

Open-Source Hardware for Heterogeneous Computing with ESP and RISC-V

Luca P. Carloni





Spring 2022 RISC-V Week Paris - May 3, 2022

The Age of Heterogeneous Computing

• State-of-the-art SoC architectures integrate increasingly diverse sets of components

 o different CPUs, GPUs, hardware accelerators, memory hierarchies, I/O peripherals, sensors, reconfigurable engines, analog blocks...

- The migration towards heterogeneous SoC architectures will accelerate, across almost all computing domains
 - IoT devices, mobile devices, embedded systems, automotive electronics, avionics, data centers and even supercomputers
- The set of heterogeneous SoCs in production in any given year will be itself heterogeneous!

o no single SoC architecture will dominate all the markets!





Heterogeneity Increases Design Complexity

- Heterogeneous architectures produce higher energyefficient performance, but make more difficult the tasks of design, verification and programming
 - at design time, diminished regularity in the system structure, chip layout
 - at runtime, more complex hardware/software and management of shared resources
- With each SoC generation, the addition of new capabilities is increasingly limited by engineering effort and team sizes
- The biggest challenges are (and will increasingly be) found in the complexity of system integration



Open-Source Hardware (OSH)

- An opportunity to reenergize the innovation in the semiconductor and electronic design automation industries
- The OSH community is gaining momentum
 - many diverse contributions from both academia and industry
 - $_{\circ}$ multi-institution organizations
 - government programs



RISC-V: The Free and Open RISC

Instruction Set Architecture





• CHIPS • CHIPS



Image Sources: https://chipsalliance.org/ https://github.com/nvdla https://www.openhwgroup.org/ https://parallel.princeton.edu/openpiton/ https://pulp-platform.org/ https://riscv.org/



4

The Open Challenge of Open-Source Hardware

- To date, however, most OSH projects are focused on the development of individual SoC components, such as a processor core or an accelerator
- This leaves open a critical challenge:

How can we realize a complete SoC for a given target application domain by efficiently reusing and combining a variety of independently developed, heterogeneous, OSH components, especially if these components are designed by separate organizations for separate purposes?



The Concept of Platform

- Innovation in SoC architectures and their design methodologies is needed to promote design reuse and collaboration
 - Architectures and methodologies must be developed together
- Platform = architecture + methodology
 - An SoC architecture enables design reuse when it simplifies the integration of many components that are independently developed
 - An SoC methodology enables design collaboration when it allows designers to choose the preferred specification languages and design flows for the various components
- An effective combination of architecture and methodology is a platform that maximizes the potential of open-source hardware
 - by scaling-up the number of components that can be integrated in an SoC and by enhancing the productivity of the designers who develop and use them



ESP : An Open-Source Platform for SoC Design

🚺 Home Release Resources 🗸 News Press Team Contact

ESP

the open-source SoC platform



The ESP Vision

ESP is an open-source research platform for heterogeneous system-on-chip design that combines a scalable tile-based architecture and a flexible system-level design methodology.



ESP provides three accelerator flows: RTL, high-level synthesis (HLS), machine learning frameworks. All three design flows converge to the ESP automated SoC integration flow that generates the necessary hardware and software interfaces to rapidly enable full-system prototyping on FPGA.

Overview



©Luca Carloni

Architecture

www.esp.cs.columbia.edu

Latest Posts



Paper accepted at MICRO 2021

Our paper "Cohmeleon: Learning-Based Orchestration of Accelerator Coherence in Heterogeneous SoCs" will be presented in October at MICRO 2021.

Read more

Published: Sep 14, 2021





Dapar published in





Outline

The ESP Architecture





The ESP Methodology

the open-source SoC platform

his 4 ml 🛟

Design Flows

RTL

Design Flows

NO HES

ESP

The ESP Vision

00

CHISEL

handlicher

Research & Teaching with ESP



©Luca Carloni

SW Library

HW IP Library

0 7 8 0

Published Sep 11, 2020

Upcoming talk at

rofessor Carloni will give a

VLSISoC 2020

Latest Posts

ESP Architecture

- RISC-V Processors
- Many-Accelerator
- Distributed Memory
- Multi-Plane NoC

The ESP architecture implements a distributed system, which is scalable, modular and heterogeneous, giving processors and accelerators similar weight in the SoC





ESP Architecture: Processor Tile

- Processor off-the-shelf
 - RISC-V CVA6-Ariane (64 bit) SPARC V8 Leon3 (32 bit)
 RISC-V IBEX (32 bit)
 - L1 private cache
- L2 private cache
 - $_{\circ}\,$ Configurable size
 - $_{\circ}$ MESI protocol
- IO/IRQ channel
 - $_{\circ}$ Un-cached
 - Accelerator config. registers, interrupts, flush, UART, ...





ESP Architecture: Memory Tile

- External Memory Channel
- LLC and directory partition
 - $_{\odot}$ Configurable size
 - $_{\circ}$ Extended MESI protocol
 - Supports coherent-DMA for accelerators
- DMA channels
- IO/IRQ channel



ESP Architecture: Accelerator Tile

- Accelerator Socket w/ Platform Services
 - Direct-memory-access
 - Run-time selection of coherence model:
 - Fully coherent
 - LLC coherent
 - Non coherent
 - $_{\circ}$ User-defined registers
 - $_{\odot}$ Distributed interrupt



ESP Accelerator Socket





ESP Software Socket

• ESP accelerator API

- Generation of device driver and unit-test application
- $_{\circ}\,$ Seamless shared memory



```
/*
* Example of existing C application with ESP
* accelerators that replace software kernels 2, 3,
* and 5. The cfg k# contains buffer and the
* accelerator configuration.
* /
 int *buffer = esp alloc(size);
 for (...) {
   kernel 1(buffer,...); /* existing software */
   esp run(cfg k2); /* run accelerator(s) */
   esp run(cfg k3);
   kernel 4(buffer,...); /* existing software */
   esp run(cfg k5);
 validate(buffer); /* existing checks */
              /* memory free */
 esp free();
```



ESP Platform Services

Accelerator tile DMA Reconfigurable coherence Point-to-point ESP or AXI interface DVFS controller	Processor Tile Coherence I/O and un-cached memory Distributed interrupts DVFS controller
Miscellaneous Tile	Memory Tile
Debug interface	Independent DDR Channel
Performance counters access Coherent DMA	LLC Slice
Shared peripherals (UART, ETH,)	DMA Handler



Outline

The ESP Architecture



Latest Posts

Upcoming talk at

Published Sep 11, 2020

VLSISoC 2020



The ESP Methodology

the open-source SoC platform

his 4 ml 🛟

Design Flows

RTL

Design Flows

NO HES

ESP

The ESP Vision

00

CHISEL

harderty

Research & Teaching with ESP



©Luca Carloni

SW Library

HW IP Library

0 7 8 0

The Pillars of the ESP Approach

- Develop platforms, not just architectures
 - A platform combines an architecture and a companion design methodology
- Move from a processor-centric to an SoC-centric perspective
 - $_{\odot}$ The processor core is just one component among many others
- Raise the level of abstraction
 - Move from RTL design to domain-specific system-level design with high-level synthesis...
 - $_{\odot}$...but keep supporting different abstraction levels and design flows
- Promote Open-Source Hardware
 - $_{\odot}$ Build libraries of reusable components
 - $_{\rm O}$ Support the integration of third-party IP components



The ESP Vision: Domain Experts Can Design SoCs



ESP Methodology In Practice

automated interactive manual (opt.) manual

Accelerator Flow

Generate accelerator

Specialize accelerator * this step is automated * for ML applications

Test behavior

Generate RTL

Test RTL

Optimize RTL





ESP Accelerator Flow

Developers focus on the high-level specification, decoupled from memory access, system communication, hardware/software interface



20

Retrospective: Latency-Insensitive Design

[Carloni et al. '99]





Relay Shell Relay Shell Station Station •∏⊥ 14 Core Core ۲J ٦ł Core Core Logic Logic CTRL CTRL dataIn dataOut main mu FF aux $\mathbf{T}\mathbf{T}$ sel mainEn auxEn stopOut stopIn control mainEn voidIn voidOut void mu 1(1) 0

Latency-Insensitive Design

- is the foundation for the *flexible synthesizable RTL representation*
- anticipates the separation of computation from communication that is proper of TLM with SystemC
 - through the introduction of the Protocols & Shell paradigm



ESP Interactive Flow for SoC Integration









Outline

The ESP Architecture



Latest Posts

Upcoming talk at

Published Sep 11, 2020

VLSISoC 2020



The ESP Methodology

the open-source SoC platform

his 4 ml 🛟

Design Flows

RTL

Design Flows

No HLS Ca

ESP

The ESP Vision

00

CHISEL

harderty

Research & Teaching with ESP



©Luca Carloni

SW Library

HW IP Library

0 7 8 0

Example of a System We Built: FPGA Prototype to Accelerate Wide-Area Motion Imagery



- **Design:** Complete design of WAMI-App running on an FPGA implementation of an ESP architecture
 - featuring 1 embedded processor,
 12 accelerators, 1 five-plane NoC,
 and 2 DRAM controllers
 - SW application running on top of Linux while leveraging multithreading library to program the accelerators and control their concurrent, pipelined execution
 - Five-plane, 2D-mesh NoC efficiently supports multiple independent frequency domains and a variety of platform services

[P. Mantovani , L. P. Carloni et al., *An FPGA-Based* Infrastructure for Fine-Grained DVFS Analysis in High-Performance Embedded Systems, DAC 2016⁴]

Seamless Integration of NVDLA Accelerators

- New design flow of general applicability to integrate third-party accelerators
 - demonstrated w/ NVIDIA Deep Learning Accelerator (NVDLA)
- Transparent accelerator integration
 - original software apps can run "as is"
- Linear performance scalability
 - when scaling up NVDLA instances with DDR channels





[D. Giri et al. "Ariane + NVDLA: Seamless Third-Party IP Integration with ESP", CARRV'20]







©Luca Carloni

IN THE CITY OF NEW YORK

ESP4ML

Open-source design flow to build and program SoCs for ML applications

- Combines and hls 4 ml
 - ESP is a platform for heterogeneous SoC design
 - hls4ml automatically generates accelerators from ML models
- Main contributions to ESP:
 - Automated integration of hls4ml accelerators
 - $_{\odot}\,$ Accelerator-accelerator communication
 - $_{\odot}\,$ Accelerator invocation API

[D. Giri, K.-L. Chiu, G. Di Guglielmo, P. Mantovani, and L. P. Carloni. "ESP4ML: Platform-Based Design of Systems-on-Chip for Embedded Machine Learning", DATE '20]





hls4ml

- Open-source tool developed by Fast ML Lab
- Translates ML algorithms into accelerator specifications that are synthesizable with highlevel synthesis tools for both **FPGA and ASIC** implementations
- Born for high-energy physics (small and ultra-low latency networks), it is gaining broad applicability and a growing community of contributors and users

hls4ml: An Open-Source Codesign Workflow to Empower Scientific Low-Power Machine Learning Devices

Farah Fahim* Luca P. Carloni Philip Harris Benjamin Hawks Giuseppe Di Guglielmo Jeffrey Krupa Columbia University Christian Herwig Dvlan Rankin New York, NY, USA James Hirschauer MIT Sergo Jindariani Cambridge, MA, USA Nhan Tran* Fermilab Batavia, IL, USA Manuel Blanco Valentin Thea Aarrestad Javier Duarte Iosiah Hester UC San Diego Hamza Javed La Jolla, CA, USA Yingvi Luo Vladimir Loncar iduarte@ucsd.edu John Mamish Maurizio Pierini Seda Orgrenci-Memik Adrian Alan Pol Northwestern University Sioni Summers Evanston, IL, USA European Organization for Nuclear Research (CERN) Geneva, Switzerland Scott Hauck Jennifer Ngadiuba Mia Liu Caltech Purdue University Shih-Chieh Hsu Pasadena, CA, USA West Lafavette, IN, USA University of Washington Seattle, WA, USA Duc Hoang Edward Kreinar Zhenbin Wu Rhodes College HawkEve360 University of Illinois at Chicago Memphis, TN, USA Herndon, VA, USA Chicago, IL, USA ABSTRACT Accessible machine learning algorithms, software, and diagnostic tools for energy-efficient devices and systems are extremely valuable across a broad range of application domains. In scientific domains, real-time near-sensor processing can drastically improve experimental design and accelerate scientific discoveries. To sup-

port domain scientists, we have developed h1s4n1, an open-source

software-hardware codesign workflow to interpret and translate

machine learning algorithms for implementation with both FPGA

and ASIC technologies. In this paper, we describe the essential

features of the h1s4nl workflow including network optimization

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed

for profit or commercial advantage and that copies hear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored.

'Also affiliated with Northwortern University

For all other uses, contact the owner/author(s).

@ 2023 Copyright held by the owner/wathor(s).

ThyAL Research Symposium 21, March 2021, San Jose, CA

techniques-such as pruning and quantization-aware trainingwhich can be incorporated naturally into the device implementations. We expand on previous hls4nl work by extending capabilities and techniques towards low-power implementations and increased usability: new Python APIs, quantization-aware pruning, end-to-end FPGA workflows, long pipeline kernels for low power, and new device backends include an ASIC workflow. Taken together, these and continued efforts in h1s4n1 will arm a new generation of domain scientists with accessible, efficient, and powerful tools for machine-learning-accelerated discovery.

KEYWORDS

hls4ml, machine learning, neural networks, tinyML, FPGA, ASIC, low-power, low-latency

ACM Reference Formal:

Farah Fahim, Benjamin Hawks, Christian Herwig, James Hirschause, Sergo Jandariami, Nhan Tran, Luca P. Carloni, Giuseppe Di Guglielmo, Philip Harris, Jeffrey Krupa, Dylan Rankin, Manuel Blanco Valentin, Josiah Hester, Yingyi Luo, John Mamith, Seda Orgrenci-Memik, Thea Aarrestad, Hama Javed, Vladimir Loncar, Maurizio Pierini, Adrian Alan Pol. Sioni Summers, Javier Duarte, Scott Hauck, Shih-Chieh Hau, Jenzañer Ngadiuba, Mia Liu, Duc



Cohmeleon: Learning-Based Orchestration of Accelerator Coherence in Heterogeneous SoCs

- Accelerator performance can vary greatly based on coherence modes
 - SoCs should support multiple coherence modes for optimal performance
- Reinforcement learning can be used to automatically manage coherence mode decisions
- With little overhead, Cohmeleon provides significant performance benefits for multiple objectives



[J. Zuckerman, D. Giri, J. Kwon, P. Mantovani, and L. P. Carloni Cohmeleon: Learning-Based Orchestration of Accelerator Coherence in Heterogeneous SoCs IEEE/ACM International Symposium on Microarchitecture (MICRO-54), 2021.] E S P ©Luca Carloni

Teaching with ESP

https://www.esp.cs.columbia.edu

🚺 Home Release Resources 🗸 News Press Team Contact

Teaching

Class projects

1) Design and integration of an accelerator with ESP

For this project each student will use ESP to design one or more accelerators and to integrate them in a system-on-chip (SoC), capable of booting Linux. Then the student will evaluate the SoC both with RTL simulation and on FPGA.

To get a more practical sense of the project, you should familiarize yourself with ESP by using the resources on this website. Specifically:

- Check out the ESP website Homepage including the short introductory video.
- Watch the 16 minutes overview video in the Documentation section.
- Watch the videos and read the guides of the relevant hands-on tutorials available in the Documentation section. Especially relevant are the "How to: setup", "How to: design a single-core SoC" guides and the "How to: design an accelerator in ... "guide that applies to your specific project.
- Explore the rest of the website to get the full picture of the ESP project.

Accelerator flows

For your project proposal, you are asked to choose which design flow you want to use to build your accelerator.

ESP offers multiple accelerator design flows: Stratus HLS flow (accelerator designed in SystemC), the Vivado HLS flow (accelerator designed in C/C++), the Catapult HLS flow (accelerator designed in C/C++) and the hIs4ml flow (accelerators designed in Keras/Pytorch/ONNX).

Other options include designing the accelerator in RTL (Verilog, VHDL, SystemVerilog, Chisel). These other options do not have full support and documentation yet. It is possible to use them, but they will require a bigger integration effort.



CSEE-4868: System-on-Chip Platforms

- Foundation course on the programming, design, and validation of SoCs with emphasis on high-performance embedded applications
- Offered at Columbia since 2011, moved to upper-level curriculum in Fall 2016
 - required course for CE BS program, elective for MS programs in CS and EE
- Course Goals
 - mastering the HW and SW aspects of integrating heterogeneous components into a complete system
 - designing new components that are reusable across different systems, product generations, and implementation platforms
 - evaluating designs in a multi-objective optimization space
 - [L. P. Carloni et al. Teaching Heterogeneous Computing with System-Level Design Methods, WCAE 2019]



Upcoming Events: The First OSCAR Workshop!

Open-Source Computer Architecture Research (OSCAR)

Saturday, June 18, 2022 - New York City (co-located with ISCA 2022)



https://oscar-workshop.github.io/

Welcome to OSCAR 2022!

OSCAR 2022 is the first edition of a new workshop on open-source hardware which addresses the wide variety of challenges encountered by both hardware and software engineers in dealing with the increasing heterogeneity of next-generation computer architectures. By providing a venue which brings together researchers from academia, industry and government labs, OSCAR promotes a collaborative approach to foster the efforts of the open-source hardware community in this direction.

©Luca Carloni

In Summary: ESP for Open-Source Hardware

- We contribute ESP to the OSH community in order to support the realization of
 - more scalable architectures for SoCs that integrate
 - more heterogeneous components, thanks to a
 - more flexible design methodology, which accommodates different specification languages and design flows
- ESP was conceived as a heterogeneous integration platform from the start and tested through years of teaching at **Columbia University**
- We invite you to use ESP for your projects and to contribute to ESP!

Home Release Resources 🗸 News Press Team Contact

https://www.esp.cs.columbia.edu

the open-source SoC platform

C Latest Posts

The ESP Vision

Overview

Architecture

NoC-based

 Available processors 64-bit Ariane (RISC-V)

32-bit Leon3 (Sparc)

32-bit Ibex (RISC-VI)

ESP accelerator design flows

· Accelerator design and integration flows

C/C++ with Xilinx Vivado HLS

 SystemC with Cadence Stratus HLS C/C++ with Mentor Catapult HLS

ESP

ESP is an open-source research platform for heterogeneous system-on-chip design that combines a scalable tile-based architecture and a flexible system-level design methodology



ESP provides three accelerator flows: RTL, high-level synthesis (HLS), machine learning frameworks. All three

software interfaces to rapidly enable full-system prototyping on FPGA

Lource

Tile-based: processor, accelerator, memory, scratchpad and I/O tiles

SoC Platform

design flows converge to the ESP automated SoC integration flow that generates the necessary hardware and

SoC	Rully coherant	Collegent Data
CPU CPU Etcathe	accelerator Social mans	coordination factor many
12 coche: 12 coche:	cathe	EURA
(bes, NoC)		_
Iter, HoC)	DMA	E.m.
(ben, NucC) (ben, NucC) (ben, NucC) (ben, NucC) (ben, NucC) (c) (c) (c) (c) (c) (c) (c) (c	Essan Social main acceler citor	Const Second second

Paper accepted at MICRO 2021

Our paper "Cohmeleon: Learning-Based Orchestration of Accelerator Heterogeneous SoCs" will be presented in October at MICRO 2021.

Read more

Published: Sep 14, 2021



Paper published in the IEEE Micro special issue on **FPGA** Computing

Our paper "Accelerator Integration for Open-Source SoC Design" has been published in the IEEE Micro special issue on FPGA Computing (Jul-Aug 2021)

IN THE CITY OF NEW YORK

Y

32

Read more



Thank you from the **ESP** team!

https://esp.cs.columbia.edu

https://github.com/sld-columbia/esp



