

# CIFER: A Cache-Coherent 12nm 16mm<sup>2</sup> SoC with Four 64-Bit RISC-V Application Cores, 18 32-Bit RISC-V Compute Cores, and a 1541 LUT6/mm<sup>2</sup> Synthesizable eFPGA

Ang Li\*, Ting-Jung Chang<sup>§1</sup>, Fei Gao\*, Tuan Ta<sup>†</sup>, Georgios Tziantzioulis\*, Yanghui Ou<sup>†</sup>, Moyang Wang<sup>†</sup>, Jinzheng Tu\*, Kaifeng Xu\*, Paul Jackson\*, August Ning\*, Grigory Chirkov\*, Marcelo Orenes-Vera\*, Shady Agwa<sup>¶2</sup>, Xiaoyu Yan<sup>†</sup>, Eric Tang<sup>†</sup>, Jonathan Balkind<sup>‡</sup>, Christopher Batten<sup>†</sup>, David Wentzloff\*  
 \*Princeton University, USA, <sup>§</sup>National Cheng Kung University, Taiwan, <sup>†</sup>Cornell University, USA, <sup>¶</sup>The University of Edinburgh, UK, <sup>‡</sup>University of California, Santa Barbara, USA

**Abstract**—This paper presents CIFER, the world’s first open-source, fully cache-coherent, heterogeneous many-core, CPU-FPGA SoC. The 12nm, 16mm<sup>2</sup> chip integrates four 64-bit, OS-capable, RISC-V application cores; three TinyCore clusters that each contain six 32-bit, RISC-V compute cores (18 in total); and an EDA-synthesized, standard-cell-based eFPGA. CIFER enables the decomposition of real-world applications and tailored execution (parallelization or specialization) per decomposed task. Our evaluation shows that: 1) the TinyCore clusters increase the throughput and energy efficiency of data- and thread-parallel tasks by up to 7.95× and 7.75× over one 64-bit core, respectively; 2) the eFPGA increases the throughput and energy efficiency of hardware-accelerable tasks by up to 9.29× and 10.62×, respectively; 3) using coherent caches for data transfer between the processors and the eFPGA increases the throughput and energy efficiency by up to 11.1× and 10.5×, respectively.

## I. INTRODUCTION

The drive for performance and energy efficiency in the post-Moore era has given rise to hardware acceleration and heterogeneous integration. However, the high design cost and programming complexity impede the broad adoption of heterogeneous system-on-chips (SoC).

This work presents **CIFER** [1] (Fig. 1), the world’s first open-source, fully cache-coherent, heterogeneous many-core, CPU-FPGA SoC. By integrating OS-capable processors, parallel compute cores, and an embedded FPGA (eFPGA), CIFER enables efficient execution of various workloads across the parallelism-specialization spectrum.

CIFER lowers the design cost and the programming barrier with the following novelties. First, it demonstrates agile hardware development facilitated by open-source hardware. CIFER was designed in seven months during the pandemic by a team of graduate students and postdocs collaborating across two institutions, due in part to the use of many open-source projects, including OpenPiton [2], BYOC [3], PyMTL3 [4], PyOCN [5], Ariane [6], and PRGA [7]. Second, the eFPGA is synthesized with off-the-shelf electronic design automation (EDA) tools and standard cell libraries. Compared to the conventional, full-custom FPGAs, CIFER’s synthesizable eFPGA is customizable in architecture, technology-agnostic, and

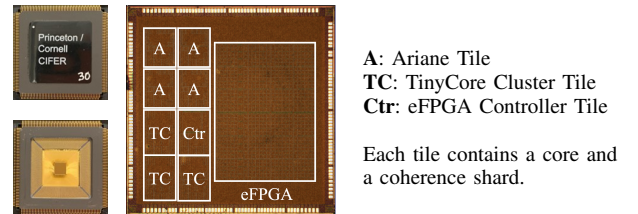


Fig. 1: CIFER Package and Die Photos

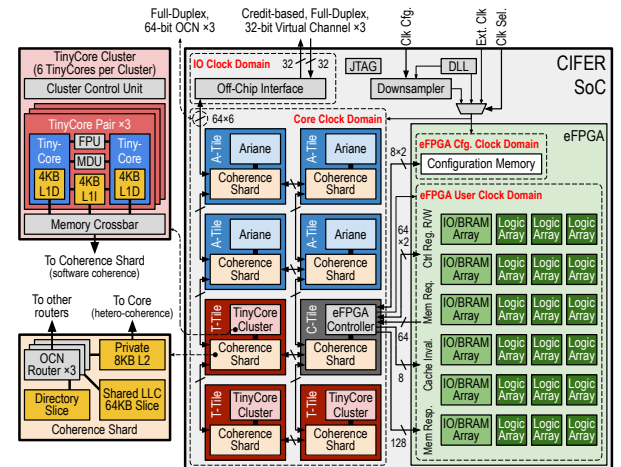


Fig. 2: CIFER SoC Architecture

flexible in physical layout. Third, CIFER implements different cache coherence schemes that are optimal for each processing unit, and unifies them within a global, bi-directionally coherent cache system.

## II. ARCHITECTURE

The CIFER architecture (Fig. 2) integrates a 2×4 mesh of tiles and an eFPGA into the distributed, coherent, OpenPiton [2] P-Mesh cache system over three packet-switched, on-chip networks (OCN) designed with PyOCN [5]. Each tile consists of a shard of the coherence system and one of the following: an Ariane core, a TinyCore cluster, or an eFPGA controller. Each coherence shard contains a private, 8KB, L2 cache and a 64KB slice of the shared, 512KB, last-level cache (LLC). Coherence between the L2 caches and the LLC is maintained in hardware with a directory-based MESI protocol.

<sup>1</sup> This work was done while Ting-Jung Chang was at Princeton.

<sup>2</sup> This work was done while Shady Agwa was at Cornell.

### A. Ariane: OS-Capable Processor

Ariane [6] is an OS-capable, 64-bit, RISC-V processor with a six-stage in-order pipeline, a 16KB L1 instruction (L1I) cache, an 8KB L1 data (L1D) cache, and a double-precision floating-point unit (FPU). Coherence between Ariane’s L1 caches and the L2 cache is maintained in hardware through adaptation to BYOC’s Transaction Response Interface (TRI) [3]. CIFER is the first silicon instantiation of BYOC.

### B. TinyCore Cluster: Thread-Level Parallel Array

Each TinyCore cluster contains six 32-bit, RISC-V cores organized into three pairs. The six cores use a MIMD execution model, where each core executes an independent stream of instructions. Each core has a six-stage, in-order issue, out-of-order write-back, late-commit, scalar pipeline. To address write-after-write and write-after-read hazards during out-of-order execution, each core supports limited register renaming with more physical registers (40 integer and 40 floating-point) than the 32 architectural registers specified in the RISC-V ISA.

Each core has a private, 4KB L1D cache, while a pair of cores share a 4KB L1I cache, an integer multiply-divide unit (MDU), and a single-precision FPU. A small L0 instruction buffer is added to each core’s front-end to minimize the latency impact of sharing the L1I cache. Coherence between the L1D caches and the L2 cache is managed explicitly in software by inserting special cache flush and invalidation instructions. In particular, a cache flush traverses the L1D cache to write back each dirty cache line, while a cache invalidation clears the valid bits of the clean cache lines. Cache invalidation requests from the L2 cache are not propagated to the L1D caches. Sharing long-latency arithmetic units and reducing coherence hardware maximize computation density in each cluster.

### C. Embedded FPGA: Reconfigurable Hardware Accelerator

The eFPGA (Fig. 3) is designed with PRGA [7]. It has 6720 multi-mode LUT6s and 18 24Kbit, dual-port, block RAMs (BRAM). Hard-wired adder/carry chains are used for efficient emulation of arithmetic operations. The BRAMs support different word sizes, e.g.,  $512 \times 48b$ ,  $1024 \times 24b$ , ...,  $24K \times 1b$ . eFPGA-emulated, “soft” accelerators can be built with an open-source, RTL-to-bitstream toolchain consisting of Yosys [8], VPR [9], and PRGA’s bitstream assembler.

The eFPGA contains three key novelties: First, the switch blocks implement a cycle-free connection pattern [10], facilitating automated, constraint-driven, area and timing optimization at the array level using off-the-shelf EDA tools. In comparison to the conventional FPGA design flow in which locally optimized blocks are tessellated in a predefined grid, this approach improves the power, performance, and area (PPA) of the synthesized FPGA by letting the EDA tools explore a larger design space. Second, the eFPGA is designed as a three-level hierarchy to balance PPA optimization and EDA runtime. The eFPGA is partitioned into two types of sub-arrays, namely logic array and IO/BRAM array, which are then composed of logic blocks. Third, the eFPGA uses a hierarchical configuration network clocked by a multi-source clock mesh. This enables fast and partial reconfiguration of the eFPGA at GHz clock frequency. In particular, each sub-array contains a bitstream router and a single-bit scanchain

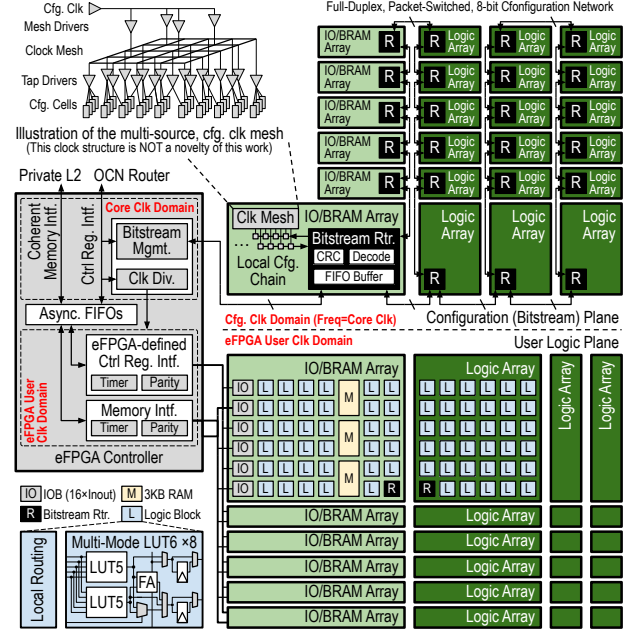


Fig. 3: CIFER eFPGA Architecture

that connects all the configuration cells with minimum routing metal usage. Bitstream segments are first sent to the bitstream routers over an 8-bit, packet-switched network, then buffered and shifted into the scanchains.

The eFPGA is integrated with the system through the eFPGA controller, the first silicon instantiation of Duet [11], which contains the following two interfaces. The control register interface allows the processors to access the eFPGA via memory-mapped I/O. The coherent memory interface is configurable at runtime to enable non-coherent, IO-coherent, or bi-directionally coherent memory accesses of the eFPGA. In bi-directionally coherent mode, cache invalidation requests from the L2 cache are forwarded into the eFPGA, allowing the accelerator to include a “soft” cache. Atomic requests from the eFPGA are also supported, enabling low-overhead synchronization in user mode. Both interfaces contain asynchronous FIFOs for clock domain crossing and are equipped with timers and parity checks to protect the OCN and the memory system from software or accelerator bugs.

### D. Heterogeneous Cache Coherence

One key contribution of CIFER is that it unifies the heterogeneous cache coherence schemes of each processing unit within a global, fully coherent cache system. This minimizes the communication overhead and maximizes the programmability of the SoC. For example, a task-parallel work-stealing runtime [12] facilitates parallel execution across the Ariane cores and the TinyCore clusters, leveraging the coherent caches and automating the insertion of cache management instructions. An eFPGA-emulated accelerator can be efficiently invoked by passing the memory addresses of the data to be processed. Depending on the computation, the accelerator can either copy a continuous chunk of data into its BRAM scratchpad or read/write memory in a random, byte-granular manner. This saves CPU cycles from explicitly managing data movement and prevents over-fetching from the eFPGA.

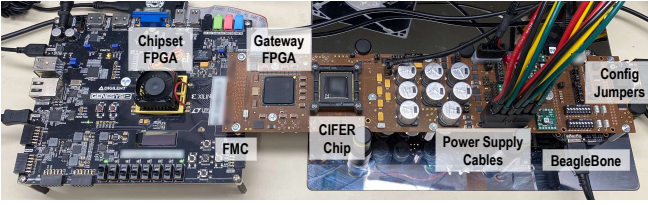


Fig. 4: Lab Evaluation Setup

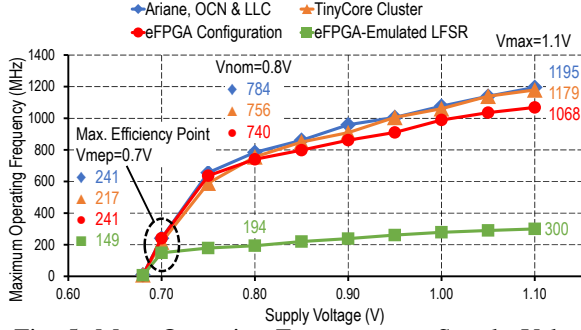
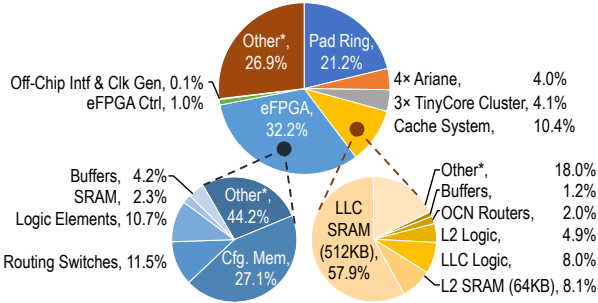


Fig. 5: Max. Operating Frequency vs. Supply Voltage



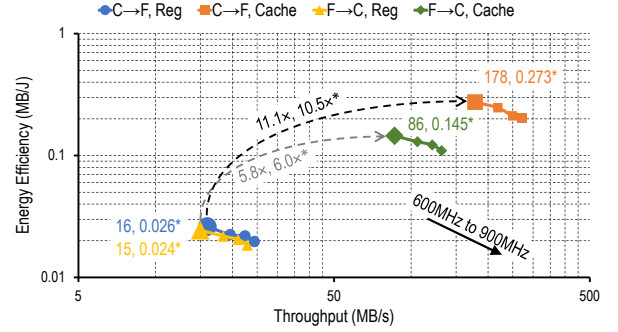
\* **Other** includes physical cells (e.g., tap cells, decap cells, boundary cells), gap areas between hard macro blocks, etc.

Fig. 6: Area Breakdown

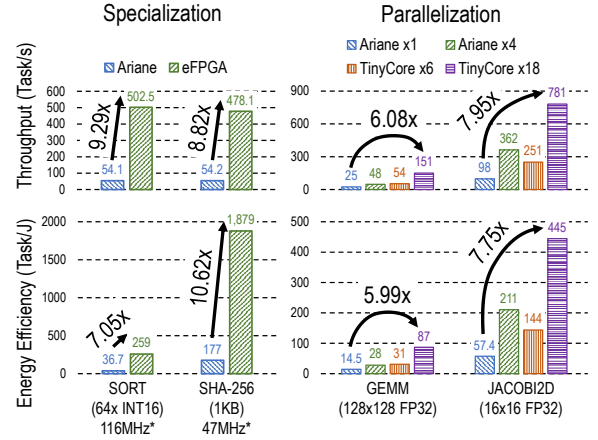
### III. EVALUATION

Fig. 4 shows our chip testing setup. Fig. 5 shows each component’s maximum operating frequency ( $F_{\max}$ ) across the range of functional supply voltages. The eFPGA’s  $F_{\max}$  depends on the emulated design, and Fig. 5 shows the  $F_{\max}$  of a 64-bit LFSR. Fig. 6 shows the area breakdown of the chip. The eFPGA’s logic and routing resources only make up a quarter of the eFPGA’s total area, while the configuration memory consumes another quarter. The eFPGA’s low area utilization is due in part to the hierarchical design and can be improved with abutted or narrow-channel macro-placement strategies.

Table I compares CIFER with other state-of-the-art CPU-FPGA SoCs targeting the edge/IoT domain. Due to tooling issues, we did not implement explicit clock-gating on the eFPGA’s configuration clock, which should be disabled except when loading the bitstream. Post-layout power analysis shows that the configuration clock sub-tree consumes about 90% of the chip’s total clock power due to the high total capacitance and short-circuit current of the clock meshes. We estimate the total power with proper clock-gating by subtracting the analyzed configuration clock power from the measured total power. Estimated numbers are shown in brackets, next to their measured counterparts in the table.



\* Throughput and energy efficiency when the processors run at 600MHz. The eFPGA runs at 1/16 of the CPUs’ clock frequency. C=CPU; F=eFPGA; Reg=memory-mapped I/O; Cache=coherent cache. Fig. 7: CPU-FPGA Communication Throughput and Energy Efficiency at Different System Clock Frequency



\*  $F_{\max}$  of the eFPGA-emulated accelerator. The processors, OCN, cache system, and the eFPGA controller runs at full speed (740MHz at 0.8V).

Fig. 8: Performance and Efficiency Gains from Offloading

CIFER runs up to 1195MHz at 1.1V. The processors provide high aggregate performance with good energy efficiency, totaling 15.54 GFLOPS at 1.1V and 53.18 GFLOPS/W (estimated as explained above) at 0.7V, outperforming the next best SoC by 8.0 $\times$  and 1.4 $\times$ . The eFPGA’s area efficiency is 1541 LUT6/mm<sup>2</sup>, outperforming the other synthesizable eFPGAs by 11.2 $\times$ , and is only 1.3 $\times$  worse than the best full-custom eFPGA. The eFPGA’s peak performance (1.92 MOPS/LUT, 126MHz at 1.1V) and energy efficiency (148.1 GOPS/W at 0.7V) are measured with a 64-point FFT that utilizes 97% of the logic blocks and 75% of the BRAMs. The 3.4 $\times$  performance gap and the 2.1 $\times$  energy efficiency gap between the full-custom eFPGA and this work can be attributed to three factors: (1) CIFER is synthesized with standard cells; (2) our eFPGA has no hardware multiply-accumulate units; and (3) this work uses an open-source RTL-to-bitstream toolchain.

Fig. 7 shows the throughput improvements and energy savings when data are transferred through the coherent caches instead of memory-mapped I/O. The improvements are due to two reasons: (1) memory-mapped I/O accesses are strictly serialized in the processor’s pipeline, while coherent caches may hide the latency of consecutive memory accesses, e.g., by buffering memory requests in the asynchronous FIFOs; (2) the eFPGA can use the L2 cache co-located in the eFPGA controller tile which runs in the fast, processors’ clock domain.

		This Work	TCAS'20 [13]	ISSCC'19 [14]	TVLSI'21 [15]	JSSC'22 [16]		
Chip	Technology	12nm FinFET	90nm BCD	40nm CMOS + 39nm MRAM	22nm FD-SOI	16nm FinFET		
	Die Area (mm <sup>2</sup> )	16	1.78	22.09	9	25		
	V <sub>nom</sub> (V <sub>min</sub> - V <sub>max</sub> )	0.8 (0.68 - 1.1)	1.2 (-)	- (1.1 - 1.3)	0.8 (0.5 - 0.8)	0.8 (0.5 - 1.05)		
	Active Power (mW)	V <sub>nom</sub>	1792	1.2	5.34	24.95	918	
	F <sub>max</sub> (MHz)	V <sub>max</sub>	<b>1195</b>	10	200	600	972	
CPU	Host	Core Type	4× Ariane	RISCY	Cortex-M0	RISCY	2× Cortex-A53	
		ISA	RV64GC	RV32I	ARMv6-M	RV32IMFC	ARMv8-A	
		CoreMark Score	V <sub>max</sub>	<b>7918</b>	31.9	466	1914	6376
	Other	Core Type	18× TinyCore	N/A	N/A	N/A	Cortex-M0	
		ISA	RV32IMAF				ARMv6-M	
		Function	Parallel Compute				Monitor	
		CoreMark Score	V <sub>max</sub>				<b>19198</b>	2265
	Total	Peak GFLOPS	V <sub>max</sub>	<b>15.54</b>	NO HW FPU	NO HW FPU	NO HW FPU	1.94
		Peak GFLOPS/W	V <sub>mep</sub>	<b>6.63 [53.18<sup>†</sup>]</b>	NO HW FPU	NO HW FPU	NO HW FPU	38.03
		IP		<b>Synthesizable w/ Std. Cells</b>	Synthesizable w/ Std. Cells	Unknown	Full-Custom Hard Macro	Full-Custom Hard Macro
eFPGA	Min. Prog. Time (μs)		<b>239.4 - 1274.8</b>	-	-	-	450	
	LUT Type & Count		6720 LUT6	48 LUT6	1176 LUT6	6000 LUT4	8760 LUT6	
	Logic Density (LUT/mm <sup>2</sup> )		1541	137	36	1505	1991	
	F <sub>max</sub> (MHz)	V <sub>max</sub>	300**	1.25	200	193	747	
	MOPS/LUT	V <sub>max</sub>	1.92 <sup>‡</sup> (INT8)	-	-	0.02 (INT32)	6.45 (INT8)	
	GOPS/W	V <sub>mep</sub>	148.1 <sup>‡*</sup> (INT8)	-	-	29.1 (INT32)	312.4 (INT8)	
Shared Memory BW (MB/s) [V <sub>max</sub> ]		CPU → eFPGA	<b>201</b>	Non-Coherent	Non-Coherent	Non-Coherent	-	
		eFPGA → CPU	<b>558</b>				486	

<sup>†</sup> Estimated power dissipation, excluding the eFPGA's configuration clock power based on post-layout power analysis

\*\* Measured when the eFPGA emulates a 64-bit LFSR

<sup>‡</sup> Measured when the eFPGA emulates an INT8-precision, complex, 64-point FFT

\* Measured power dissipation in the eFPGA's user clock domain

TABLE I: Comparison to the State of the Art

Fig. 8 shows the throughput and energy efficiency gains by offloading four representative edge applications to their preferred compute unit. SORT and SHA-256 use eFPGA-emulated accelerators, while GEMM and JACOBI2D use the TinyCore clusters. The execution time is measured from when an Ariane core initiates a task to when the same core reads back all the results. All the control overhead is included, while the data transfer overhead is mitigated by overlapping compute with ad hoc, coherent memory accesses. At nominal voltage (0.8V), the eFPGA outperforms the Ariane-only baseline by up to 9.29× in throughput and 10.62× in energy efficiency; the TinyCore clusters improve the performance and energy efficiency by up to 7.95× and 7.75×, respectively.

#### IV. CONCLUSION

This paper presents CIFER. Through cache-coherent integration of OS-capable processors, parallel many-core arrays, and an eFPGA, CIFER improves performance and energy efficiency on a wide range of workloads across the parallelism-specialization spectrum. The heterogeneous cache coherence scheme minimizes communication overhead and maximizes the programmability of the SoC. The EDA-synthesized, standard-cell-based eFPGA's area efficiency, peak performance, and energy efficiency are approaching those of full-custom eFPGAs.

#### ACKNOWLEDGMENTS

This material is based on research sponsored by the Air Force Research Laboratory (AFRL) and Defense Advanced Research Projects Agency (DARPA) under agreement No. FA8650-18-2-7852. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of Air Force Research Laboratory (AFRL) and Defense Advanced Research Projects Agency (DARPA) or the U.S. Government.

#### REFERENCES

- [1] T.-J. Chang *et al.*, "CIFER: A 12nm, 16mm<sup>2</sup>, 22-Core SoC with a 1541 LUT6/mm<sup>2</sup> 1.92 MOPS/LUT, Fully Synthesizable, CacheCoherent, Embedded FPGA," in *CICC '23*.
- [2] J. Balkind *et al.*, "OpenPiton: An Open Source Manycore Research Framework," in *ASPLOS '16*, pp. 217–232.
- [3] —, "BYOC: A "Bring Your Own Core" Framework for Heterogeneous-ISA Research," in *ASPLOS '20*, p. 699–714.
- [4] S. Jiang *et al.*, "Mamba: Closing the Performance Gap in Productive Hardware Development Frameworks," in *DAC '18*, pp. 1–6.
- [5] C. Tan *et al.*, "PyOCN: A Unified Framework for Modeling, Testing, and Evaluating On-Chip Networks," in *ICCD '19*, pp. 437–445.
- [6] F. Zaruba and L. Benini, "The Cost of Application-Class Processing: Energy and Performance Analysis of a Linux-Ready 1.7-GHz 64-Bit RISC-V Core in 22-nm FDSOI Technology," *IEEE TVLSI*, vol. 27, no. 11, pp. 2629–2640, 2019.
- [7] A. Li and D. Wentzlauff, "PRGA: An Open-Source FPGA Research and Prototyping Framework," in *FPGA '21*, p. 127–137.
- [8] C. Wolf, "Yosys open synthesis suite," <https://yosyshq.net/yosys/>, 2020.
- [9] K. E. Murray *et al.*, "VTR 8: High-Performance CAD and Customizable FPGA Architecture Modelling," *ACM TRETS*, vol. 13, no. 2, 2020.
- [10] A. Li *et al.*, "Automated Design of FPGAs Facilitated by Cycle-Free Routing," in *FPL '20*, pp. 208–213.
- [11] —, "Duet: Creating Harmony between Processors and Embedded FPGAs," in *HPCA '23*, pp. 745–758.
- [12] M. Wang *et al.*, "Efficiently Supporting Dynamic Task Parallelism on Heterogeneous Cache-Coherent Systems," in *ISCA '20*, pp. 173–186.
- [13] F. Renzini *et al.*, "A Fully Programmable eFPGA-Augmented SoC for Smart Power Applications," *IEEE TCAS-I*, vol. 67, no. 2, pp. 489–501, 2020.
- [14] M. Natsui *et al.*, "12.1 An FPGA-Accelerated Fully Nonvolatile Microcontroller Unit for Sensor-Node Applications in 40nm CMOS/MTJ-Hybrid Technology Achieving 47.14μW Operation at 200MHz," in *ISSCC '19*, pp. 202–204.
- [15] P. D. Schiavone *et al.*, "Arnold: An eFPGA-Augmented RISC-V SoC for Flexible and Low-Power IoT End Nodes," *IEEE TVLSI*, vol. 29, no. 4, pp. 677–690, 2021.
- [16] S. K. Lee *et al.*, "SMIV: A 16-nm 25-mm<sup>2</sup> SoC for IoT With Arm Cortex-A53, eFPGA, and Coherent Accelerators," *IEEE JSSC*, vol. 57, no. 2, pp. 639–650, 2022.