

OSCAR@ISCA 2026

Early Stage SoC Architectural Concept Definition: *Open-Source Toolset from the EPOCHS Project*

June 28, 2026



Pradip Bose

IBM T. J. Watson Research Center

{pbose}@us.ibm.com

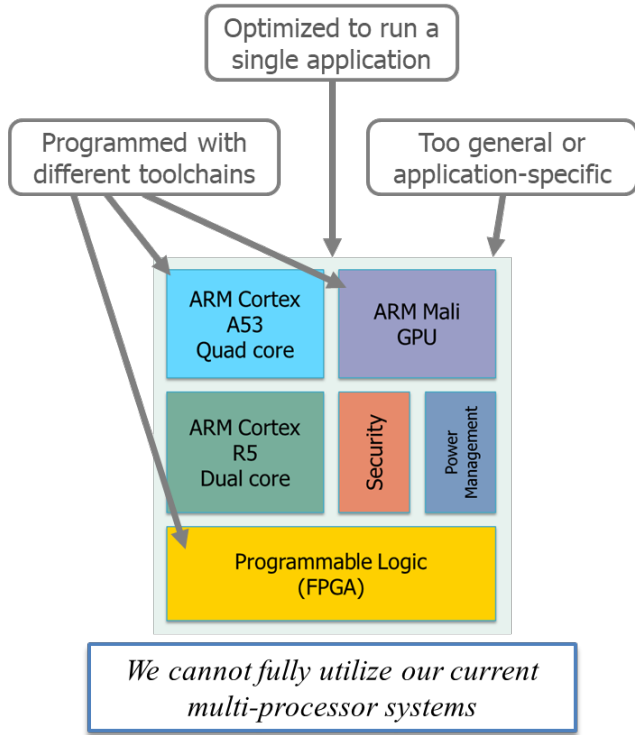
Acknowledgement

- Thanks to the many IBM colleagues who contribute to and support different aspects of this work + our esteemed university collaborators at **Harvard**, **Columbia**, and **UIUC** (Profs. David Brooks, Vijay Janapa Reddi, Gu-Yeon Wei, Luca Carloni, Ken Shepard, Sarita Adve, Vikram Adve, Sasa Misailovic) + many brilliant graduate students and postdocs!
- Special thanks to **Dr. Thomas Kazior** and **Dr. Thomas Rondeau**, final and first Program Managers of the DARPA MTO DSSoC Program, respectively.
- Also, thanks to John Wohlbier et al. at CMU-SEI, Jesse Jimenez, Gena Osborn et al. at C5ISR, John Marsh, Chris Earl et al. at DARPA for their constructive feedback and support.

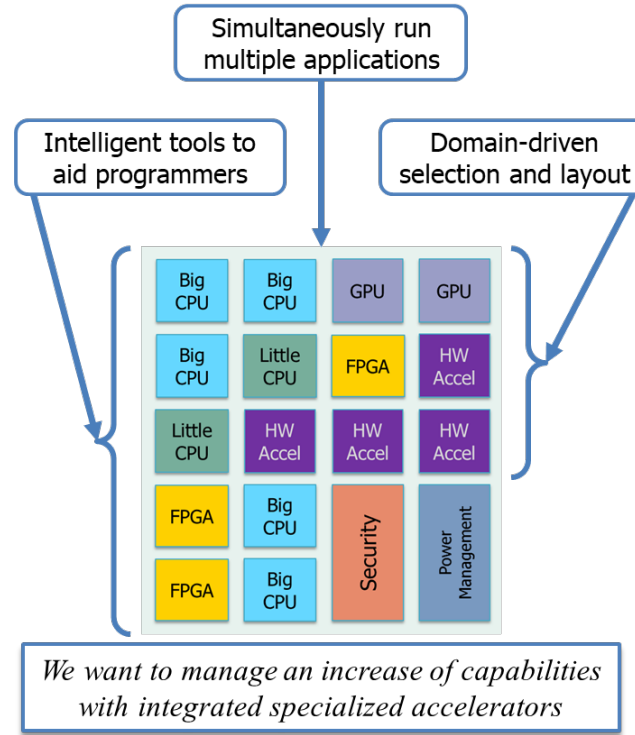
This research was developed with funding from the Defense Advanced Research Projects Agency (DARPA). The views, opinions and/or findings expressed are those of the authors and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government.

The DSSoC (Domain-Specific SoC) Vision

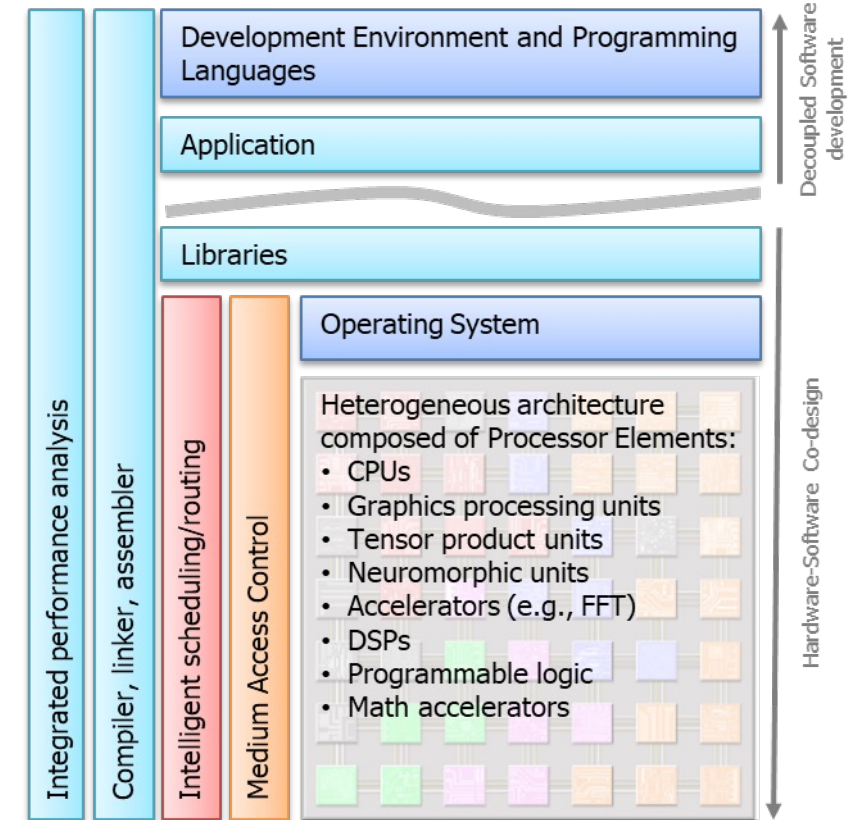
Status Quo (2018)



DSSoC End Goal (2023)



DSSoC's Full-Stack Integration

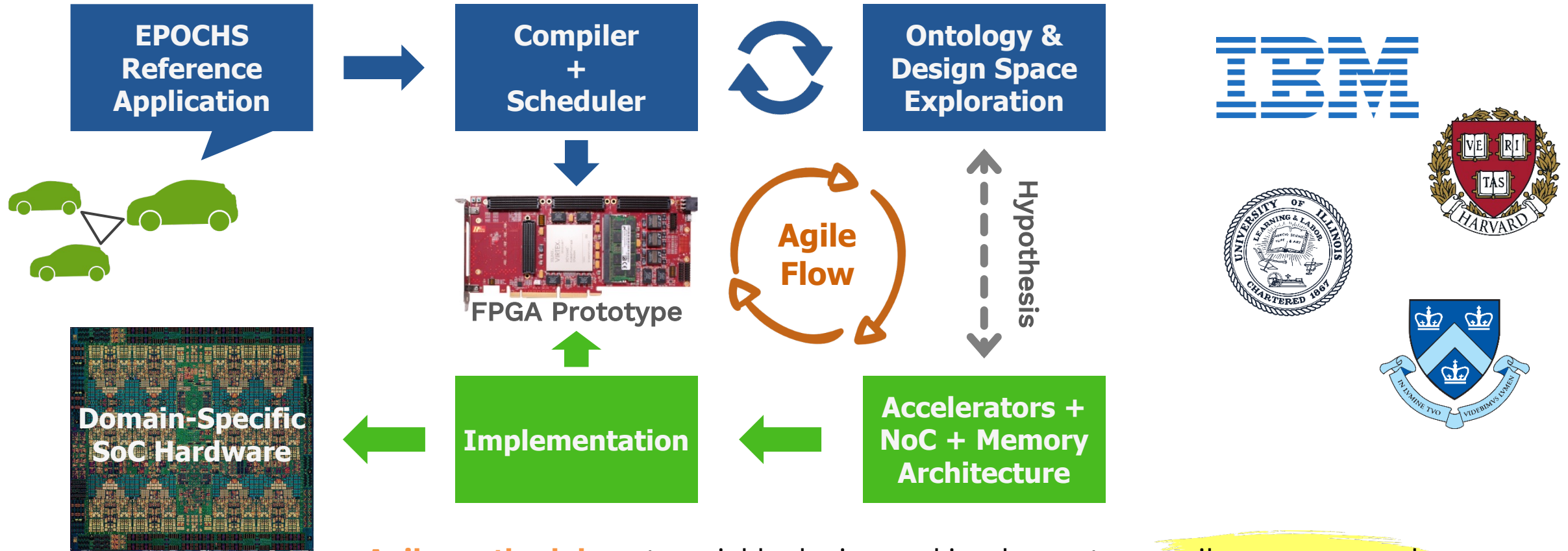


- Software-Hardware Co-Design [Efficiency]
- Agile SoC Design [Design Productivity]
- Easily Programmable [User Friendly]

➔ The three driving principles

Major slide content: courtesy Tom Rondeau (DARPA Prog. Mgr)

EPOCHS Agile Flow Methodology

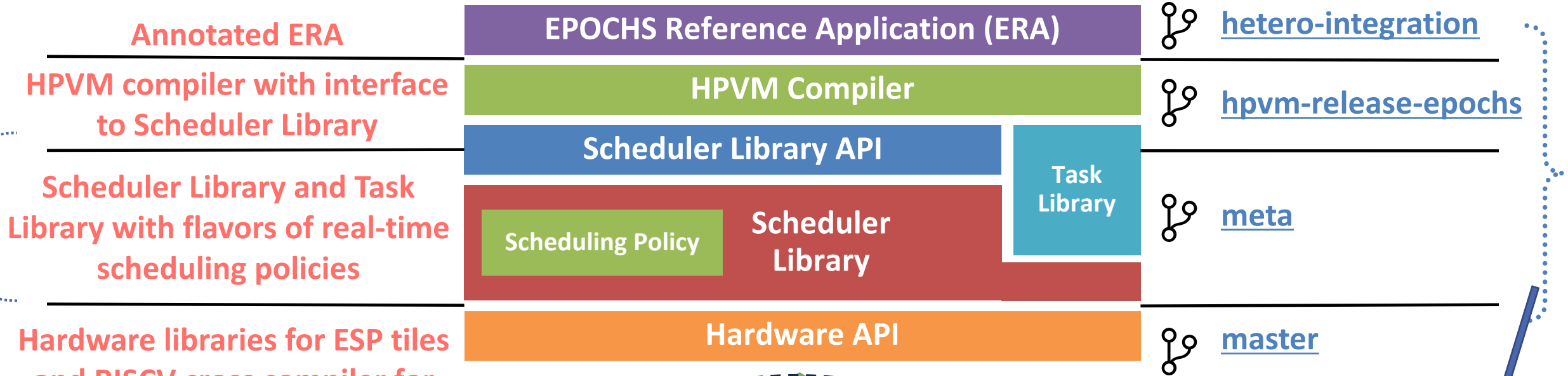


Agile methodology to quickly design and implement an easily programmed domain-specific SoC for real-time cognitive decision engines in connected vehicles

Putting It All Together: EPOCHS Appliance



Programmable! Programmable! Programmable!



Annotated ERA

HPVM compiler with interface to Scheduler Library

Scheduler Library and Task Library with flavors of real-time scheduling policies

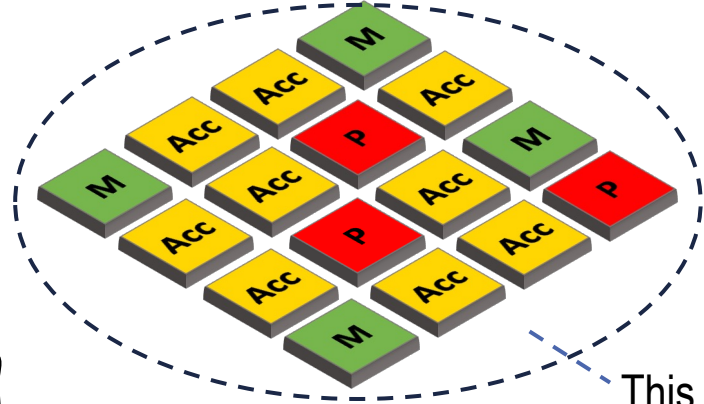
Hardware libraries for ESP tiles and RISC-V cross compiler for the RISC-V Ariane core

SoC-agnostic agile system programmability (open-source tools)



“Smart” Scheduler

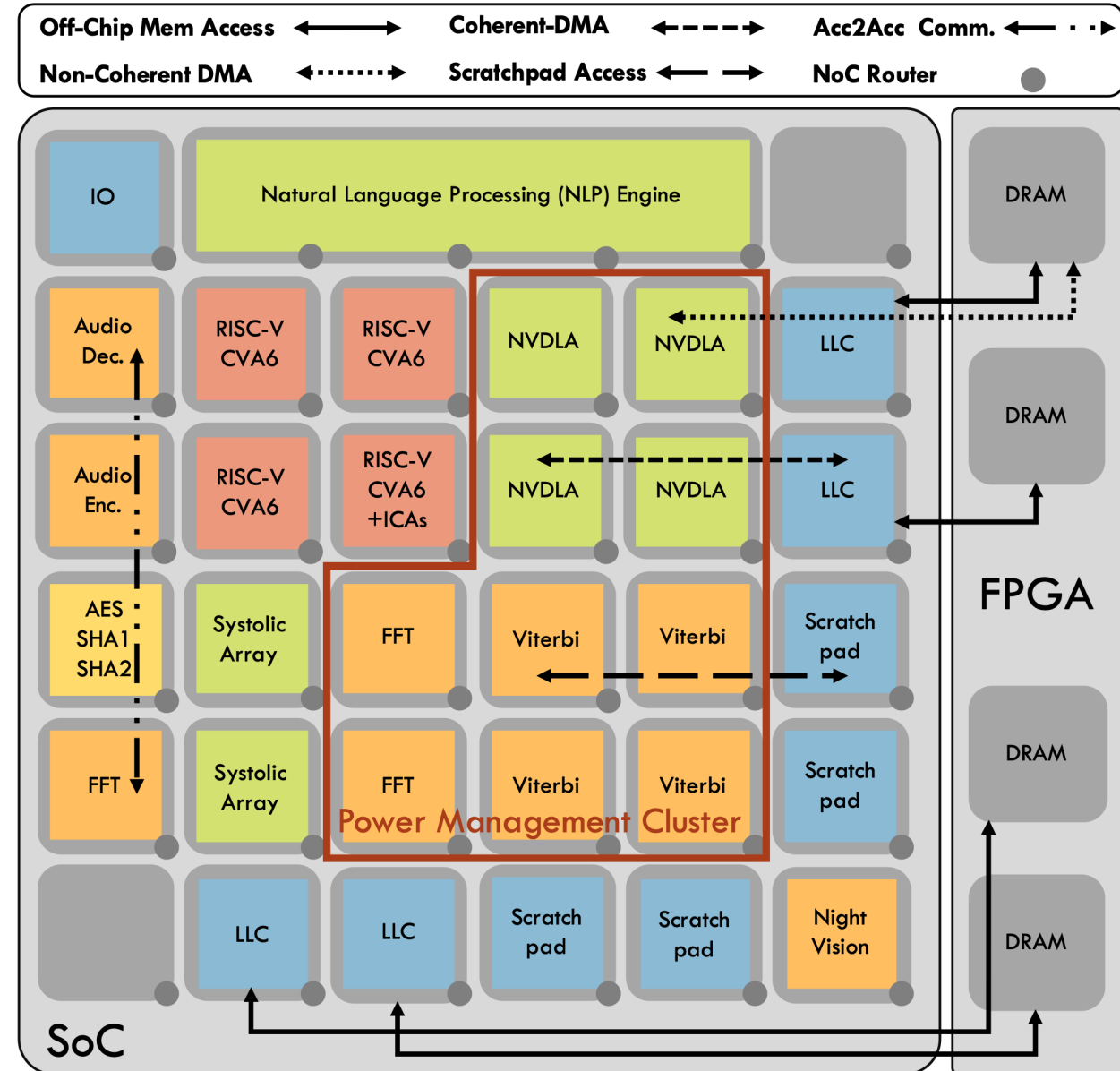
github.com/IBM/scheduler-library



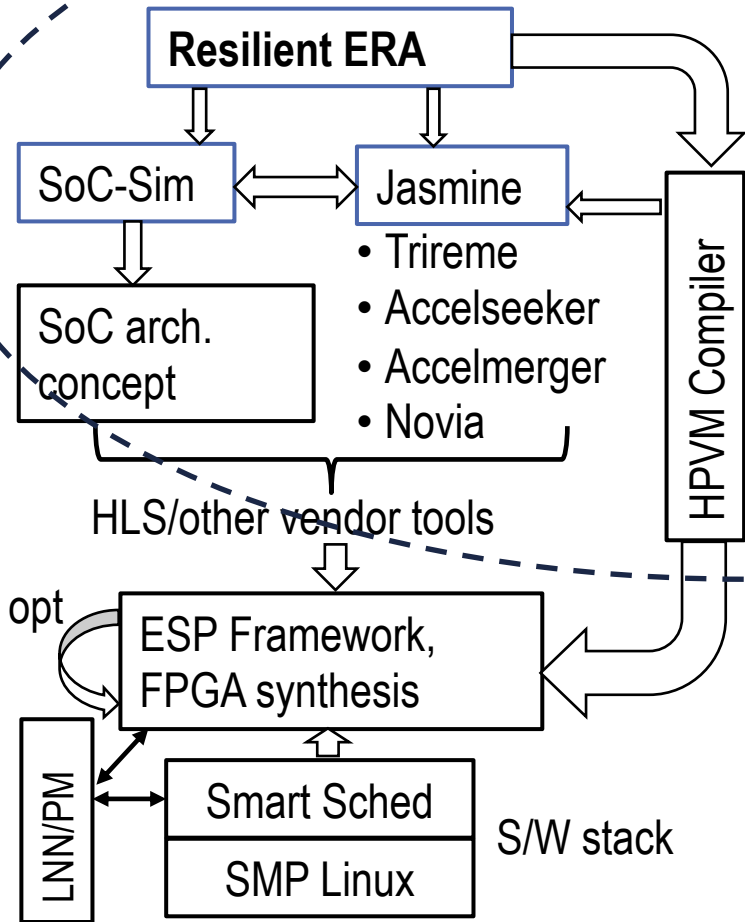
This is a cartoon to help highlight the s/w stack!
The next slide shows the real hardware SoC

EPOCHS-1 Chip Highlights

- **64** mm² SoC designed in **12** nm FinFET
- **35** clock domains; **23** power domains
- **8.4** MB on-chip SRAM memory
- Tile-based SoC architecture
- **34** tiles connected by a **6**-plane 2-D mesh NoC
- The **74** Tbps NoC provides flexible orchestration of data
- **23** accelerators of **14** different types
- **10** accelerators compose a cluster demonstrating a **novel distributed hardware power management scheme**
- Designed by a small team of PhD students, postdocs, and industry researchers in **3** months with **ESP**, an open-source platform for agile SoC design



Documented Toolchain (Open-Source)



- [SoC-Sim \(ARTEMIS\)](#) development, built around the [FARSI](#) tool from Harvard and [STOMP](#) simulator from IBM

- **Jasmine ontology toolset** – open-source, documented

- ✓ [Trireme](#)
- ✓ [Accelseeker](#)
- ✓ [Accelmerger](#)
- ✓ [NOVIA](#)

Focus, in this talk, on this pre-RTL SoC architecture definition process

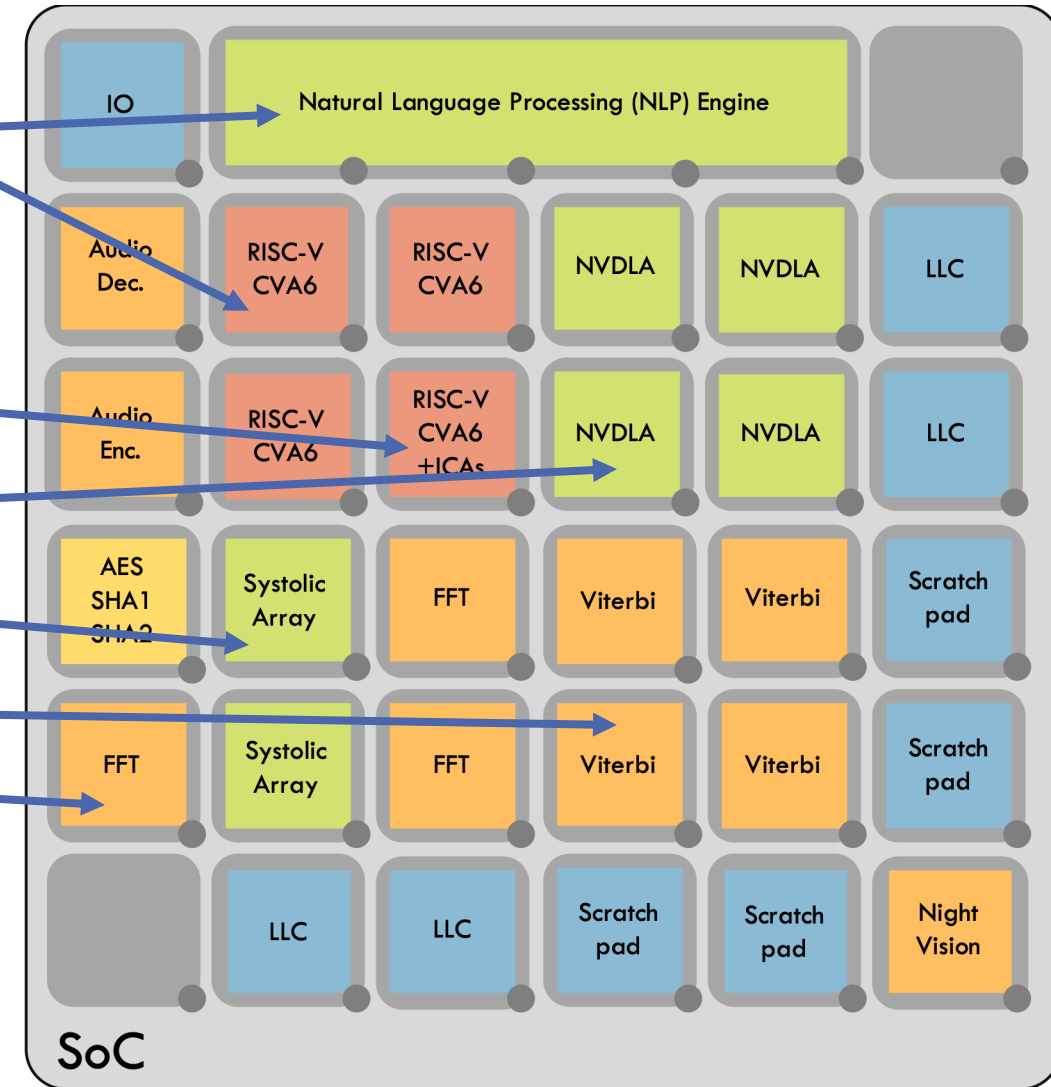
This is the aspect of design that we don't hear about much. Yet, this is where computer architects must deliver so that there are no post-silicon shortfalls in power-performance or other metrics!

- Software stack

- ✓ [Scheduler library \(SL\)](#)
- ✓ [HPVM compiler integration](#)
- ✓ [SMP Linux bring-up](#)
- ✓ [ESP](#) from Columbia Univ. (Prof. Luca Carloni's group).

Questions that may arise

- Why are there **four** RISC-V (CVA-6) Ariane processor cores?
- And, why this 4-tile spanning NLP engine?
- Why is there **one** core with inline compute accelerators (ICAs)?
- **Four** NVDLA engines, **plus two** systolic array PEs ?
- **Four** Viterbi and **three** FFT accelerators?
- And, so on, including questions about the scratch pad and LLC blocks



Honest answer: our concept-phase pre-RTL microarchitecture definition toolset was not ready in time! Our choices were *ad hoc*!

But that would not an acceptable answer for a real product chip!

- Every piece of the silicon real-estate must be justified in light of target workloads and metrics!
- So, if one has to design and deploy an edge-AI heterogeneous SoC with tight SWaP-C and real-time performance requirements, one must invest into the pre-RTL SoC microarchitecture definition toolset in advance and use it to ensure a competitive, performant product.
- Therefore, as we attempt to design **future** domain-specific SoCs (and SiPs), it is worth examining what we were able to develop in terms of the pre-RTL modeling & microarchitecture optimization tools.

Early-Stage SoC Architecture Definition

Two Key Challenges

a) How does one systematically “discover” the “right” accelerator types and counts?

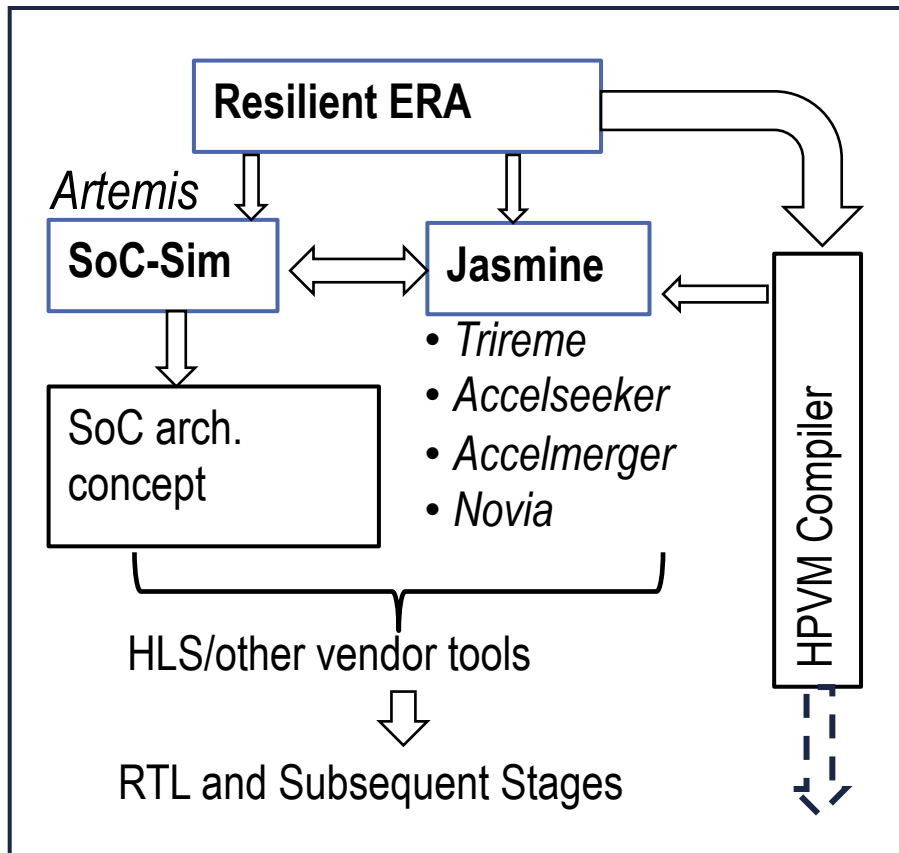
The Jasmine toolset; the tools listed are:

- ✓ [Trireme](#)
- ✓ [Accelseeker](#)
- ✓ [Accelmerger](#)
- ✓ [NOVIA](#)

And, once the desired h/w primitives are known,

b) How does one iteratively “optimize” the placement of CPUs and accelerator PE tiles in order to maximize a specified performance or power-performance metric?

SoC-Sim (aka [Artemis](#))



Heterogeneous Parallel Virtual Machine (HPVM) – Compiler Framework

Hashim Sharif, Yifan Zhao, Adel Ejeh,
Akash Kothari, A. Rafae Noor, Leon Medvinsky,
Vikram Adve, Sasa Misailovic, Sarita Adve (UIUC)



Goal: Programmability for Heterogeneous Parallel Systems

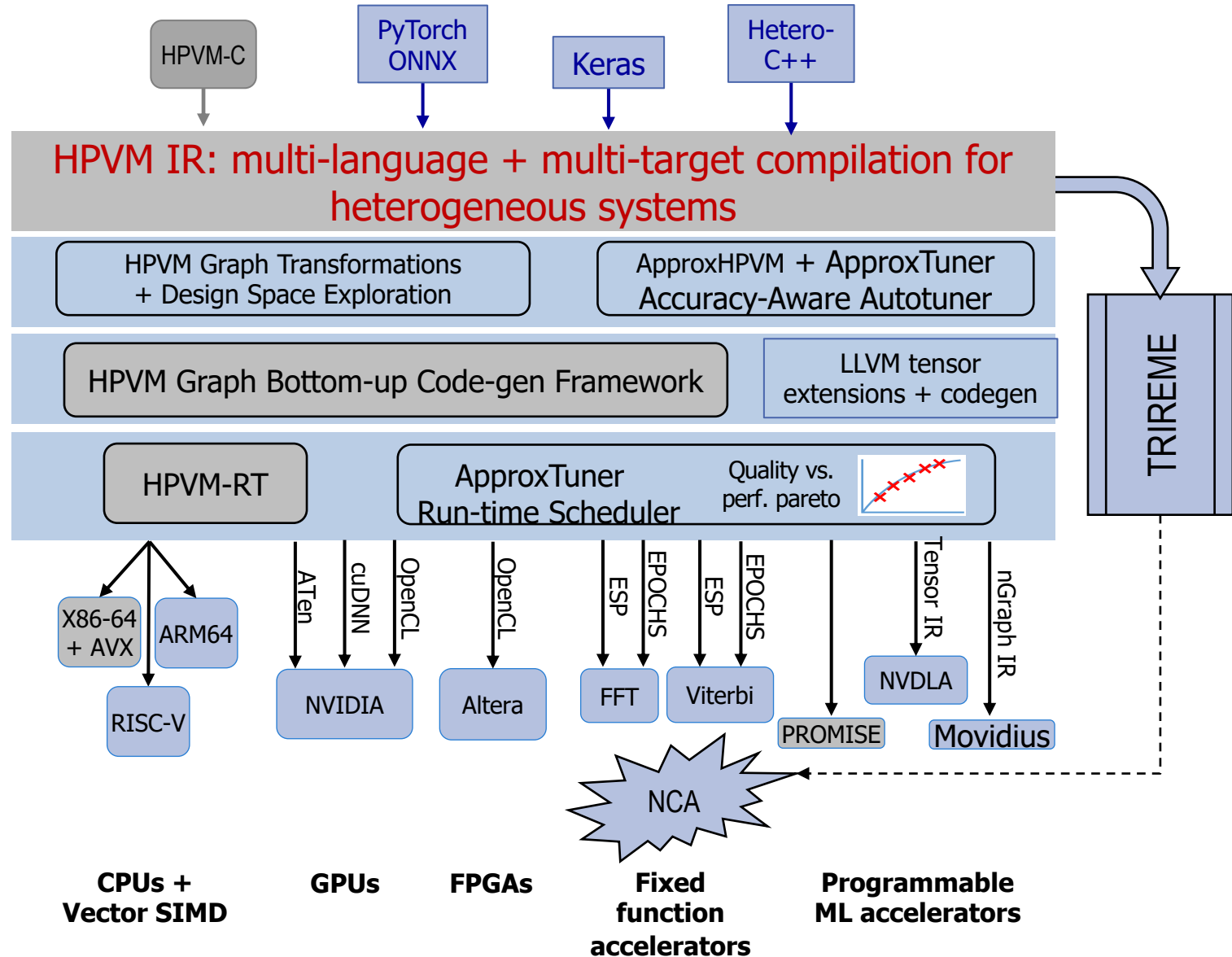
HPVM enables:

1. Multiple parallel languages, DSLs
2. Retargetable parallel compiler IR, optimizations, codegen
3. Large perf + energy gains for edge ML workloads
4. Automated hardware DSE

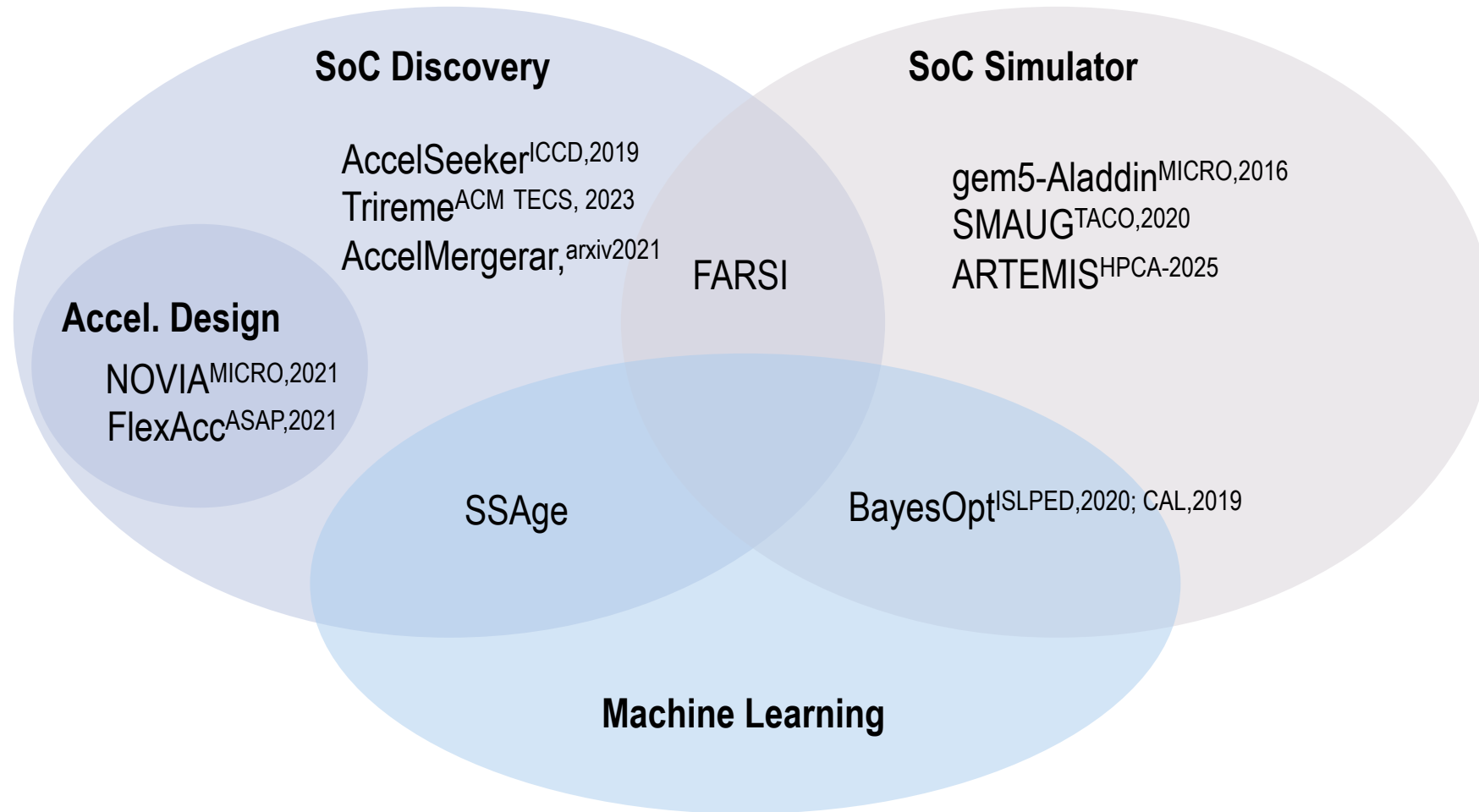
Papers
Kotsifakou+, PPOPP 2018
Srivastava+, ISCA 2018
Kang+, MICRO 2019
Sharif+, OOPSLA 2019
Sharif+, PPOPP 2021
Sharif+, in review, 2021
Zacharapolous+, in review. 2021

Software
Jan. 2020: HPVM 0.5
April 2021: HPVM 1.0

<https://publish.illinois.edu/hpvm-project/>

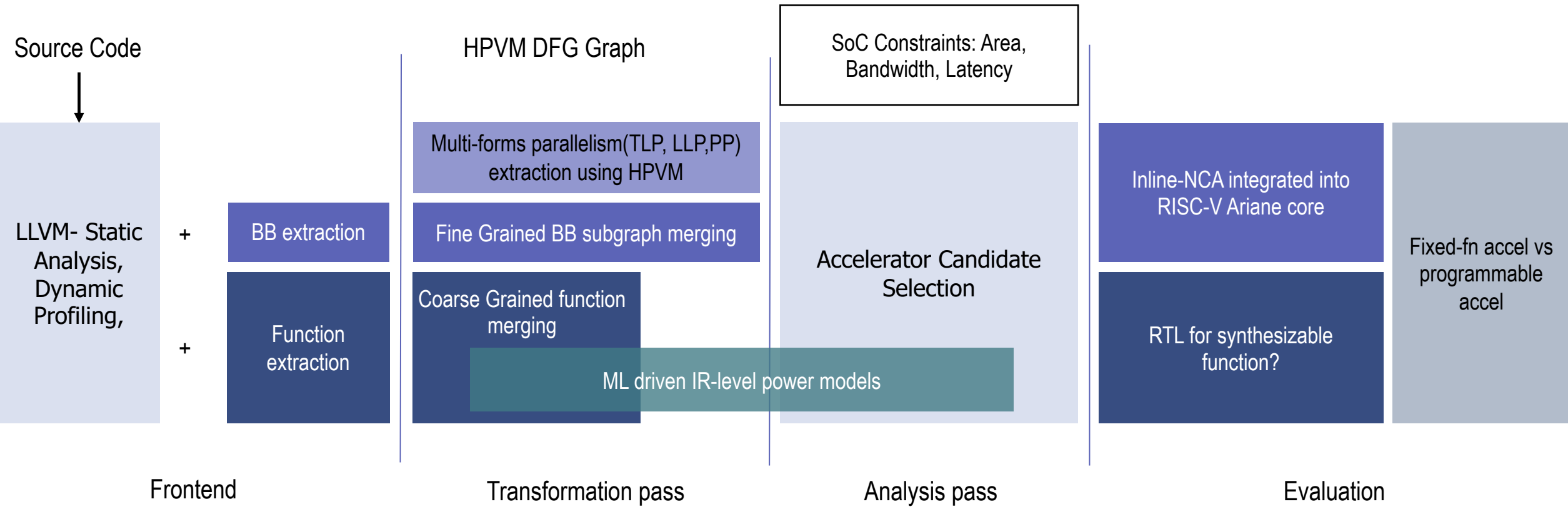


Landscape of Ontology Tools: *Jasmine*



We actually attempted a whole lot more than what finally proved to be appropriate for the problem at hand 😊!

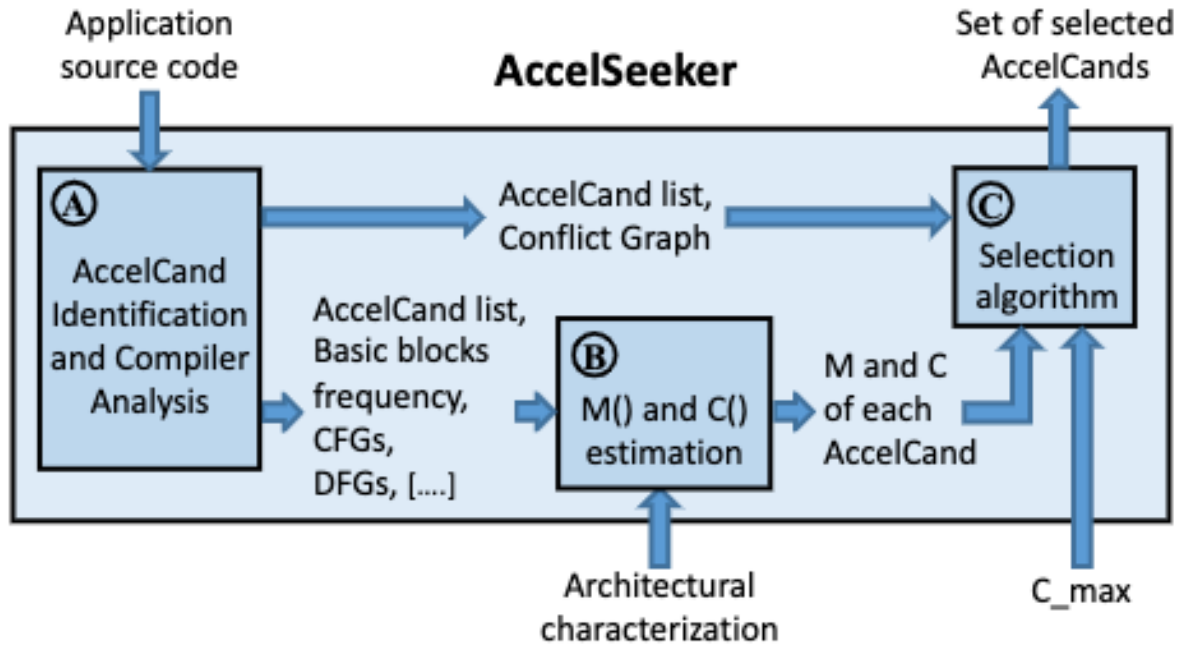
Jasmine Toolflow: *SoC Accelerator Discovery*



- **AccelSeeker**[ICCD,2019](#)
 - Identify basic-block level parallelism for HW acceleration
- **Trireme**[ACM TECS, 2023](#)
 - Leverages HPVM data flow graph to identify multiple forms of parallelism for further hardware acceleration
- **NOVIA**[MICRO,2021](#)
 - Identify and extract NCA (inline fine-grained accelerators integrated to core) for less area overhead and small startup time
- **AccelMerger**[arxiv,2021](#)
 - Coarse grained, loosely coupled accelerator
 - Identify opportunities to merge kernels for resource-efficient acceleration
- **SSAge**
 - ML driven IR-level power models for merging accelerators

The tools that seemed to make a difference

AccelSeeker* identifies HW acceleration candidates



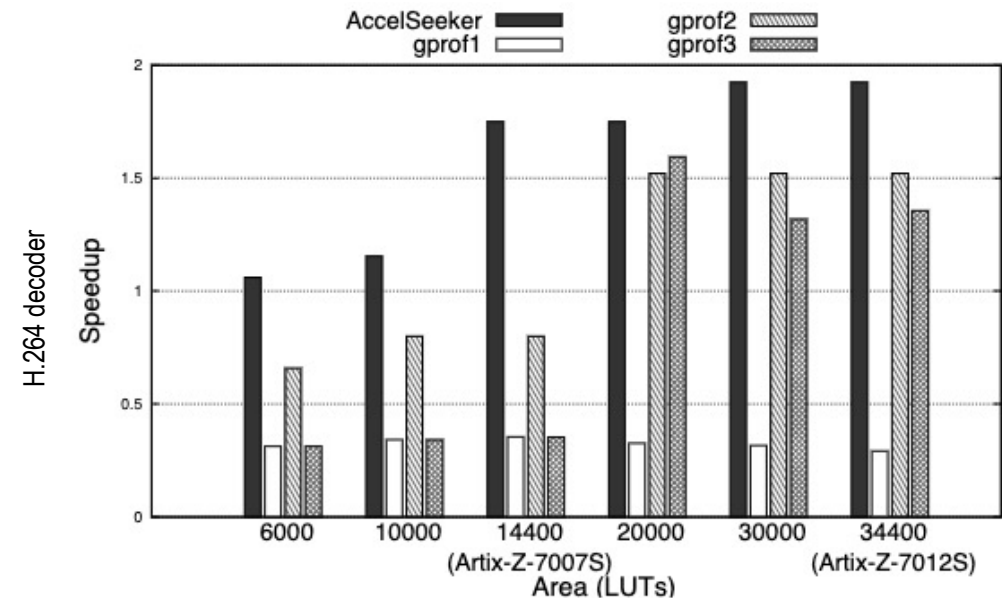
- Automated HW/SW codesign
- What parts of Application should be Accelerated?
- Identify basic-block level parallelism for HW acceleration.

MERIT $M()$

SW latency
 HW Computation Latency
 HW Communication Latency
 Runtime execution Frequency

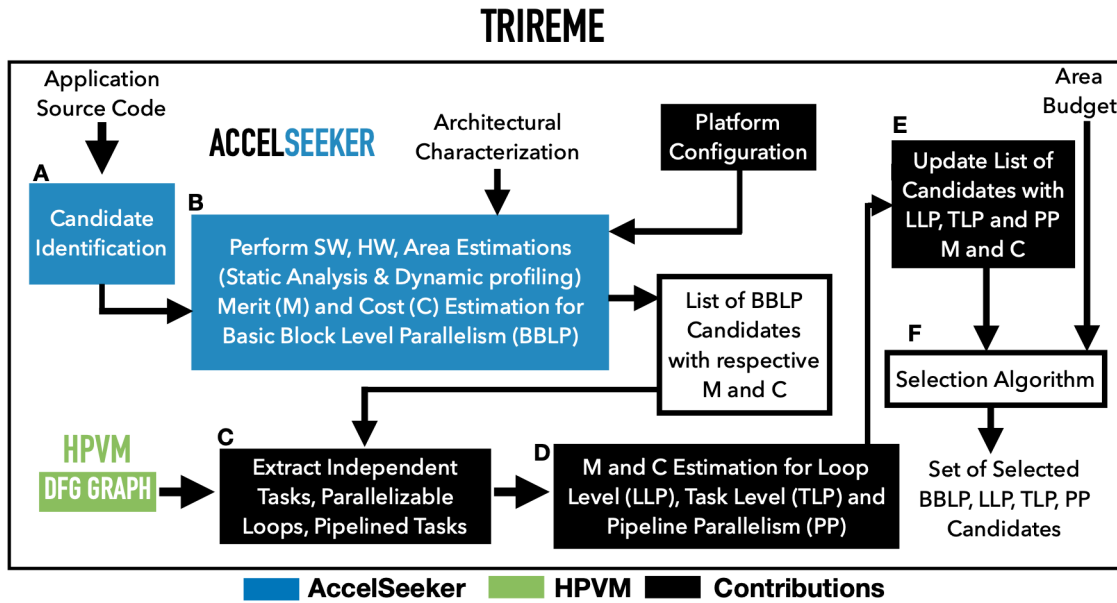
COST $C()$

HW Area
 Resources: LUTs, DSPs, BRAMs

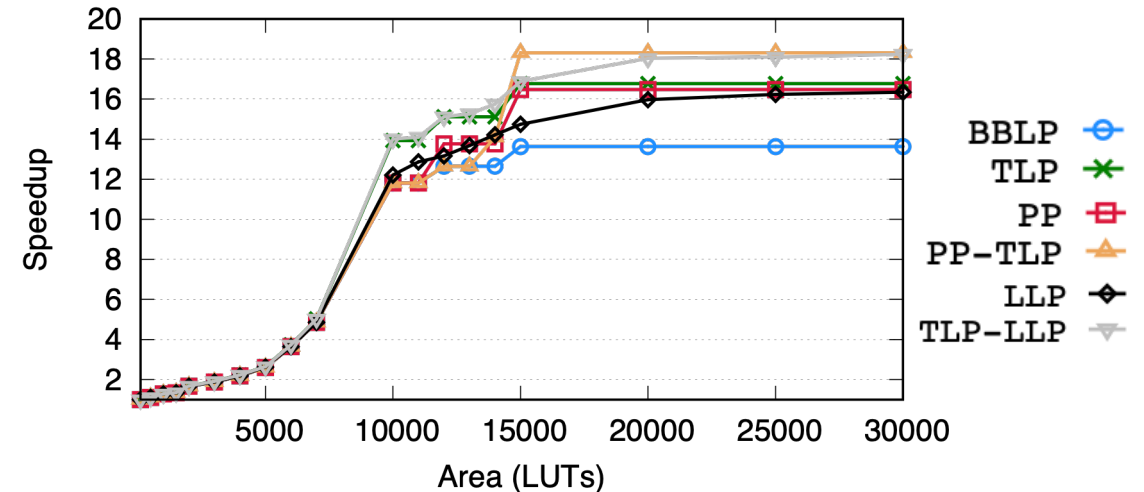
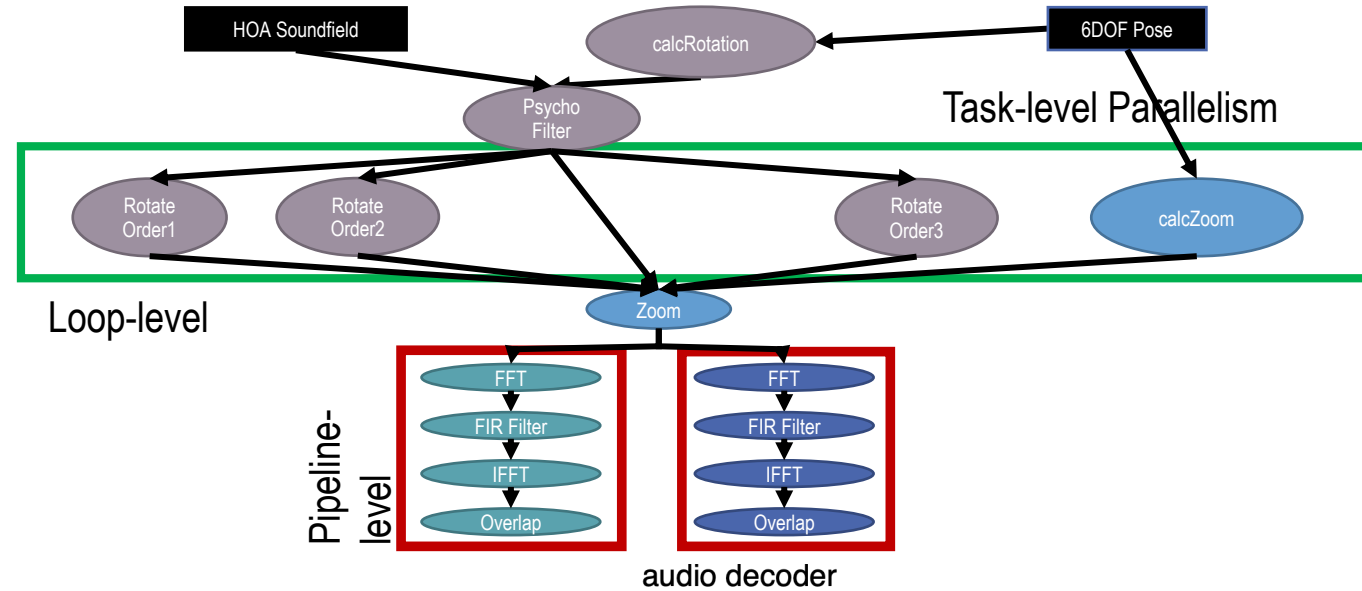


* Zacharopoulos, Georgios, et al. "Compiler-assisted selection of hardware acceleration candidates from application source code." IEEE International Conference on Computer Design (ICCD), 2019.

Trireme*: Automatic Identification of Accelerators using HPVM



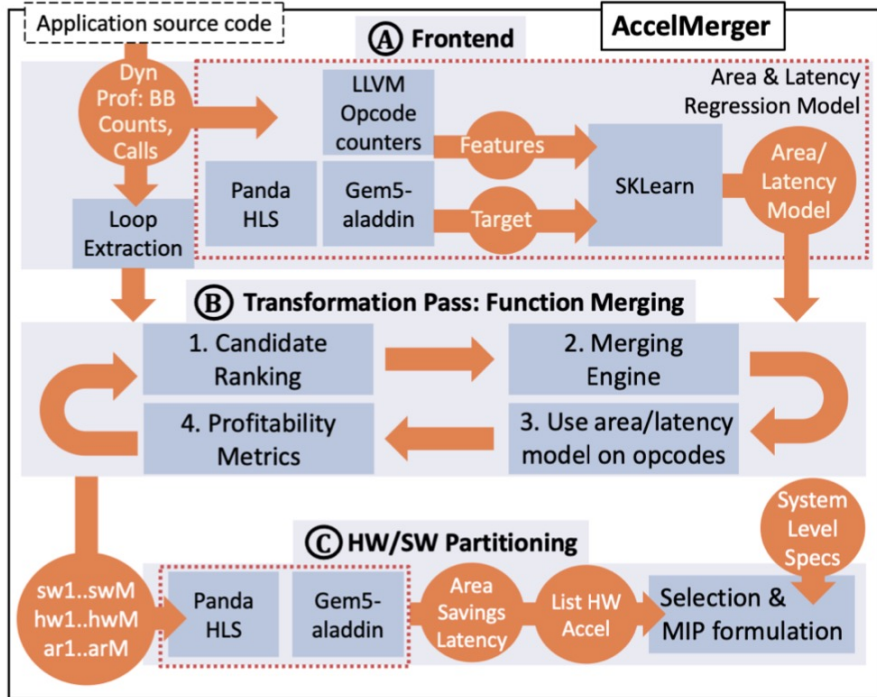
- Leverages HPVM to identify different types of parallelism
 - Task Parallelism (TLP)
 - Loop Parallelism (LLP)
 - Pipeline Parallelism (PP)
- Incorporates them into Metric and Cost estimation



* Zacharopoulos, Georgios, et al. Trireme: “Exploration of Hierarchical Multi-level Parallelism for Hardware Acceleration,” ACM TECS, vol. 22, issue 3, May 2023.

Accelmerger*: Automatic Generation of Coarse Grained Merged Accelerators

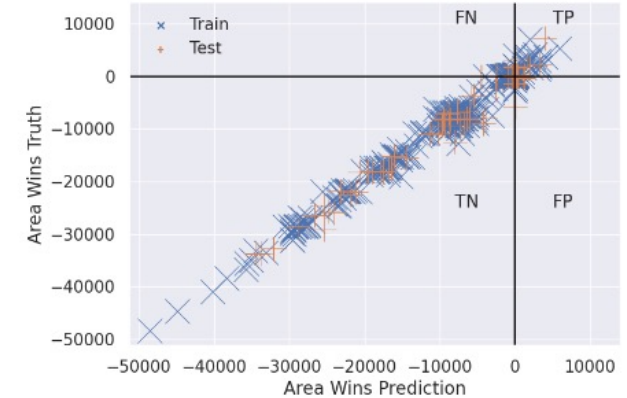
Methodology



```

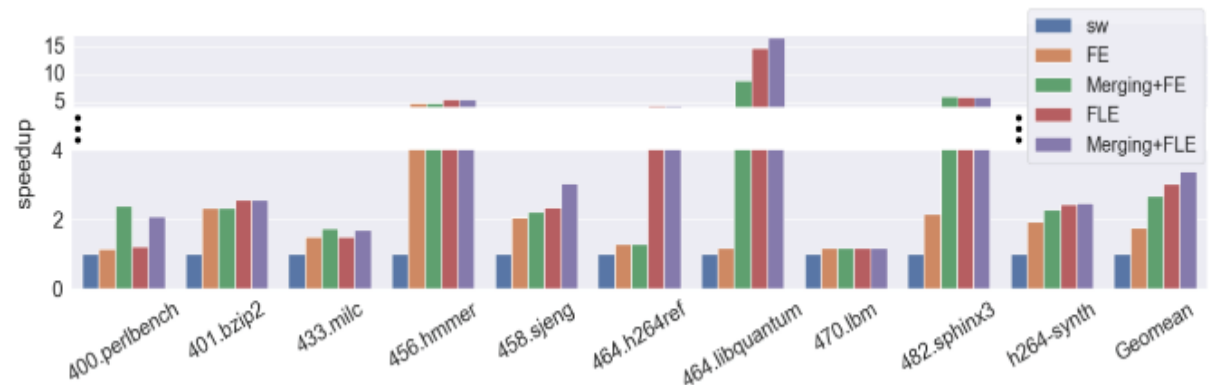
1 int f1(int a, int b, bool sum) {
2   int c;
3   if (sum)
4     c = a + b;
5   else
6     c = a * b;
7   c = f3(c, a);
8   return c;
9 }
10
11 int f2(int a, int b, int d, bool mult) {
12   int c;
13   if (mult)
14     c = a * b;
15   else
16     c = d + b;
17   return c;
18 }
19
20 Merged Function
21
22 int f12(int a, int b, int d, bool sum,
23        bool mult, bool f_sel) {
24   int c;
25   bool cond = f_sel : sum, not mult;
26   if (cond) {
27     int op = f_sel : a, d;
28     c = op + b;
29   } else
30     c = a * b;
31   if (f_sel)
32     c = f3(c, a);
33   return c;
34 }

```



- Functions are merged when the cost model estimates a profitable circuitry reduction.
- HW/SW partitioning enables speedup estimates over SW-only executions.

Results

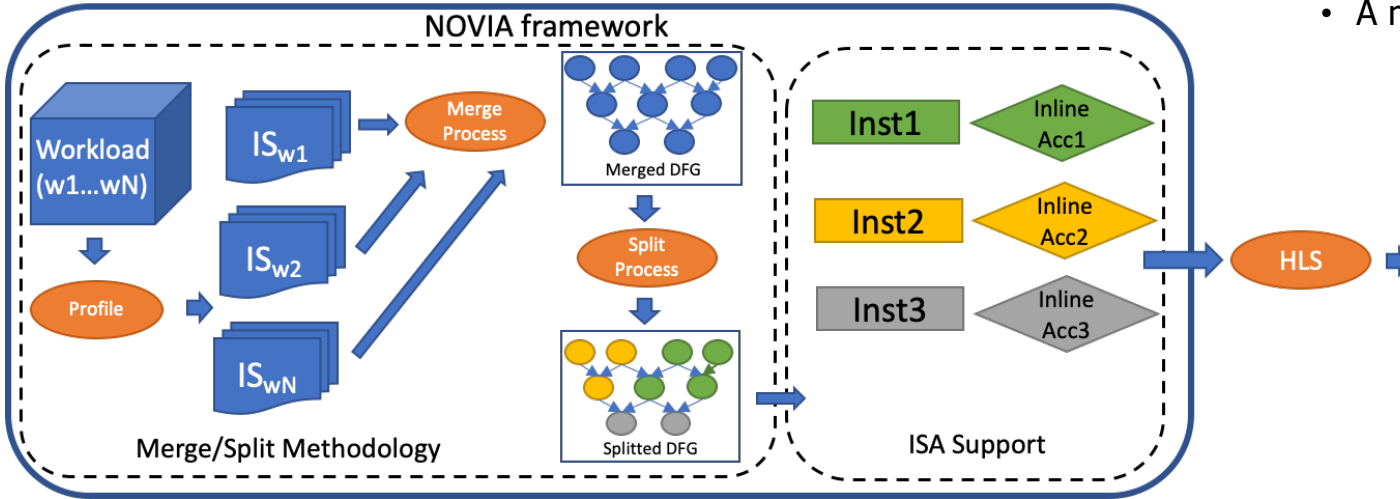


Identify opportunities to merge kernels for resource-efficient acceleration yielding non-conventional accelerators (NCAs), using an accurate ML resource model.

* Brumar, Iulian, et al. "Early DSE and Automatic Generation of Coarse Grained Merged Accelerators." *arXiv preprint arXiv:2111.09222* (2021)

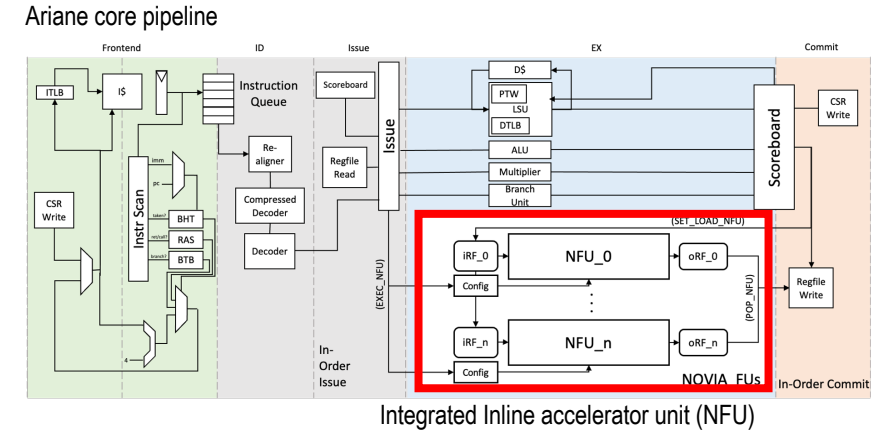
NOVIA*: A Framework for Discovering Non-Conventional Inline Accelerators

Systematic Methodology



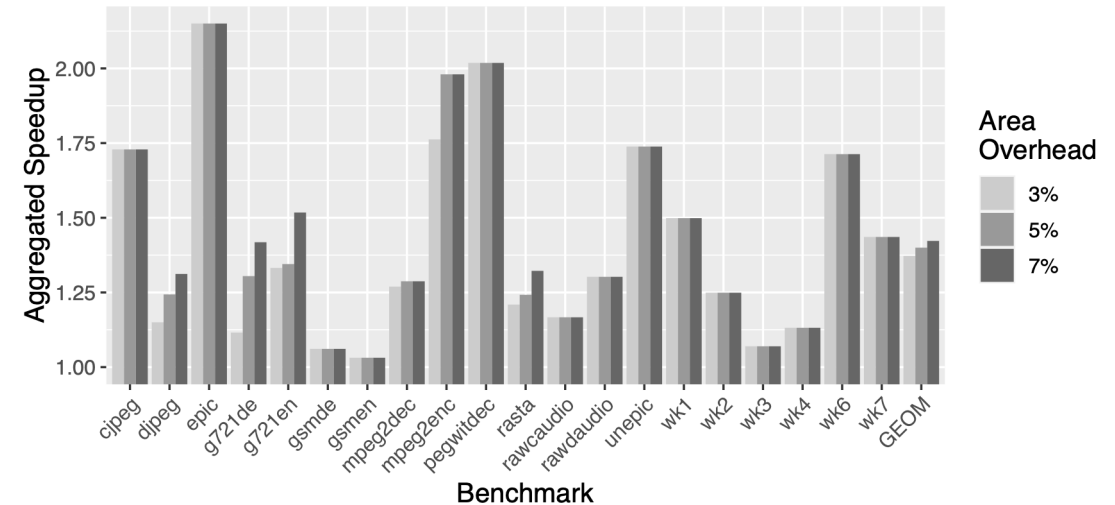
Integration with Ariane Cores

- A modular interface for any proposed inline accelerators



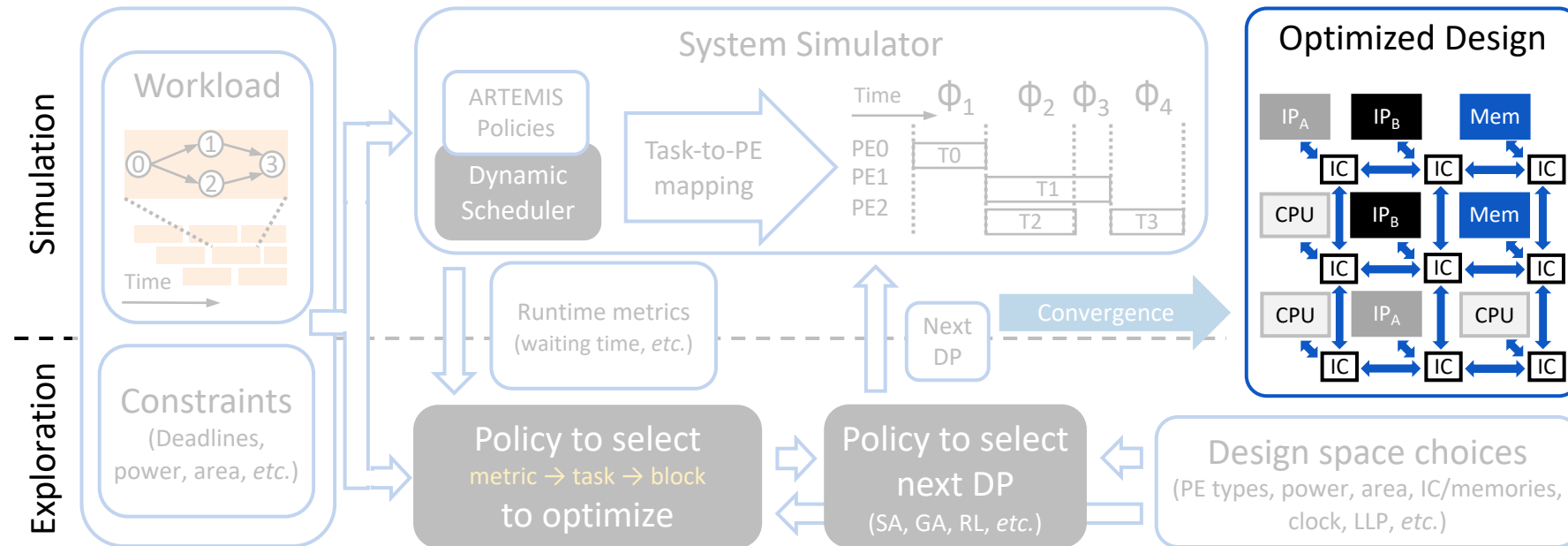
- Offload accelerators:**
 - Large enough computation to compensate for invocation overheads
 - Offload accelerator start-up can restrict the functions that can be speed up
- Inline accelerators:**
 - Faster start-up times
 - allowing them to provide speedup on smaller pieces of the workload
 - Close to core pipeline
 - Customized functional units
 - Minimal integration effort

Results



* Trilla, David, et al. "NOVIA: A Framework for Discovering Non-Conventional Inline Accelerators." *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture, 2021.*

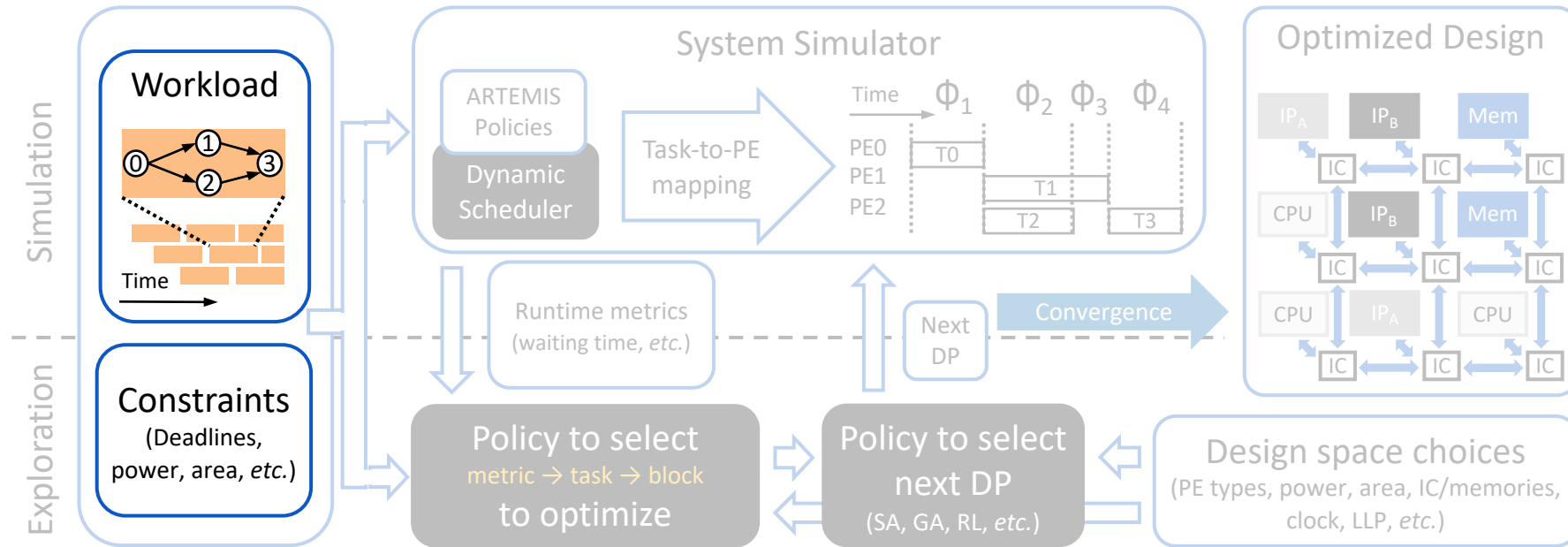
The ARTEMIS* Framework



- ARTEMIS consists of an exploration and a simulation framework
 - Begins with a seed SoC DP and incrementally transforms it to an optimized DP that meets the pre-specified constraints

*Subhankar Pal, Aporva Amarnath, Behzad Boroujerdian, Augusto Vega, Alper Buyuktosunoglu, John-David Wellman, Vijay Janapa Reddi, Pradip Bose, "ARTEMIS: Agile Discovery of Efficient Real-Time Systems-on-Chips in the Heterogeneous Era," HPCA 2025.

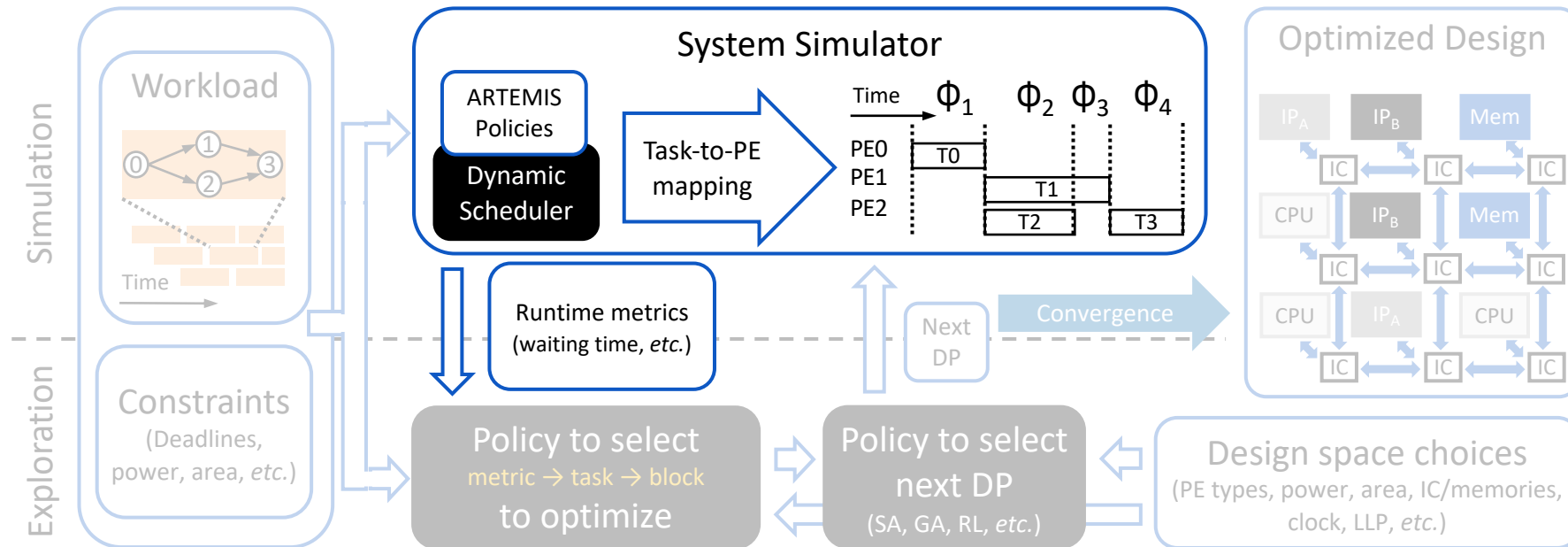
The ARTEMIS* Framework



- DAG representation of the workload and application and system constraints are sent as inputs to the framework

*Subhankar Pal, Aporva Amarnath, Behzad Boroujerdian, Augusto Vega, Alper Buyuktosunoglu, John-David Wellman, Vijay Janapa Reddi, Pradip Bose, "ARTEMIS: Agile Discovery of Efficient Real-Time Systems-on-Chips in the Heterogeneous Era," HPCA 2025.

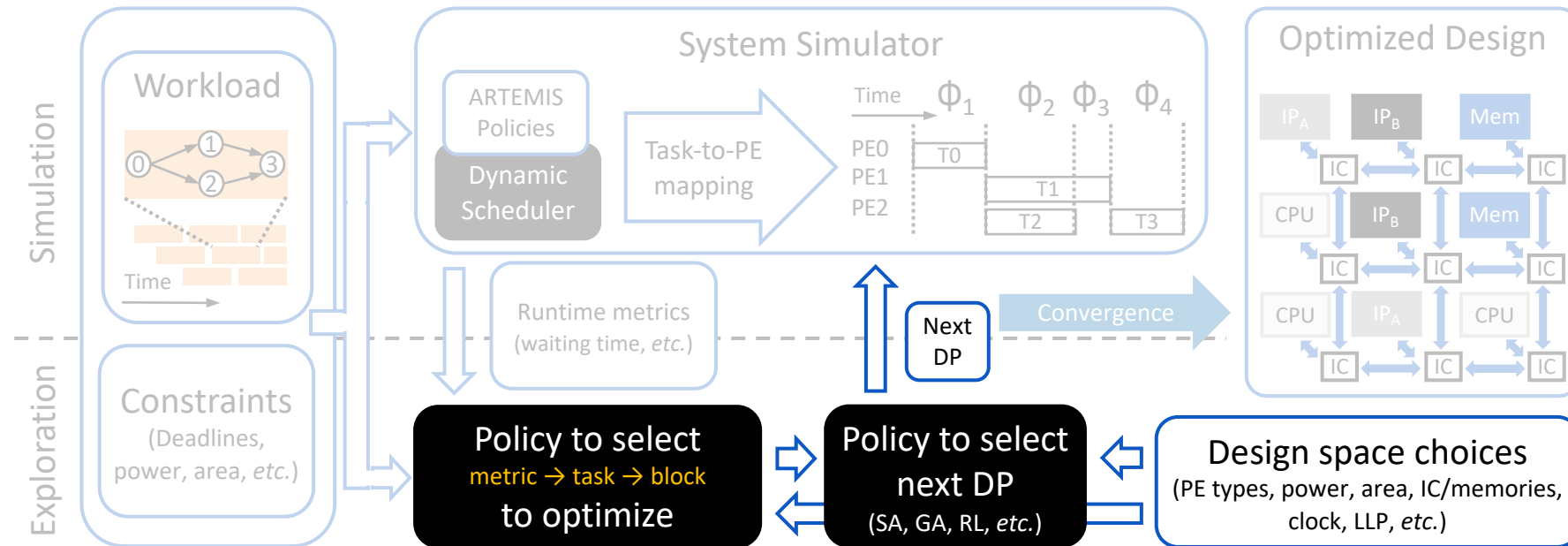
The ARTEMIS* Framework



- The simulator uses an SoC roofline model to perform phase-driven simulation
 - We use a dynamic scheduler in the simulator to perform fast and efficient task-to-PE mapping and identify bottlenecks in the system
- Several runtime statistics (e.g., task waiting times and task deadlines) are collected and fed into the explorer

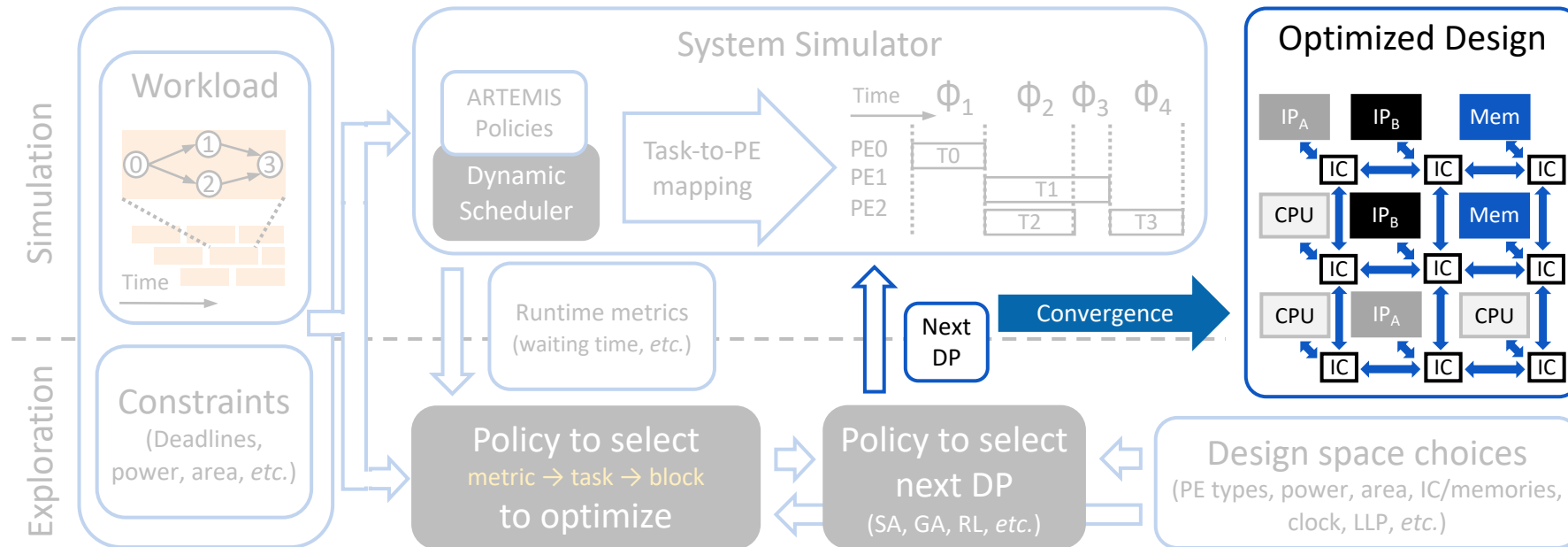
*Subhankar Pal, Aporva Amarnath, Behzad Boroujerdian, Augusto Vega, Alper Buyuktosunoglu, John-David Wellman, Vijay Janapa Reddi, Pradip Bose, "ARTEMIS: Agile Discovery of Efficient Real-Time Systems-on-Chips in the Heterogeneous Era," HPCA 2025.

The ARTEMIS* Framework



- The explorer uses RT-aware policies to iteratively select ① the metric to optimize (latency/power/area), ② the task to optimize to improve this metric, and ③ the hardware block to improve upon
- A library of pre-characterized PE/IC/memory blocks are fed in as inputs
- The next design point is selected based on an architecture-aware simulated annealing process

The ARTEMIS Framework



- DSE continues until ARTEMIS encounters a DP that meets deadlines for all DAGs, and has power and area within the constraints

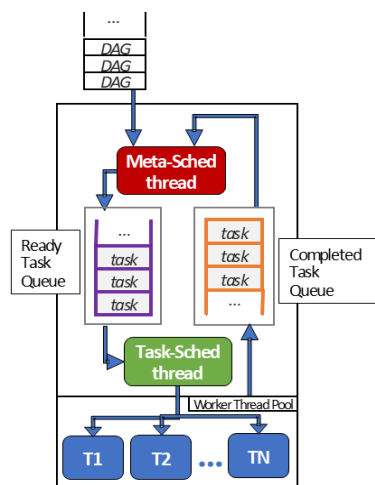
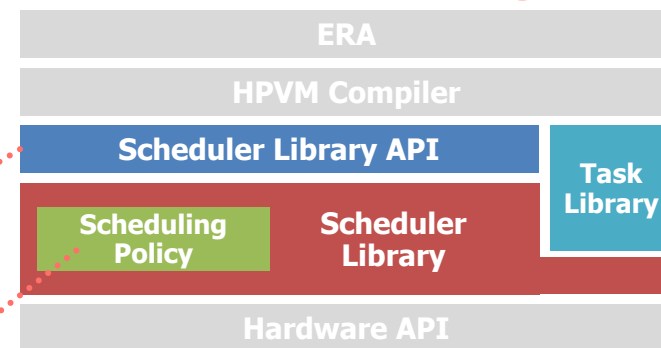
EPOCHS Appliance: The Smart Scheduler

EPOCHS Scheduler

github.com/IBM/scheduler-library

- Provides “smart” scheduling capabilities to applications in a hardware-independent manner
 - Schedule tasks across heterogeneous processing elements
 - Improve metrics of interest (throughput, criticality, efficiency)
 - Implemented as a user-level library (Scheduler Library)
- Exposes an API that abstracts applications from the scheduler
- User-specified scheduling policies can be easily “plug” into the scheduler

Full software stack integration



Two-level scheduling: **META** policies (level 1)

- Assign rank based on our HetSched paper [1]

Two-level scheduling: **TASK** policies (level 2)

(choose the PE on HW where task should be executed)

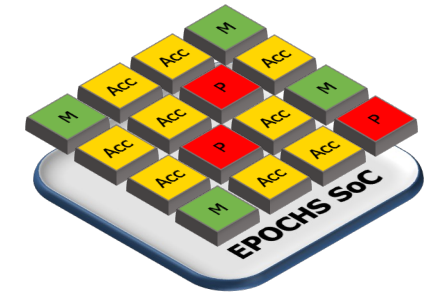
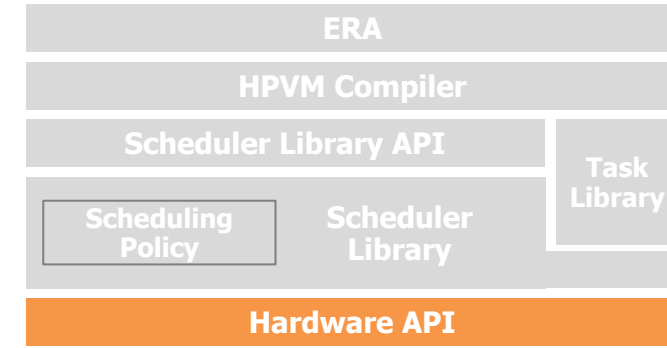
[1] Amarnath, Aporva, et al. “Heterogeneity-Aware Scheduling on SoCs for Autonomous Vehicles”, IEEE CAL (2021).

EPOCHS Appliance: The Agile SoC Flow

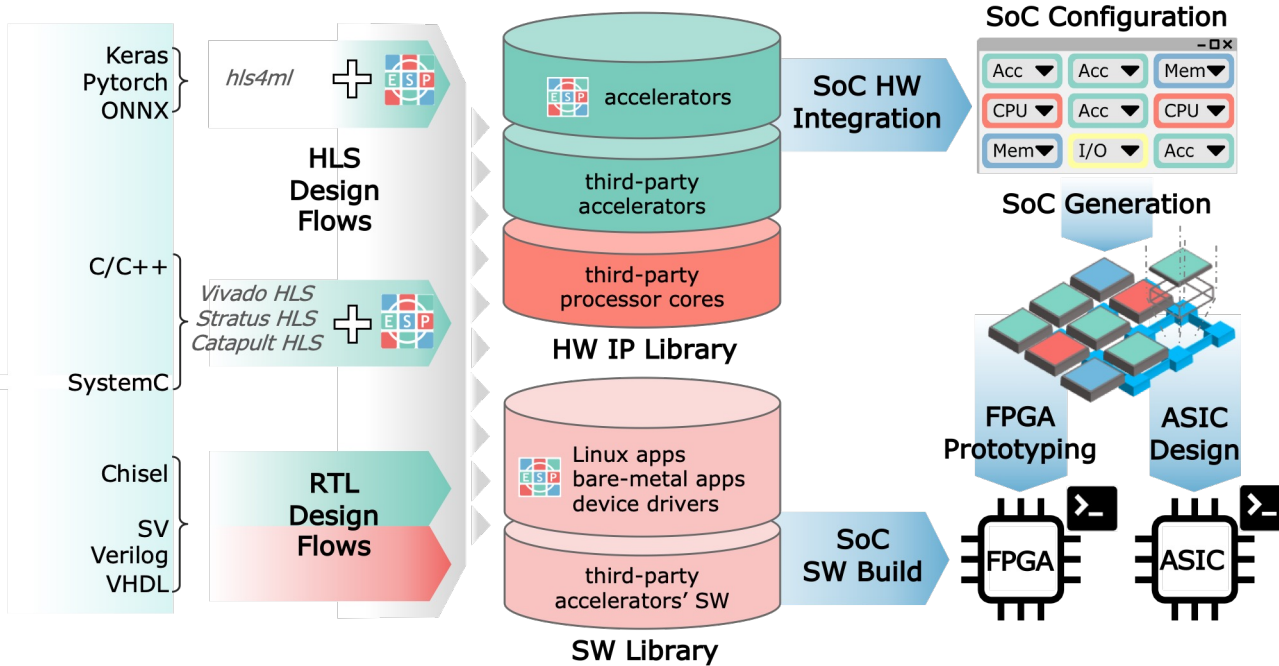
ESP: the open-source agile flow for system-on-chip (SoC) design

- Seamless integration of SoC components with NoC and platform services
- Push-button generation of complete SoC RTL for ASIC physical design
- **Two EPOCHS full SoC ASIC chip tapeouts to date**
- Rapid FPGA prototyping with software build for early application development

Full software stack integration



www.esp.cs.columbia.edu



Accomplishments Summary



EPOCHS-0 SoC tapeout

- 4x4 SoC fabricated



1.52 GHz peak @ 1.0 V
 2.43mW @ 0.5V
 1.83W @ 1.0 V
ESSCIRC-2022

Scaled-out EPOCHS-1 SoC tapeout

- 6x6 SoC with new accelerators
- Novel distributed h/w power manager



Significant design cost mitigation

- 10x-100x reduction in person-years

Hardware-agnostic programming of heterogeneous SoCs

- HPVM compiler, smart scheduler...

Open-source ecosystem for collaboration

ERA: github.com/IBM/era HPVM: gitlab.engr.illinois.edu/llvm/hpvm-release

Mini-ERA: github.com/IBM/mini-era STOMP: github.com/IBM/stomp

ESP: www.esp.cs.columbia.edu Scheduler: github.com/IBM/scheduler-library

Spandex: github.com/sld-columbia/esp/tree/master/rtl/caches

Benefits of acceleration

| | FFT | Viterbi |
|-------------|------|---------|
| Performance | 71x | 20x |
| Energy | 233x | 56x |

→ Tape-out: Sept. 2021, respin: Nov. 2022

ISSCC-2024

A 12nm Linux-SMP-Capable RISC-V SoC with 14 Accelerator Types, Distributed Hardware Power Management and NoC-Based Data Orchestration

Maico Cassel dos Santos^{1*}, Tianyu Jia^{2*}, Joseph Zuckerman^{1*}, Martin Cochet^{3*}, Davide Giri¹, Erik Loscalzo¹, Karthik Swaminathan³, Thierry Tambe², Jeff Jun Zhang², Alper Buyuktosunoglu³, Kuan-Lin Chiu¹, Giuseppe Di Guglielmo¹, Paolo Mantovani¹, Luca Piccolboni¹, Gabriele Tombesi¹, David Trilla³, John-David Wellman³, En-Yu Yang², Aporva Amarnath³, Ying Jing¹, Bakshree Mishra⁴, Joshua Park², Vignesh Suresh⁴, Sarita Adve⁴, Pradip Bose³, David Brooks², Luca P. Carloni¹, Kenneth L. Shepard¹, Gu-Yeon Wei²

* These authors have equal contributions.

¹ COLUMBIA UNIVERSITY IN THE CITY OF NEW YORK ² HARVARD UNIVERSITY ³ IBM Research ⁴ UNIVERSITY OF ILLINOIS

© 2024 IEEE International Solid-State Circuits Conference 14.5. A 12nm Linux-SMP-Capable RISC-V SoC with 14 Accelerator Types, Distributed Hardware Power Management and NoC-Based Data Orchestration 1 of 21

2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)

ISCA-2024

BlitzCoin: Fully Decentralized Hardware Power Management for Accelerator-Rich SoCs

Martin Cochet¹, Karthik Swaminathan¹, Erik Loscalzo², Joseph Zuckerman², Maico Cassel dos Santos², Davide Giri², Alper Buyuktosunoglu¹, Tianyu Jia³, David Brooks³, Gu-Yeon Wei³, Kenneth Shepard², Luca P. Carloni², and Pradip Bose¹
¹IBM Research, Yorktown Heights, NY ²Columbia University, New York, NY ³Harvard University, Cambridge, MA



ERI Summit 2023

Davide Giri (1990-2021): Key Contributor



Thank you for everything, Davide!



Karthik Swaminathan (1984 – 2026): KEY CONTRIBUTOR



Thank you from colleagues and friends 🙏



- **Star performer** in Pradip's team at IBM ("Efficient & Resilient Systems")
 - ✓ At IBM Research (fulltime) since 2014.
 - ✓ Summer intern before that.
- **Most recent key accomplishments:**
 - ✓ Key contributor to the DoW RAMP-C funded (SARA-1) AI/FHE hardware accelerator SoC design and implementation
 - [IBM Research Accomplishment award: 2025](#)
 - ✓ Novel decentralized hardware power management (DHPM) architecture and design – part of EPOCHS-1 SoC (DARPA)
 - **BlitzCoin** paper at ISCA-2024 ("Top Pick", IEEE Micro, 2025)
 - **EPOCHS-1 SoC** at ISSCC-2024; journal paper, JSSCC, 2026)
 - [IBM Research accomplishment award \(EPOCHS/DSSoC\): 2024](#)
 - ✓ Novel energy-efficient AI hardware research and development
 - Bit-error robustness for energy-efficient AI/ML systems:
BERRY: DAC 2023, MuBERRY: ASPLOS 2024 (joint work, GaTech/IBM)
 - ✓ SERMiner and other methodology innovations in support of IBM's POWER and mainframe processor products: DSN-2021.
- **Ph.D in Computer Science and Engineering**, Penn State University (2014)
- **B.Tech + M.Tech in Electrical Engineering**, IIT Madras (2008)



[ERI Summit 2023](#)



Thank You!

Pradip Bose + **many colleagues** at IBM T. J. Watson Research Center
Sarita Adve, Vikram Adve, Sasa Misailovic, University of Illinois at Urbana-Champaign
Luca Carloni, Ken Shepard, Columbia University
David Brooks, Gu-Yeon Wei, Vijay Janapa Reddi, Harvard University
plus many talented students and postdoctoral fellows