

HeroSDK: Streamlining Heterogeneous RISC-V Accelerated Computing From Embedded To High-Performance Systems

International Conference on Computer Design (ICCD)
Milan, Italy, 19-11-2024

Cyril Koenig¹

cykoenig@iis.ee.ethz.ch

Björn Forsberg²

bjorn.forsberg@ri.se

Luca Benini^{1,3}

lbenini@iis.ee.ethz.ch

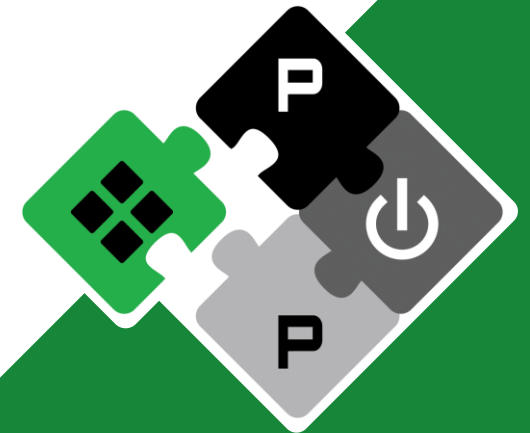
¹ IIS, ETH Zürich, Switzerland

² RISE, Sweden

³ DEI, University of Bologna, Italy

PULP Platform

Open Source Hardware, the way it should be!



@pulp_platform 

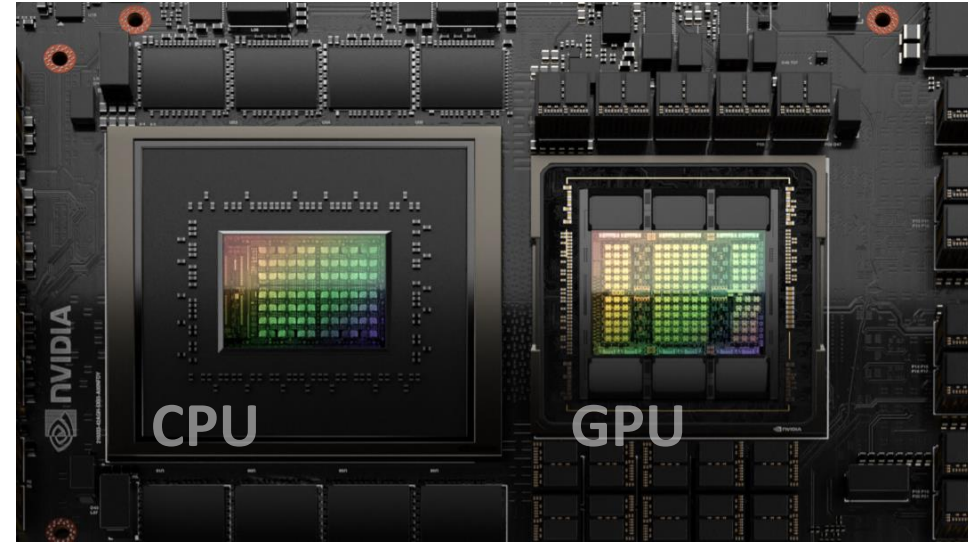
pulp-platform.org 

youtube.com/pulp_platform 

Where do we find heterogeneous systems?



www.apple.com



www.nvidia.com

- Embedded Heterogeneous SoCs (HeSoC)
- Heterogeneous HPC platforms
- 9 of the 10 most powerful supercomputers are heterogeneous



This heterogeneity brings challenges



- **Benefits**

- Energy efficiency
- Peak performance

- **Challenges**

- Synchronization
- Communication
- Programmability
- ISA and ABI differences
- Portability

These challenges are reduced with RISC-V



- **Benefits**

- Energy efficiency
- Peak performance
- Large design space exploration of accelerators (*open source ISA*)

- **Challenges**

- Synchronization
- Communication
- Programmability
- ISA ~~and ABI~~ differences
- Portability *simplified (same memory model, same base ISA, same non-ISA IP specs)*

RISC-V heterogeneous platforms are already there



| | OpenPiton [1] | Cohort [2] | Vortex [3] | HeroSDK (Our work) |
|------------------|---------------|-----------------------|--------------|------------------------------|
| Host | CVA6 | CVA6 | x86 | CVA6/SG2042 |
| Accelerator | MIAOW GPU | Dataflow accelerators | Vortex GPGPU | Snitch/Spatz clusters |
| On-chip offload | NoC | Crossbar | <i>None</i> | Crossbar |
| Off-chip offload | Chip Bridge | <i>None</i> | PCIe | PCIe |
| Programming | Custom | C++ Boost library | OpenCL | OpenMP |

[1] J. Balkind et al., "OpenPiton: An Open Source Hardware Platform for Your Research", IEEE MICRO 40 (2020)

[2] T. Wei et al., "Cohort: Software-Oriented Acceleration for Heterogeneous SoCs", ASPLOS'20 (2020)

[3] B. Tine et al., "Vortex: Extending the RISC-V ISA for GPGPU and 3D-Graphics", MICRO '21 (2021)

High level heterogeneous programming for both on-chip and off-chip accelerators

HeroSDK: Contributions



- A generic **OpenMP offloading runtime for RV based HeSoC** targeting **portability** across RV-accelerators.
- Its implementation on **two emulated HeSoC running Linux**.
- A Linux module to **export this framework from HeSoC to PCIe based accelerator cards with 100% code reuse**.
- We demonstrate matrix-vector multiplication with **5.9× speedup** with **19% offload-to-work ratio** in plain C-code with embedded accelerator.

Original Snitch-based demonstrator



- **Snitch embedded platform**

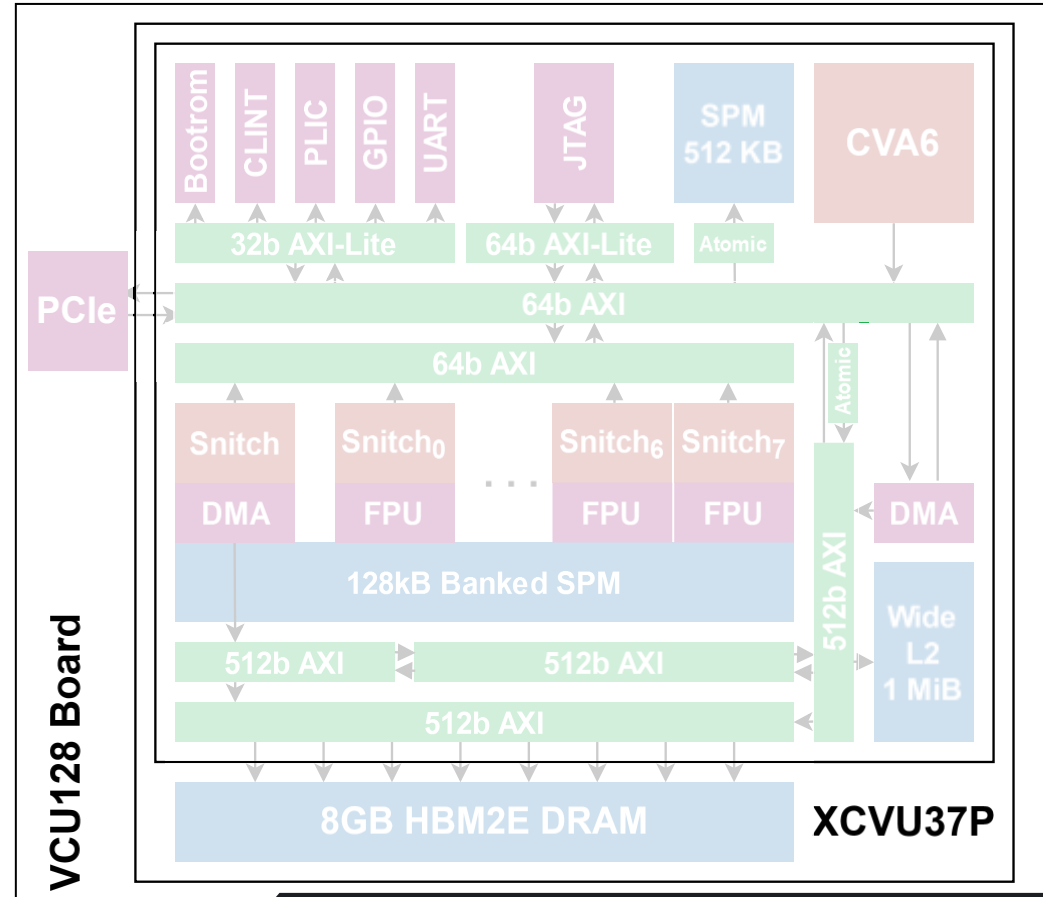
- RV64G single core host
- 8-core RV32IMAFD cluster
- Cluster scratchpad and DMA

- **Snitch accelerator card**

- 8-cores RV32IMAFD cluster
- Cluster scratchpad and DMA
- PCIe endpoint

- **+ SG2042 [1] RISC-V CPU**

- RV64G 64-cores host



www.github.com/pulp-platform/occamy

Original Snitch-based demonstrator



- **Snitch embedded platform**

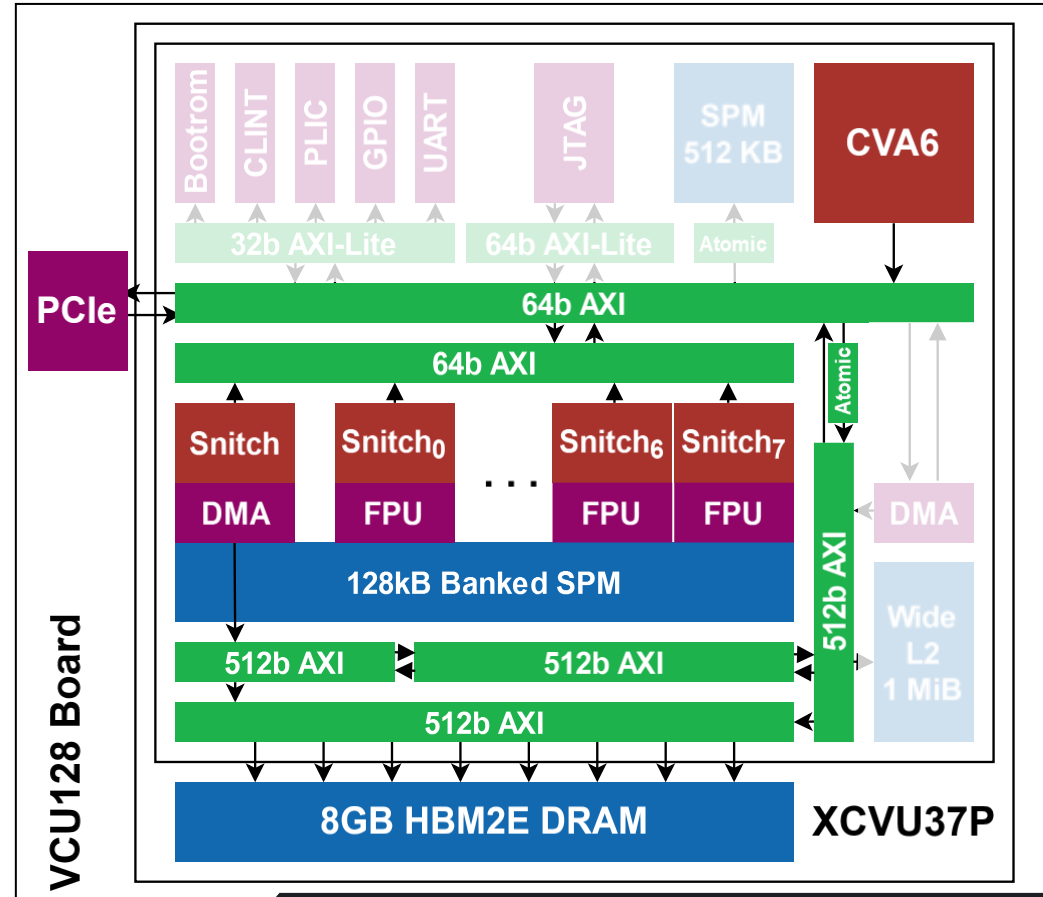
- RV64G single core host
- 8-core RV32IMAFD cluster
- Cluster scratchpad and DMA

- **Snitch accelerator card**

- 8-cores RV32IMAFD cluster
- Cluster scratchpad and DMA
- PCIe endpoint

- **+ SG2042 [1] RISC-V CPU**

- RV64G 64-cores host



 www.github.com/pulp-platform/occamy

Four demonstrator platforms



RV32IMAFD

Embedded

- **Snitch embedded platform**

- RV64G single core host
- 8-core RV32IMAFD cluster
- Cluster scratchpad and DMA

- **Snitch accelerator card**

- 8-core RV32IMAFD cluster
- Cluster scratchpad and DMA
- PCIe endpoint

- **Spatz embedded platform**

- RV64G single core host
- 2-core RV32IMAFDV* cluster
- Cluster scratchpad and DMA

- **Spatz accelerator card**

- 2-core RV32IMAFDV cluster
- Cluster scratchpad and DMA
- PCIe endpoint

RV32IMAFDV*

High performance

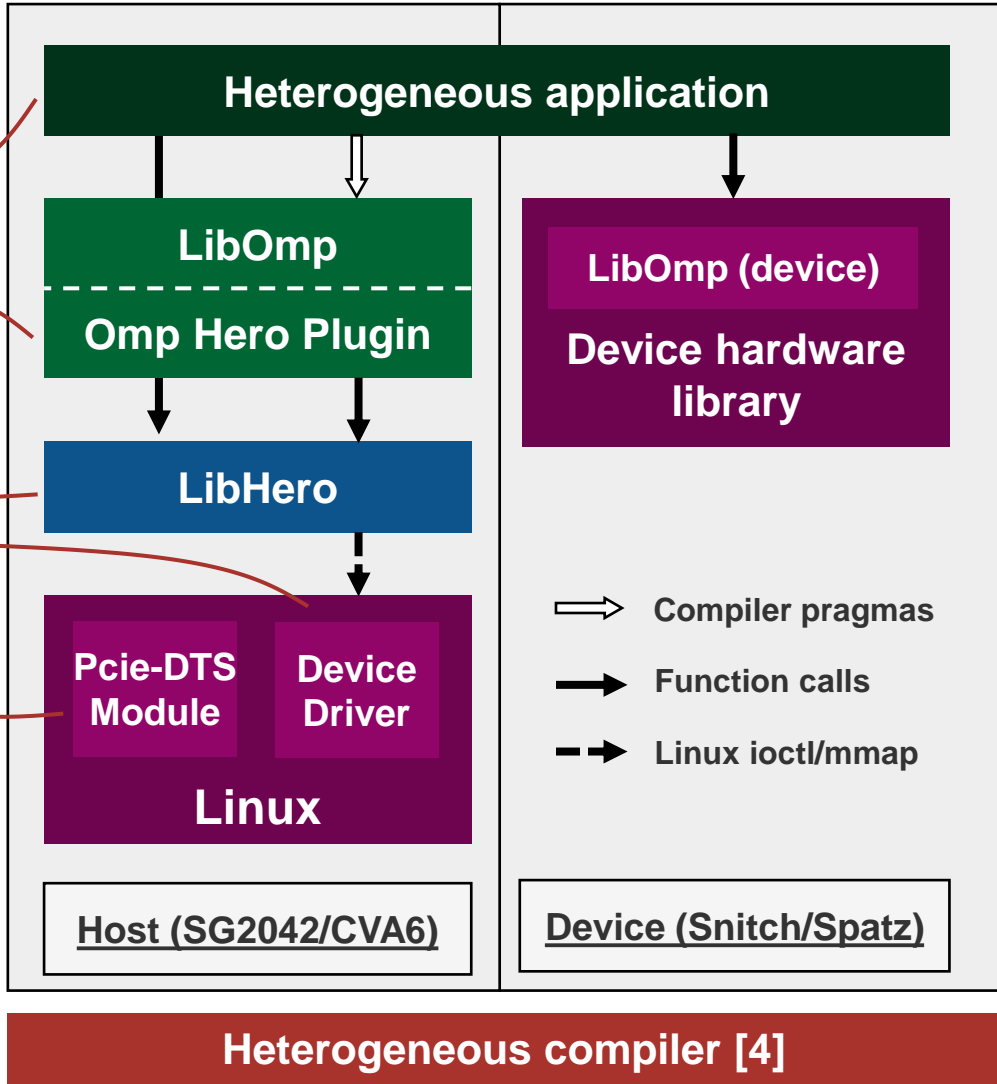
The HeroSDK software stack



- Your application using OpenMP
- Reusable OpenMP offload plugin
- Middle layer user-space library
- Device driver (kernel module)
- Extension to PCIe cards via device tree overlays [5]

```

pcie@706200000
├── ranges
├── pci@0,0
│   ├── ranges
│   ├── dev@0,0
│   │   ├── ranges
│   │   ├── accelerator-soc
│   │   │   ├── ranges
│   │   │   └── snitch-cluster@10000000
│   └── ...
    
```



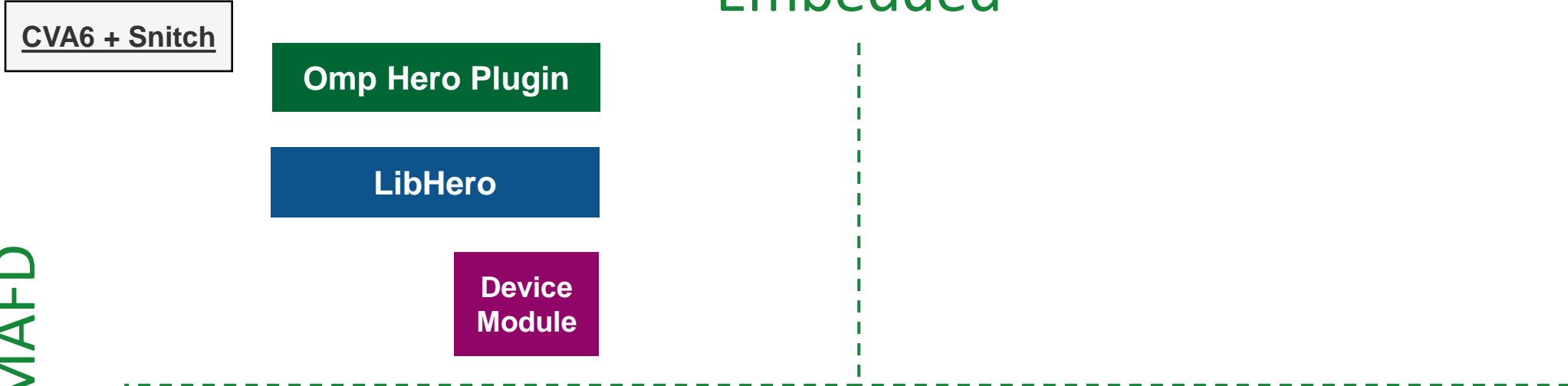
[4] Kurth et al. HEROv2: Full-Stack Open-Source Research Platform for Heterogeneous Computing
 [5] PCIe node introduced in Linux 6.6-rc5 (Lizhi Hou)

First support for RV32IMAFD

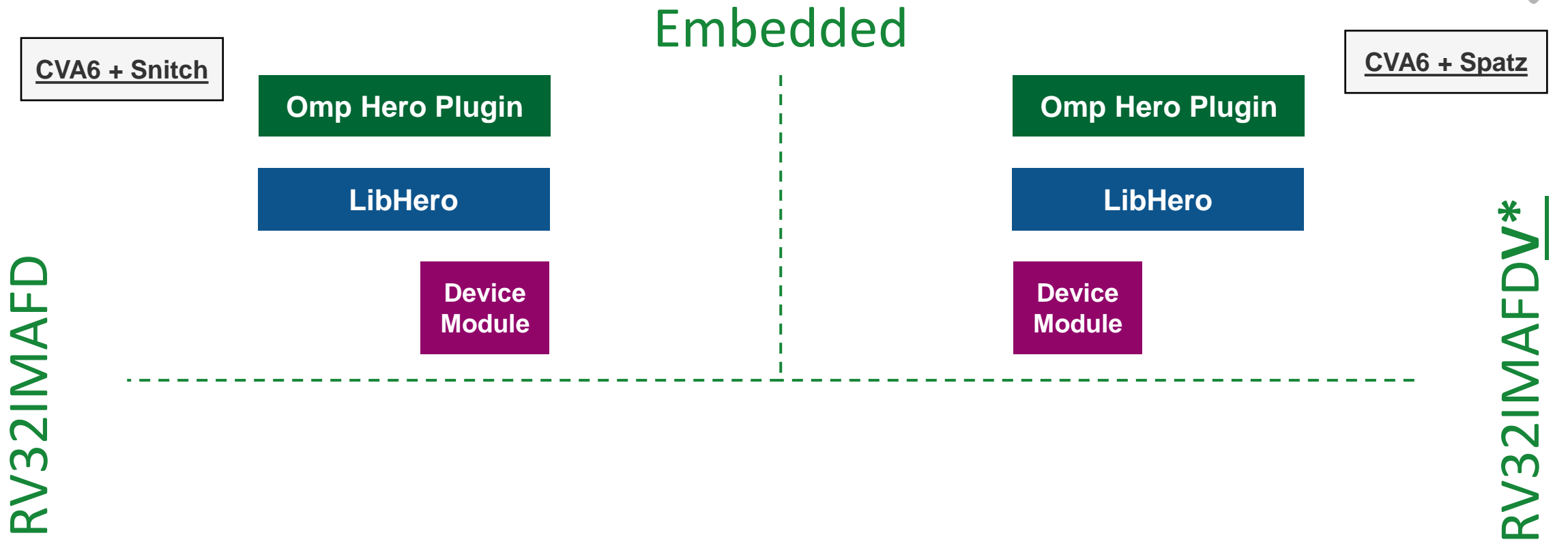


RV32IMAFD

Embedded

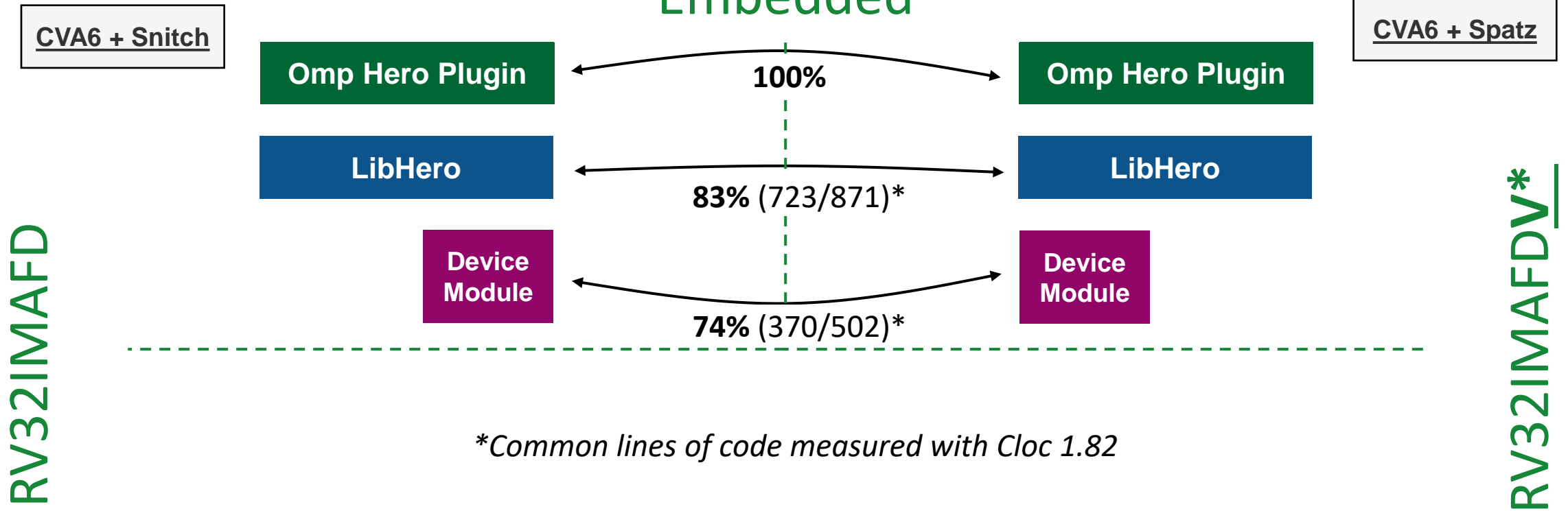


Extension to RV32IMAFDV



*Zve64d

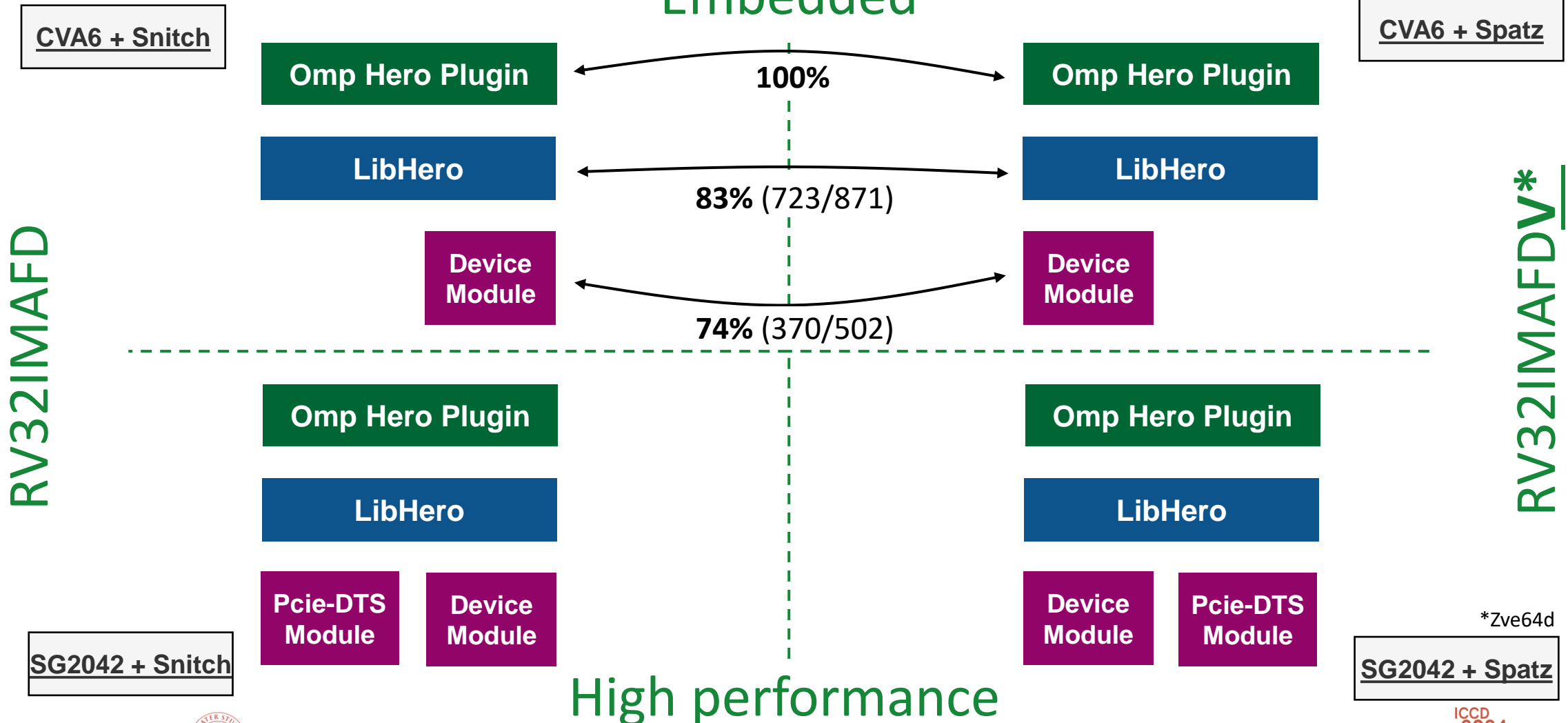
Extension to RV32IMAFDV with high reusability



**Common lines of code measured with Cloc 1.82*

*Zve64d

Extension to PCIe platform



RV32IMAFD

RV32IMAFDV*

CVA6 + Snitch

CVA6 + Spatz

SG2042 + Snitch

SG2042 + Spatz

Extension to PCIe platform with 100% reusability



RV32IMAFD

RV32IMAFDV*

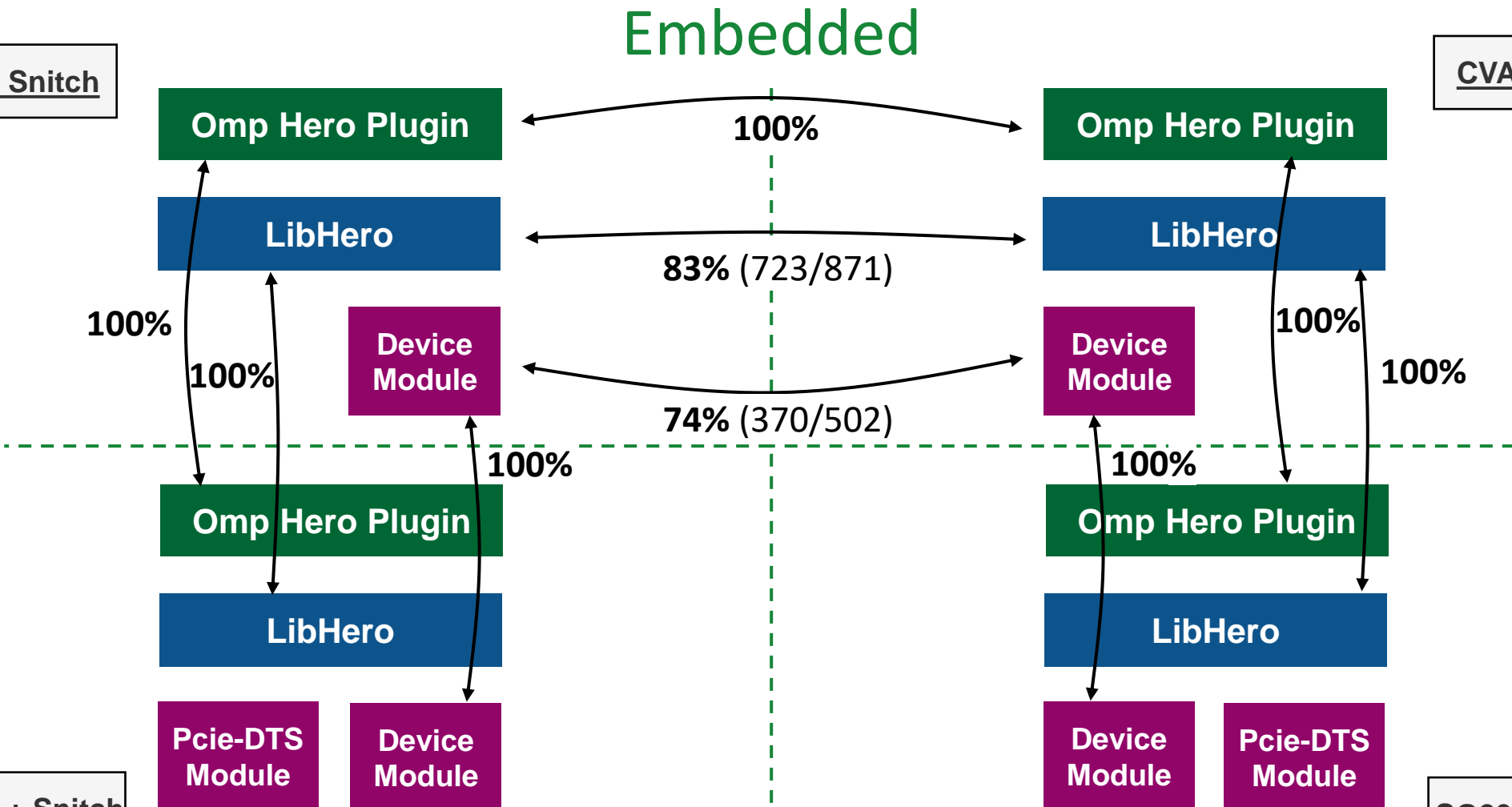
CVA6 + Snitch

CVA6 + Spatz

SG2042 + Snitch

SG2042 + Spatz

*Zve64d



Measuring offload latency of HeroSDK



Three micro kernels to measure offload latency

a) No-op

b) Input read

c) Input read and Output write

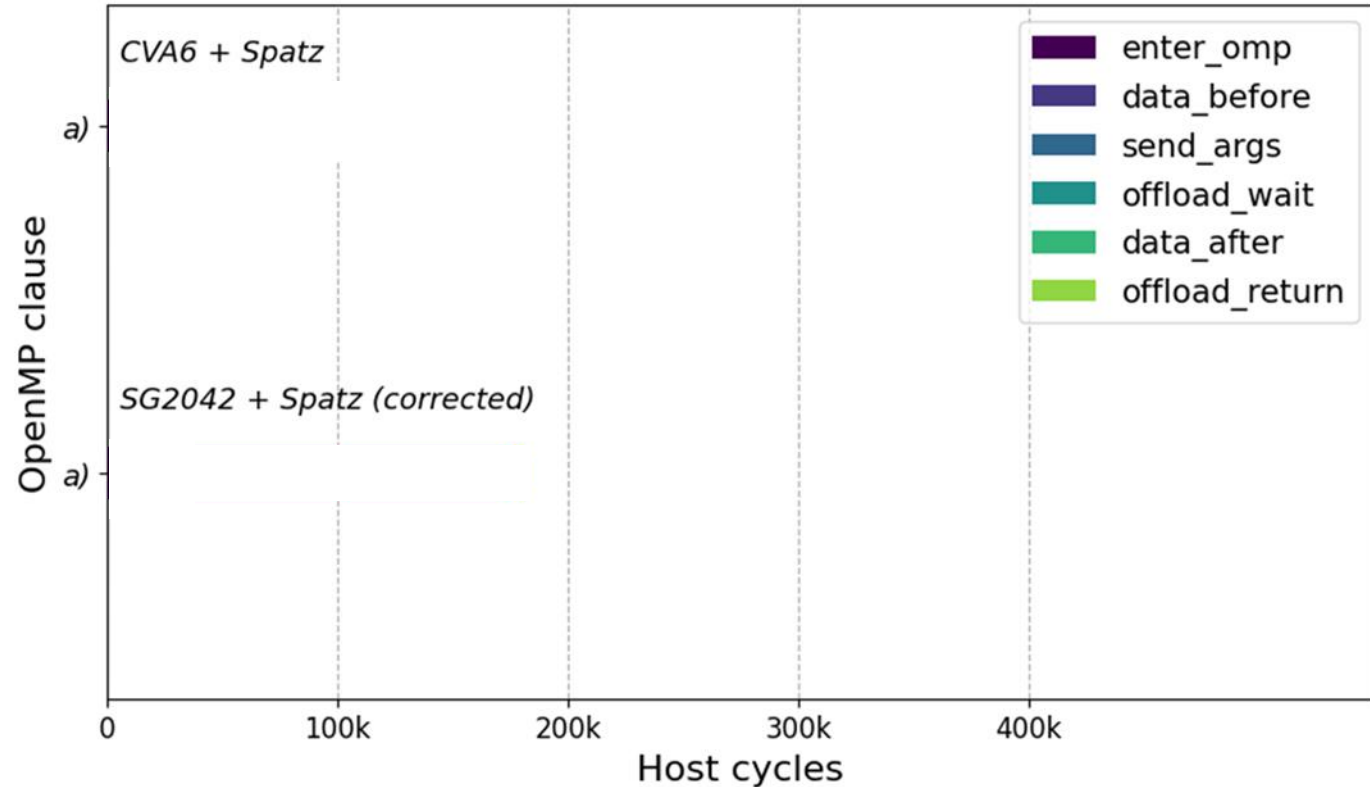
Process kernel info on the host

> faster on SG2042 (triple issue C910 core)

Less than 100k cycle OpenMP offload overhead

Synchronous fork/join on device

> requires PCIe access on SG2042



Offload benchmarks



Three micro kernels to measure offload latency

a) No-op

b) Input read

c) Input read and Output write

Process kernel infos on the host

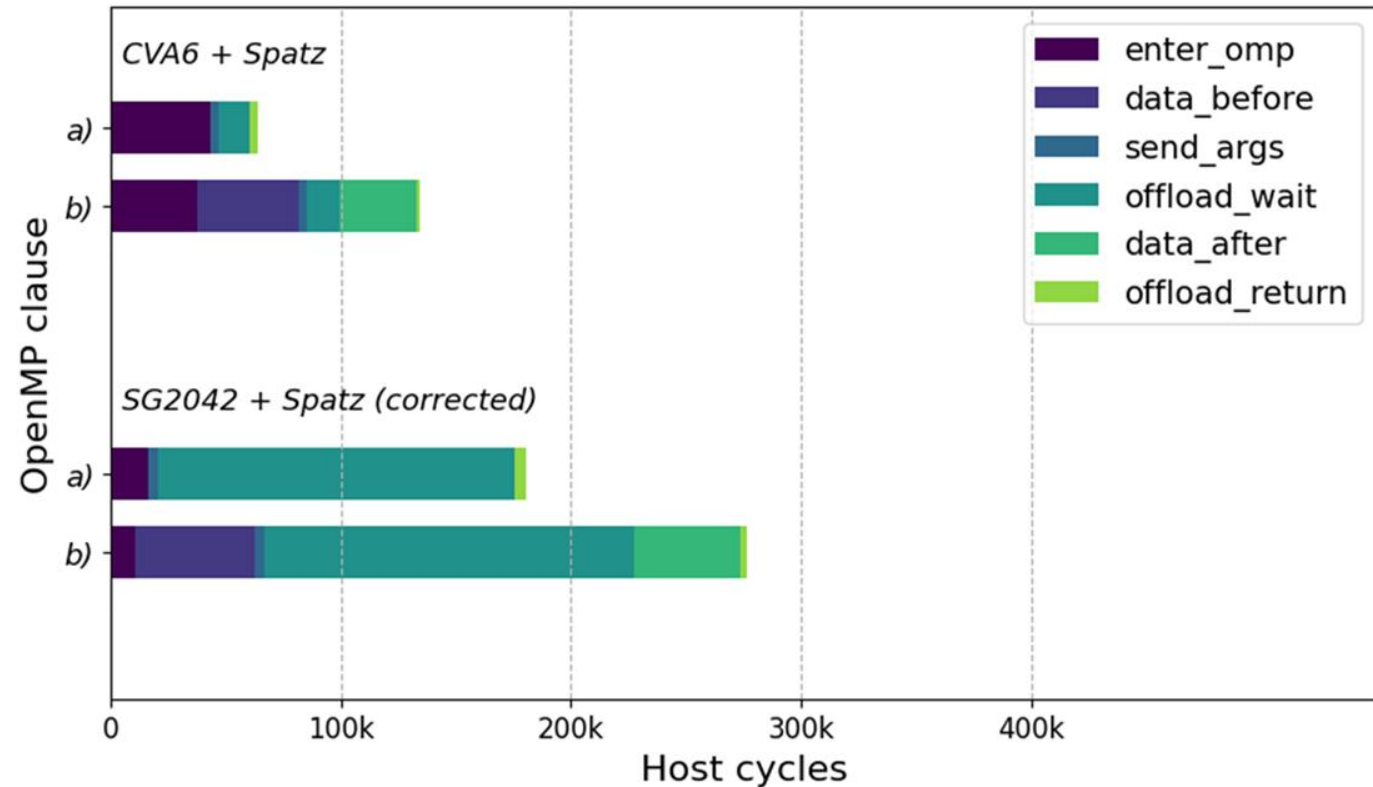
→ Faster on SG2042 (triple issue C910 core)

→ Allocate and send inputs to device

Synchronous fork/join on device

→ Requires PCIe access on SG2042

→ De-allocate



Offload benchmarks



Three micro kernels to measure offload latency

- a) No-op
- b) Input read
- c) Input read and Output write

Process kernel infos on the host

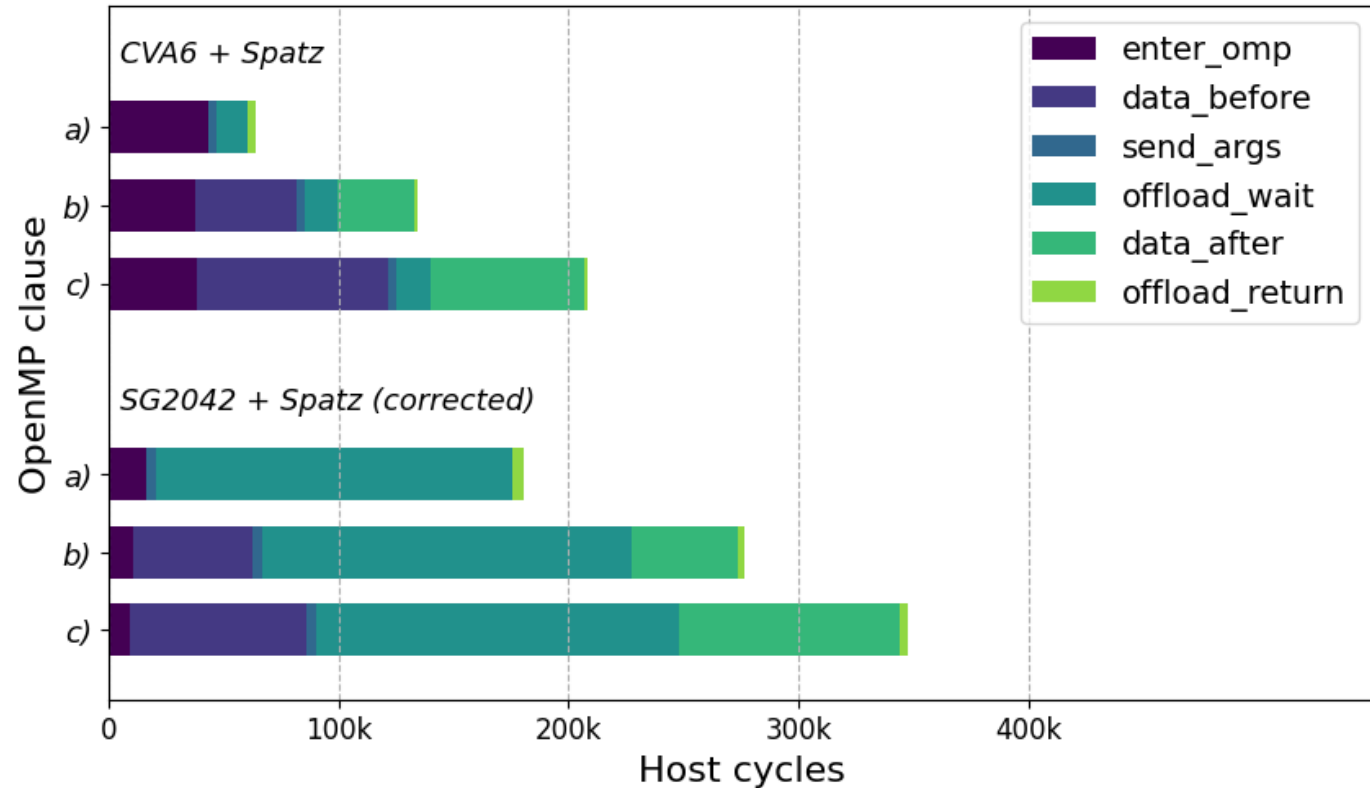
→ Faster on SG2042 (triple width C910 core)

O(100k) offload overhead
This impacts the choice of kernels to offload

Synchronous fork/join on device

→ Requires PCIe access on SG2042

→ De-allocate and read outputs from device



Matrix vector toy example



- **Matrix vector multiplication**

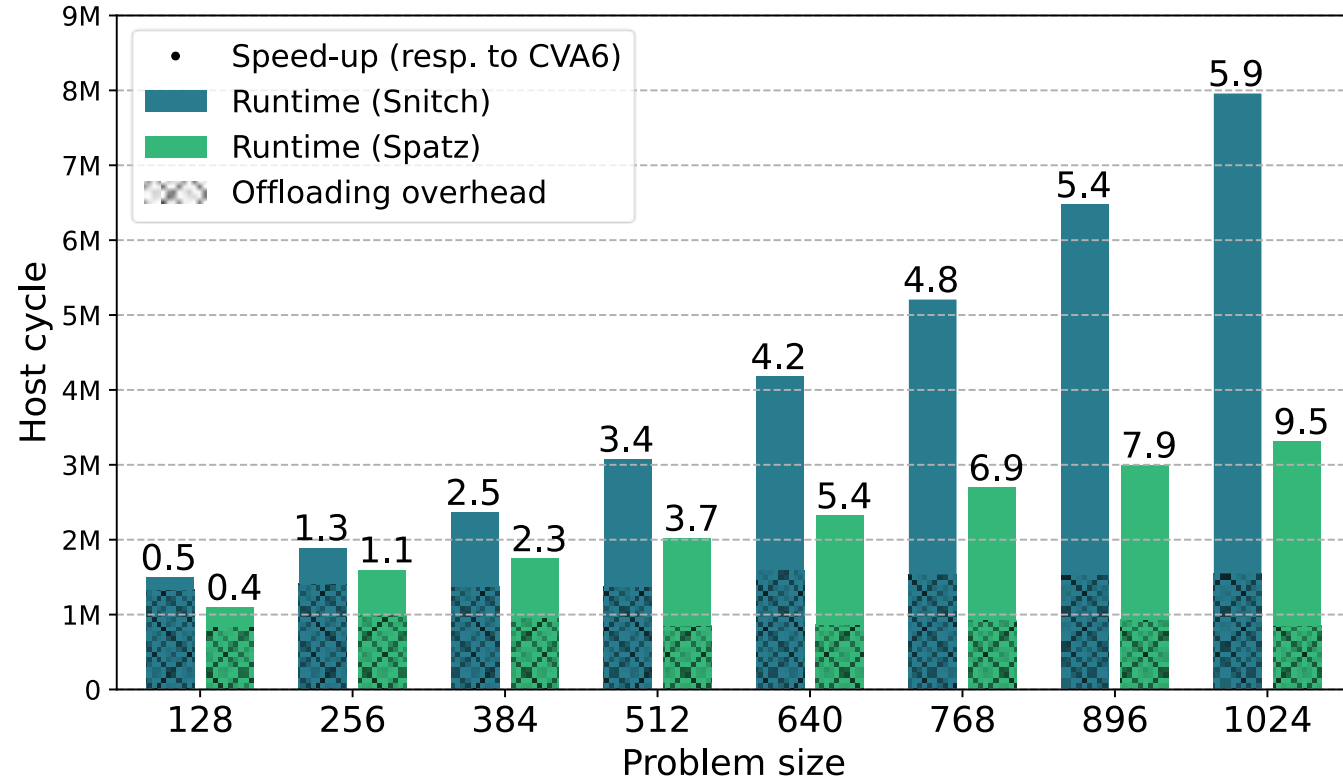
- **5.9× speedup**
19 % offload-to-work ratio

(Snitch)

- **9.5× speedup**
25 % offload-to-work ratio

(Spatz)

Accelerated matrix-vector multiplication runtimes



RISC-V unlocks heterogeneous programming across applications and devices!



- **HeroSDK** proposes a framework for all-RV heterogeneous programming
- It can be easily extended to new devices
- It suits embedded or performance software stack



www.github.com/pulp-platform/hero-tools

Cyril Koenig
Björn Forsberg
Luca Benini

cykoenig@iis.ee.ethz.ch
bjorn.forsberg@ri.se
lbenini@iis.ee.ethz.ch

Q&A

Institut für Integrierte Systeme – ETH Zürich

Gloriastrasse 35
Zürich, Switzerland

DEI – Università di Bologna

Viale del Risorgimento 2
Bologna, Italy



Demonstrated Spatz platforms



- **Spatz embedded platform**

- RV64G single core host
- 2-core RV32IMAFDV cluster
- Cluster scratchpad and DMA

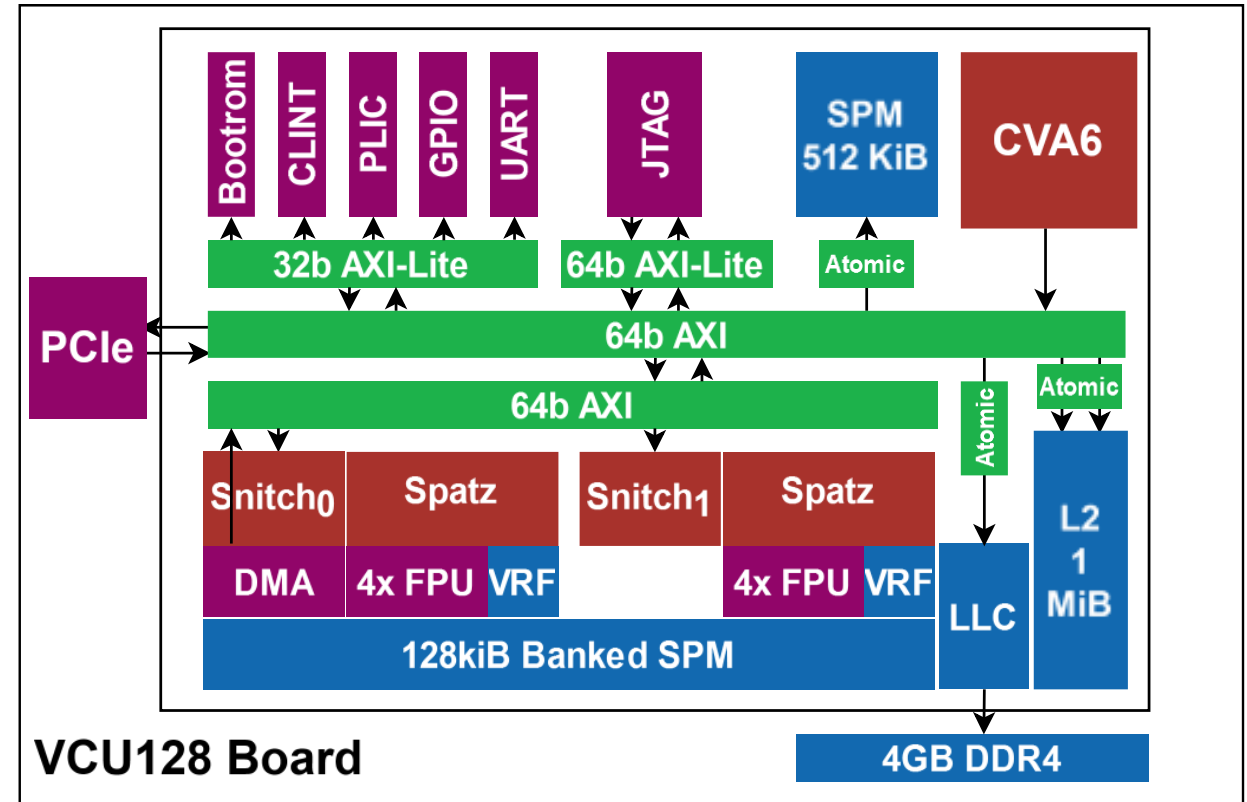
- **Spatz accelerator card**


- 2-core RV32IMAFDV cluster
- Cluster scratchpad and DMA
- PCIe endpoint

+

- **SG2042 RISC-V CPU**

- RV64G 64-cores host

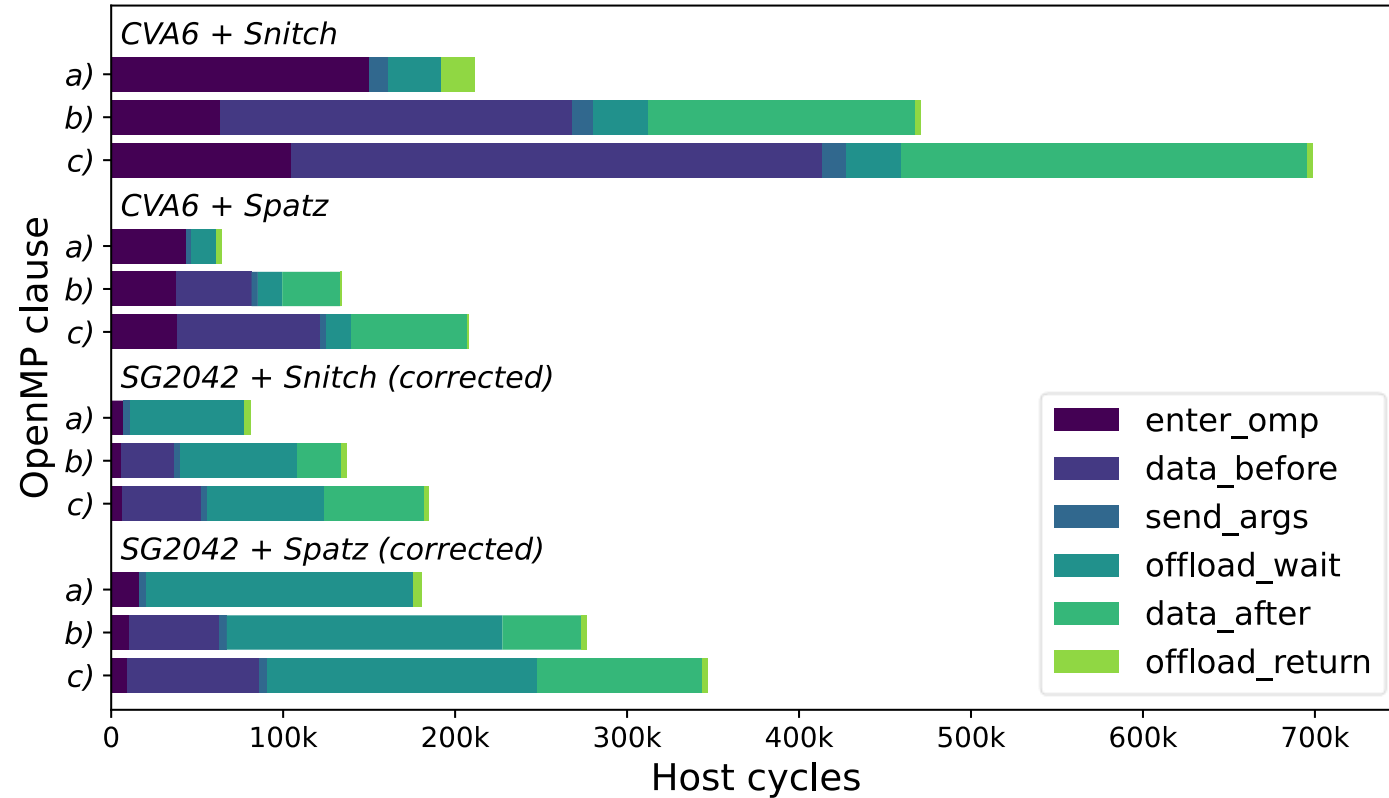


 www.github.com/pulp-platform/carfield

Offload benchmark



```
1 #pragma omp target device(1)
2   asm volatile ("nop");
3
4 #pragma omp target device(1) map(to: a)
5   volatile uint32_t b = a;
6
7 #pragma omp target device(1) map(to: a) map(tofrom: b)
8   b = a;
```



Example application



```
1
2 void matvec( float *W_p, float *X_p, float *Y_p,
3             unsigned int N ) {
4     // Offload to manager accelerator core
5     #pragma omp target device(1) map(to: W_p, X_p, Y_p, N)
6     {
7         // Allocate and move input data to local memory
8         X_l1 = snitch_l1_alloc(N);
9         W_rows_l1 = snitch_l1_alloc (CORES * N);
10        snitch_dma_memcpy (X_l1, X_p, N)
11        // Tile matrix rows
12        for ( offset = 0; offset += CORES; offset < N ) {
13            snitch_dma_memcpy(W_rows_l1, X_p[offset], N*CORES)
14            snitch_dma_wait();
15            // Execute in parallel
16            #pragma omp parallel for
17            for ( int p = 0; p < CORES; p++ )
18                Y_p[offset+p] = dotp(X_l1, &W_rows_l1[p*n], N)
19        }
20    }
21    return;
22 }
```


Heterogeneous compiler

