up project funded by the Italian Space Agency (ASI) and the Ministry of University and Research (MUR) under contract n. 2024-5-E.0 - CUP n. 153D24000060005.



















