

Hardware Cache Coherence Analysis in Zynq UltraScale+ MPSoC: Characterization of CCI-400 Performance and APU-RPU Communication

Luigi Caucci, Simone Cecere

Department of Electrical Engineering and Information Technology
University of Naples Federico II
Naples, Italy

Email: {l.caucci, simo.cecere}@studenti.unina.it

Abstract—The efficient management of cache coherence in heterogeneous multiprocessor systems remains a critical challenge for embedded system designers. This paper presents a comprehensive experimental analysis of the ARM CoreLink CCI-400 Cache Coherent Interconnect within the AMD Xilinx Zynq UltraScale+ MPSoC architecture, focusing specifically on Application Processing Unit to Real-time Processing Unit communication scenarios. Through rigorous experimental validation on the Kria KR260 platform, we develop a hybrid Linux-based testing framework that combines the Remoteproc subsystem with kernel module instrumentation to definitively characterize cache coherence behavior and quantify software cache management overhead. Our stress testing methodology, comprising 100,000 consecutive read operations, demonstrates complete cache incoherence with 100% stale data observations under standard Board Support Package configurations. We precisely measure cache invalidation overhead ranging from 1.66 μ s for single cache lines to 63.4 μ s for 1024 cache lines, exhibiting linear scaling at approximately 60 ns per 64-byte cache line. Based on validated theoretical models with prediction errors consistently below 3%, we estimate potential performance improvements of 4–10 \times with functional hardware coherence. Through systematic investigation of the complete configuration prerequisite chain, we identify three critical missing components that prevent CCI-400 activation. Through Tightly Coupled Memory baseline characterization at 0.61 μ s write latency, we validate our theoretical models and confirm 5.7 \times cache overhead that CCI-400 would eliminate. TCM measurements isolate 2.89 μ s pure cache management cost, proving CCI-400 represents the scalable coherence solution for DDR-based communication that prevent CCI-400 activation on commercial platforms: PMUFW broadcast configuration macros, First Stage Bootloader snoop port initialization, and Vivado hardware design parameters. This work provides both a rigorous characterization of non-coherent baseline performance and a validated experimental methodology applicable to cache coherence verification in heterogeneous MPSoC systems.

Index Terms—Cache coherence, CCI-400, Zynq UltraScale+ MPSoC, heterogeneous systems, ARM Cortex-A53, ARM Cortex-R5F, Remoteproc, real-time systems, APU-RPU communication

I. INTRODUCTION

Modern heterogeneous embedded systems integrate diverse processing elements with fundamentally different architectural characteristics and operational requirements. The AMD Xilinx Zynq UltraScale+ MPSoC exemplifies this architectural

paradigm by combining a quad-core ARM Cortex-A53 Application Processing Unit operating at 1.33 GHz under Linux, a dual-core ARM Cortex-R5F Real-time Processing Unit operating at 533 MHz for deterministic tasks, programmable logic fabric, and the ARM CoreLink CCI-400 Cache Coherent Interconnect [1]–[4]. When these heterogeneous processing elements access shared memory regions, maintaining cache coherence becomes essential for ensuring both functional correctness and system performance.

Figure 1 illustrates the complete architectural organization of the Zynq UltraScale+ MPSoC, showing the relationship between the APU cluster with its hierarchical cache structure, the RPU subsystem with Tightly Coupled Memory, and the CCI-400 interconnect topology. Without proper cache synchronization mechanisms, a fundamental incoherence scenario emerges: the APU writes data that remains resident in its L2 cache due to write-back policies, while the RPU simultaneously reads from main DDR memory and obtains stale values. This situation compromises both functional correctness and real-time determinism guarantees required for safety-critical applications [5], [6].

Three primary architectural approaches address cache coherence challenges in heterogeneous embedded systems. Software-managed coherence employs explicit cache flush and invalidate operations through CP15 coprocessor instructions, requiring minimal hardware support but introducing significant performance overhead and demanding rigorous programming discipline [2]. Alternatively, designating shared memory regions as non-cacheable eliminates coherence concerns entirely by forcing all accesses directly to main memory, but this approach results in memory access latencies typically 10–100 \times slower than cached operations [1]. The third approach, hardware-managed coherence through the CCI-400 interconnect, implements the AMBA ACE protocol to provide transparent hardware-managed coherence via snoop-based mechanisms, offering optimal performance characteristics but requiring complex multi-layer system configuration and proper initialization sequences [4], [7].

Despite the theoretical advantages offered by hardware-managed coherence mechanisms, practical implementation on

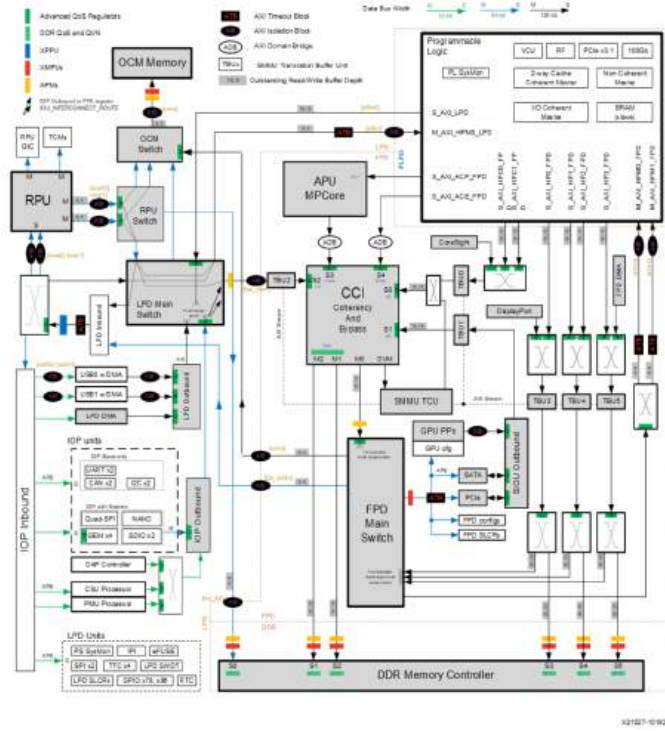


Fig. 1: Zynq UltraScale+ MPSoC architecture illustrating the APU (quad Cortex-A53 with L1/L2 cache hierarchy), RPU (dual Cortex-R5F with Tightly Coupled Memory), and CCI-400 interconnect topology. Source: AMD Xilinx UG1085 [1].

commercial embedded platforms presents substantial challenges. The documentation landscape remains fragmented: the Zynq Technical Reference Manual [1] provides comprehensive architectural specifications but lacks concrete end-to-end configuration examples, while ARM’s CCI-400 Technical Reference Manual [4] thoroughly describes the interconnect without addressing SoC-specific integration requirements. Most existing research and application notes either utilize exclusively the APU or exclusively the RPU, or resort to non-cacheable memory mappings for inter-processor communication [8], thereby avoiding the coherence problem rather than solving it.

This work systematically addresses these documentation and implementation gaps through comprehensive experimental investigation. Our principal contributions include the first systematic analysis of CCI-400 configuration challenges on commercial Kria KR260 platforms utilizing standard Board Support Package components, development of a rigorous kernel module-based testing methodology that eliminates false positives inherent in userspace memory mapping approaches, definitive coherence verification through 100,000-iteration stress testing demonstrating complete cache incoherence, quantitative performance characterization measuring cache invalidation overhead at approximately 60 ns per 64-byte cache

line, validated theoretical performance models predicting 4–10 \times speedup potential with functional CCI-400 hardware coherence, complete prerequisite documentation identifying three critical missing requirements, and creation of a reusable experimental framework with nanosecond-resolution timing capabilities applicable to broader heterogeneous computing research.

The remainder of this paper proceeds as follows. Section II examines related work in cache coherence for heterogeneous systems and identifies critical documentation gaps. Section III presents comprehensive background on Zynq architecture, cache coherence protocols, and CCI-400 operational principles. Section IV describes our iterative research methodology and the evolution from bare-metal approaches to successful Linux-based implementations. Section V details system implementation including device tree configuration, kernel module development, and RPU firmware structure. Section VI presents our experimental framework incorporating shared timer synchronization and definitive coherence verification methodology. Section VII reports comprehensive experimental results including stress testing outcomes and quantitative performance measurements. Section VIII provides theoretical analysis predicting expected coherent scenario performance. Section IX discusses methodology validation, limitations, and practical implications. Section X concludes with future research directions and broader community impact.

II. RELATED WORK

The intersection of cache coherence, heterogeneous computing, and embedded systems has received considerable attention in both academic literature and industrial practice, yet significant gaps remain in practical implementation guidance for commercial platforms.

A. Cache Coherence in Heterogeneous Systems

Cache coherence protocols for homogeneous multiprocessor systems represent a mature field with well-established theoretical foundations and practical implementations. The classical MESI protocol and its variants have been extensively studied and deployed in systems ranging from high-performance computing to mobile processors [5]. However, heterogeneous systems introduce additional complexity: processors with different instruction set architectures, distinct cache hierarchies, and fundamentally different operational characteristics must maintain coherent views of shared memory while operating in different clock domains.

Research examining ARM big.LITTLE coherence mechanisms provides valuable insights into coherence maintenance across asymmetric processor clusters, but these studies target fundamentally different SoC implementations with distinct boot sequences and firmware requirements than the Zynq architecture. Studies of cache coherence from security perspectives, including cache timing attacks and side-channel analysis, demonstrate important security implications but do not address functional correctness requirements in heterogeneous embedded systems. Performance analyses utilizing CCI-

400 interconnects focus primarily on performance monitoring counter capabilities rather than comprehensive configuration prerequisite documentation.

B. Zynq Platform Research

Academic research utilizing Zynq UltraScale+ platforms frequently employs either APU-exclusive or RPU-exclusive configurations. Published work addressing inter-processor communication commonly resorts to non-cacheable memory region mappings [8], thereby circumventing coherence challenges entirely rather than addressing them directly. This approach, while pragmatic for many applications, leaves significant performance potential unrealized and fails to provide guidance for scenarios where hardware coherence activation is essential for meeting real-time deadlines or power budgets.

Industrial application notes and whitepapers typically provide high-level architectural overviews without the detailed prerequisite analysis necessary for successful implementation. Community forums contain extensive discussions of coherence-related topics but exhibit characteristic patterns: many coherence configuration questions remain unanswered or receive incomplete responses, proposed solutions frequently suggest manual cache flush and invalidate operations rather than hardware coherence enablement, device tree configuration snippets are shared without experimental verification of functionality, and prerequisite dependency chains are rarely made explicit or systematically documented [9].

C. Documentation Analysis

AMD Xilinx documentation resources, including the Technical Reference Manual [1] and Software Developers Guide [10], collectively exceed 2000 pages and provide comprehensive architectural specifications. However, these documents lack practical end-to-end configuration examples for cache coherence scenarios. The TRM thoroughly describes individual component specifications, while the SDG covers Remoteproc framework usage and boot flow sequences but provides minimal guidance specific to CCI-400 configuration. Online wiki resources [9] offer fragmented device tree code snippets without adequately addressing prerequisite dependencies or validation procedures.

ARM technical documentation [2]–[4] thoroughly documents individual IP components but deliberately avoids SoC-specific integration details, which fall under the SoC vendor’s documentation responsibility. The CCI-400 TRM comprehensively describes register interfaces and protocol behavior while implicitly assuming that firmware-layer prerequisites are properly handled by bootloader implementations.

D. Identified Gaps

Our comprehensive analysis reveals several critical gaps in the available documentation ecosystem. First, no documentation comprehensively combines Linux APU operation, Remoteproc-based RPU management, and CCI-400 configuration with experimental validation procedures. Second, relationships and dependencies between PMUFW configuration,

FSBL initialization, and device tree specifications are not clearly documented or systematically analyzed. Third, no established procedures exist for definitively verifying whether hardware coherence mechanisms are functionally active beyond simple functional testing that may produce false positives. Fourth, CCI-400 Performance Monitoring Unit behavior within Zynq UltraScale+ is not documented in vendor materials. Finally, limitations of commercial development platforms like Kria regarding coherence enablement are not explicitly stated in product documentation.

This work systematically addresses these identified gaps through comprehensive implementation, rigorous testing methodology development, and detailed documentation of findings applicable to the broader embedded systems community.

III. BACKGROUND

This section provides essential background on the Zynq UltraScale+ architecture, cache coherence fundamentals, and the ARM CoreLink CCI-400 interconnect that forms the foundation for understanding our experimental methodology and results.

A. Zynq UltraScale+ Architecture

The Zynq UltraScale+ MPSoC represents AMD Xilinx’s third-generation heterogeneous SoC architecture. The XCZU5EV variant deployed on the Kria KR260 platform [1] integrates multiple distinct processing subsystems designed for complementary operational roles.

The Application Processing Unit comprises four 64-bit ARM Cortex-A53 processor cores [2] operating at 1.33 GHz. Each core incorporates 32 KB L1 instruction cache, 32 KB L1 data cache organized as 4-way set-associative, and access to a shared 1 MB L2 cache organized as 16-way set-associative. The APU subsystem implements write-back cache policies by default and maintains internal cluster coherence via ACE-compliant interfaces. All cache structures operate with 64-byte cache line granularity, a fundamental parameter affecting all cache management operations.

The Real-time Processing Unit consists of two 32-bit ARM Cortex-R5F processor cores [3] operating at 533 MHz. Each core incorporates 32 KB L1 instruction cache organized as 2-way set-associative, 32 KB L1 data cache organized as 4-way set-associative, and 128 KB Tightly Coupled Memory providing deterministic single-cycle access characteristics. Notably, the RPU subsystem does not include an L2 cache, resulting in a simplified cache hierarchy particularly suitable for hard real-time applications where temporal predictability is paramount [6]. The RPU architecture operates with architectural independence from the APU and can execute either bare-metal firmware or real-time operating systems.

This heterogeneous architecture enables workload partitioning strategies where computationally intensive non-real-time tasks execute on the APU under Linux, while deterministic real-time control loops execute on the RPU with guaranteed deadline satisfaction. However, this partitioning introduces

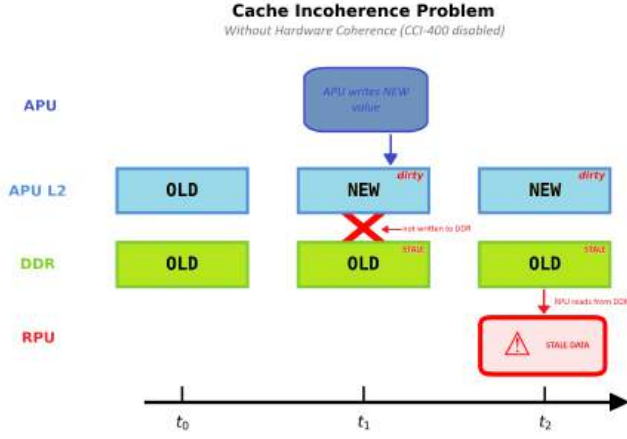


Fig. 2: Timeline illustration of cache incoherence scenario. At time t_1 , the APU writes a NEW value that remains resident in L2 cache due to write-back policy. At t_2 , the RPU reads from DDR and observes the stale OLD value, demonstrating absence of hardware cache coherence.

the fundamental challenge of maintaining coherent views of shared data structures accessed by both processing subsystems.

B. Cache Coherence Fundamentals

Cache coherence protocols ensure that all processors within a multiprocessor system maintain consistent views of shared memory locations [5]. The classical MESI protocol defines four mutually exclusive states for each cache line: Modified indicates the cache line has been modified and represents the only valid copy requiring write-back before replacement, Exclusive indicates a clean cache line representing the only cached copy that can transition to Modified without external notification, Shared indicates a clean cache line potentially present in other processor caches, and Invalid indicates the cache line contains invalid data requiring fetch from memory or another cache if accessed.

ARM's implementation extends this foundation with the MOESI variant utilized by the ACE protocol [7], introducing the Owned state to enable more efficient sharing of modified data without requiring immediate write-back to main memory. This extension provides performance benefits in scenarios where multiple processors read recently modified data before the originating processor writes back to main memory.

Figure 2 illustrates a fundamental cache incoherence scenario that occurs without proper hardware coherence mechanisms. At time t_1 , the APU writes a new value that remains resident in its L2 cache due to write-back policy without immediate propagation to DDR. Subsequently at time t_2 , the RPU reads directly from DDR and observes the stale old value. This temporal sequence definitively demonstrates the complete absence of hardware cache coherence mechanisms and motivates our experimental methodology for characterizing the performance impact of software-managed cache operations.

C. AMBA ACE Protocol

The AMBA ACE (AXI Coherency Extensions) protocol [7] extends the baseline AXI4 interconnect protocol with three additional transaction channels specifically designed for coherence maintenance. The AC channel transmits snoop requests from the interconnect to caching bus masters, the CR channel returns snoop responses containing current cache line state information from bus masters to the interconnect, and the CD channel enables direct cache-to-cache data transfer without memory subsystem involvement when one cache holds modified data required by another processor.

ACE transaction types include operations such as Read-Shared for obtaining shared copies, ReadUnique for obtaining exclusive ownership, ReadClean for obtaining clean copies, CleanShared for writing back modified data while retaining a shared copy, CleanInvalid for writing back and invalidating, and MakeInvalid for invalidating copies in other caches. These operations enable fine-grained coherence state management optimized for various access patterns.

Memory region attributes specify coherence domain behavior according to the ARMv8 architecture [11]. Outer Shareable regions maintain coherence across all system agents including external masters, Inner Shareable regions maintain coherence within specific processor clusters or domains, and Non-shareable regions provide no hardware coherence guarantees requiring software to manage consistency explicitly.

D. CCI-400 Architecture

The ARM CoreLink CCI-400 Cache Coherent Interconnect [4] implements the ACE protocol within the Zynq UltraScale+ architecture at base address 0xFD6E0000. The interconnect architecture comprises up to five ACE-compliant slave interfaces for coherent bus masters such as APU clusters, RPU subsystems, and GPU units, two ACE-Lite master interfaces connecting to DDR memory controllers and peripheral subsystems, a centralized Snoop Control Unit responsible for coordinating all coherence transactions, integrated hardware performance monitoring counters for observing coherence-related events, and support for multiple coherence domain configurations with shareability attribute enforcement.

When a coherent bus master initiates a memory access transaction, a precisely defined sequence occurs. The memory request enters through an ACE slave interface port designated S0 through S4. The Snoop Control Unit analyzes the transaction type and target address to determine coherence requirements. The SCU determines which other bus masters may cache the requested address based on address range tracking. Snoop requests are broadcast to all relevant slave interface ports connected to caching agents. Each snooped master examines its cache tags and responds with cache line state information indicating Hit or Miss combined with Clean or Dirty status. If modified data exists in another cache, direct cache-to-cache transfer occurs via the CD channel, bypassing main memory entirely. The final coherent response is assembled and returned to the requesting master. Memory transactions to DDR proceed only if necessary, specifically

when no cache holds a valid copy or when write-back to memory is required for consistency.

An important architectural limitation of the CCI-400 concerns its snoop distribution mechanism. The CCI-400 employs broadcast-based snooping, meaning every coherence transaction must query all connected bus masters regardless of whether they actually cache the requested address. This architecture scales adequately for small numbers of masters, typically two to five agents, but exhibits efficiency degradation beyond this range due to unnecessary snoop traffic. Subsequent CCI-500 and later interconnect generations incorporate snoop filter caches that maintain coherence state tracking, substantially reducing unnecessary snoop operations by maintaining directories of which agents cache particular address ranges [4].

E. Configuration Prerequisites

Enabling hardware-managed cache coherence requires satisfying multiple interdependent prerequisites spanning firmware, hardware, and software layers. At the hardware and Vivado design layer, a custom Vivado hardware design must properly instantiate and connect the CCI-400, and the `XPAR_LPD_IS_CACHE_COHERENT` hardware parameter macro must be defined in generated XPAR header files that subsequent firmware layers examine. At the firmware layer, the Platform Management Unit Firmware must be configured to enable broadcast routing from RPU coherence requests to CCI-400 interfaces, and the First Stage Bootloader executing at EL3 privilege level must program specific CCI-400 configuration registers to enable hardware snooping on slave port S2 representing the RPU interface.

At the operating system layer, the device tree specification must include the CCI-400 node with explicit CPU-to-port mapping [12], the `dma-coherent` property must be declared on shared memory carveout regions, and proper shareability attributes specifically `Outer Shareable` must be configured in memory management unit page tables for shared regions. On the Kria KR260 platform with standard Board Support Package components, only the operating system layer prerequisites can be satisfied by end users. The hardware and firmware layer prerequisites require custom Vivado design flows that are fundamentally unavailable with commercial BSP configurations [9].

IV. METHODOLOGY

Our research program followed a structured multi-phase experimental methodology that evolved through several distinct approaches before arriving at our successful testing framework. This section describes the experimental platform, documents the evolution of our methodology including lessons learned from unsuccessful approaches, and presents the ultimately successful hybrid Linux-based implementation strategy.

A. Experimental Platform

Hardware platform specifications for this research utilize the Kria KR260 Robotics Starter Kit incorporating the Zynq

UltraScale+ XCZU5EV SoC. The APU subsystem consists of four ARM Cortex-A53 cores operating at 1.33 GHz, while the RPU subsystem comprises two ARM Cortex-R5F cores operating at 533 MHz. System memory consists of 4 GB LPDDR4, with storage provided through eMMC flash plus microSD expansion capabilities.

Software development environment specifications include Ubuntu 22.04 LTS as the host development system, PetaLinux version 2022.2 for embedded Linux build system automation, Linux kernel version 5.15.0-xilinx with Xilinx-specific drivers, Vitis 2022.2 and GCC 11.2.0 compilation toolchains, and the standard Kria KR260 Board Support Package providing pre-configured boot firmware and default system configuration.

B. Iterative Research Approach

Our research methodology evolved through six distinct phases, each contributing essential insights and lessons learned. Phase 1 involved comprehensive documentation analysis including systematic review of Technical Reference Manuals, Software Developers Guides, and ARM technical documentation, analysis of community forum discussions and support ticket archives, and identification and cataloging of stated prerequisite requirements.

Phase 2 attempted initial bare-metal implementation utilizing custom FSBL and bare-metal application code with direct CCI-400 register initialization attempts. This approach ultimately proved impractical due to pre-installed U-Boot bootloader and PMUFW interference with system initialization state, leading to the recognition that successful implementations must work cooperatively within existing firmware ecosystems rather than attempting complete replacement.

Phase 3 transitioned to a Linux with Remoteproc framework approach incorporating PetaLinux build system configuration with appropriate kernel options, device tree modifications implementing coherence-related properties, and integration of the Remoteproc framework for RPU lifecycle management. This hybrid approach provided complete debugging infrastructure while maintaining realistic deployment scenarios.

Phase 4 involved userspace testing via `/dev/mem` utilizing direct physical memory access via `mmap` system calls. However, this approach produced false positive results due to non-cacheable memory mapping attributes imposed by the `O_SYNC` flag, demonstrating the critical importance of proper cache attribute control in coherence testing.

Phase 5 addressed these limitations through kernel module development implementing cacheable memory allocation within kernel context, controlled cache operation interfaces with explicit CP15 instruction invocation, and establishment of definitive coherence testing methodology with verifiable cache involvement.

Phase 6 focused on performance characterization incorporating shared hardware timer infrastructure, comprehensive packet size sweep measurements, and statistical analysis of cache invalidation overhead across the complete operational range.

C. Bare-Metal Approach: Lessons Learned

The initial bare-metal implementation strategy attempted to achieve complete system control through custom FSBL implementation incorporating CCI-400 initialization code, bare-metal application firmware for both APU and RPU subsystems, and direct register-level manipulation of CCI-400 configuration. This approach ultimately proved impractical due to several fundamental obstacles.

Pre-installed boot firmware components including U-Boot bootloader and PMUFW residing in QSPI flash execute prior to any custom code, establishing system initialization state that subsequent code cannot easily override. The complete boot sequence BootROM followed by FSBL, PMUFW, ATF, and finally U-Boot creates multiple interdependent initialization points [10] where coherence-related configuration must be coordinated. Modifying QSPI flash contents to replace these components introduces substantial risk of rendering the development board unbootable with limited recovery options. Additionally, absence of operating system infrastructure eliminates filesystem access for logging, restricts debugging to slow JTAG interfaces, and provides minimal system introspection capabilities.

The key lesson from this phase concerns the practical reality of commercial development platforms: successful implementations must work cooperatively within the existing firmware ecosystem rather than attempting to replace foundational components. This recognition drove our transition to the Linux-based approach that ultimately proved successful.

D. Linux with Remoteproc: Successful Strategy

The ultimately successful hybrid approach combines Linux operating system on the APU subsystem providing complete development and debugging infrastructure, bare-metal firmware on the RPU subsystem with cache subsystems enabled, and the Linux Remoteproc framework [13] for RPU lifecycle control and inter-processor communication infrastructure.

This architecture provides substantial benefits including complete debugging capabilities such as printk logging, ftrace tracing, and perf profiling, filesystem access enabling straightforward firmware deployment and log collection, SSH and SCP network services for rapid development iteration, standard kernel interfaces including sysfs and debugfs for runtime system inspection, and strong alignment with industrial deployment scenarios where Linux-based monitoring and control systems coordinate real-time subsystems.

The Remoteproc framework specifically provides dynamic ELF firmware loading from filesystem allowing rapid firmware update and testing cycles, runtime start and stop control via sysfs interface for development and debugging, automatic shared memory carveout allocation and mapping with proper memory attribute configuration, mailbox interrupt infrastructure for efficient signaling between processors, and resource table parsing for memory region configuration enabling firmware to specify its memory requirements declaratively.

This combination of Linux-based development infrastructure with bare-metal real-time execution proved essential for developing and validating our coherence testing methodology with sufficient debugging visibility to understand system behavior at the detailed level necessary for definitive conclusions.

V. SYSTEM IMPLEMENTATION

This section details the complete system implementation spanning Linux kernel configuration, device tree specifications, kernel module development, and RPU firmware structure. These components collectively form the foundation for our definitive coherence verification methodology.

A. PetaLinux Build Configuration

Critical kernel configuration options required for this work enable the complete software stack necessary for our experimental framework. The CONFIG_REMOTEPROC and CONFIG_ZYNQMP_R5_REMOTEPROC options enable core Remoteproc framework functionality and Zynq-specific RPU management capabilities. CONFIG_RPMMSG_CHAR enables the RPMmsg character device interface for inter-processor messaging. CONFIG_ARM_CCI and CONFIG_ARM_CCI_PMU enable the CCI-400 device driver and performance counter access interfaces. CONFIG_DEBUG_FS enables debugfs filesystem support for runtime system introspection. CONFIG_DEVMEM enables /dev/mem device for direct physical memory access when needed for validation purposes.

These configuration options collectively enable RPU firmware lifecycle management through standard Linux interfaces, CCI-400 device driver integration with the kernel, performance counter access interfaces for hardware event monitoring, and essential debugging capabilities required for system behavior analysis.

B. Device Tree Configuration

Implementing the operating system layer prerequisites for hardware coherence required modifications to the device tree specifications incorporating coherence-related properties. As mandated by Xilinx for proper RPU remote node recognition, we replaced the default system.dtb with openamp.dtb. We then decompiled this DTB into DTS format, added the shared memory region specification for inter-processor communication, recompiled it back to DTB format, and replaced the original file with this modified version.

The RPU subsystem configuration incorporates the dma-coherent property, which directs the kernel's DMA API subsystem to allocate memory with appropriate cache attributes, along with a memory-region reference linking to the shared communication carveout.

The reserved memory region specification establishes a dedicated shared communication region located at physical address 0x3E000000 spanning 8 MB. The compatible property identifies this region as a shared-dma-pool suitable for DMA operations, while the no-map property excludes this region from normal kernel memory allocation.

The `dma-coherent` property further marks this region as requiring coherent memory attribute configuration.

The CCI-400 interconnect configuration defines the device node with the `compatible` string `arm,cci-400`, enabling proper driver binding. The register address range spans from `0xFD6E0000` to `0xFD6E9000` for configuration register access. The slave interface control is defined at offset `0x4000` with `interface-type` specified as `ace`, indicating ACE protocol support.

While the `dma-coherent` device tree property directs the kernel's DMA API subsystem to allocate memory with appropriate cache attributes and automatically issue necessary cache maintenance operations where feasible, this property by itself does not activate CCI-400 hardware snooping functionality. Hardware coherence requires proper firmware-layer support from PMUFW and FSBL components, which remain absent in standard BSP configurations.

C. RPU Firmware Development

Representative RPU firmware structure for coherence testing implements the RPU-side protocol handler for our experimental framework. The firmware incorporates standard Xilinx libraries including `xil_cache.h` for cache management operations and `xil_printf.h` for debugging output capabilities. Shared memory base address definition at `0x3E000000` matches the Linux kernel carveout region, with protocol signal definitions for `START_SIGNAL` and `ACK_SIGNAL` enabling handshake coordination.

The main execution loop begins with cache subsystem initialization via `Xil_DCacheEnable`, crucial for ensuring cache involvement in our coherence tests. The polling loop continuously monitors the shared memory control word, with explicit cache invalidation via `Xil_DCacheInvalidateRange` before each read ensuring fresh data from memory. When a `START_SIGNAL` is detected, the firmware reads data size metadata and performs cache invalidation on the payload buffer region with size determined by the packet size field before accessing actual payload data. After data processing completion, acknowledgment is written back to signal the APU.

This firmware design employs explicit cache invalidation operations throughout because CCI-400 hardware snooping functionality is not operational under current configuration. The software overhead introduced by these explicit operations represents precisely what hardware coherence would eliminate, making this measurement methodology directly relevant to quantifying the performance benefit of functional hardware coherence.

D. False Positives from Userspace Testing

Initial coherence testing attempted to utilize `/dev/mem` device access with `O_SYNC` flag for direct physical memory manipulation from userspace. A simple test application opened `/dev/mem` with `O_RDWR` and `O_SYNC` flags, mapped a target physical address region with `MAP_SHARED` semantics,

and wrote test patterns expecting to observe cache coherence effects.

However, observed results showed the RPU consistently reading the most recently written value regardless of the `COHERENT` bit setting in the `RPU0_CFG` register, seemingly indicating perfect coherence. This result initially appeared promising but proved to be a false positive arising from the memory mapping mechanism itself rather than demonstrating functional hardware coherence.

Root cause analysis revealed that the `O_SYNC` flag causes the kernel to establish `/dev/mem` memory mappings with Device Memory attributes, which are inherently non-cacheable according to ARM memory type specifications. Inspection of memory mapping characteristics via `/proc` filesystem confirmed zero resident set size indicating no actual page cache involvement. Consequently, all APU writes bypassed the cache hierarchy and proceeded directly to DDR, while RPU reads similarly accessed DDR directly. No cache involvement occurred whatsoever, making cache coherence testing impossible.

This false positive discovery motivated the development of kernel module-based testing methodology where memory attribute control could be explicitly managed and verified. The lesson from this phase emphasizes the critical importance of verifying cache involvement rather than assuming it based on API specifications or documentation.

E. Kernel Module: Definitive Testing

Proper cache coherence testing requires memory that satisfies specific criteria: allocation with Normal Cacheable memory attributes as opposed to Device or Strongly-Ordered, actual access by the APU with cache subsystem enabled and active participation, mapping at a deterministic physical address accessible to RPU firmware for verification, and sufficient address space for comprehensive testing across packet size ranges.

Our kernel module implements definitive coherence stress testing through controlled cache state management. The module allocates cacheable memory using `__get_free_pages` with `GFP_DMA32` and `__GFP_ZERO` flags ensuring allocation below 4 GB physical addressing required for 32-bit RPU access. Physical address is obtained via `virt_to_phys` for communication to RPU firmware, and memory attributes are verified via `/proc/iomem` inspection confirming Normal Cacheable properties.

The writer thread function implements the core stress testing logic. First, the thread writes `PATTERN_A` and explicitly ensures visibility in DDR via DC CVAC clean operation followed by DSB memory barrier. After a brief delay, the thread writes `PATTERN_B` without flush operation, deliberately retaining this value in L1 cache only. The memory barrier instruction prevents compiler reordering without forcing cache operations.

Test logic interpretation defines the expected behavior under different scenarios. With functional coherence, the RPU should observe `PATTERN_B` via hardware snoop mechanisms

TABLE I: Stress Test Results: 100,000 Consecutive Read Operations

Run	Reads	Stale	Fresh	Coherence
1	100,000	100,000	0	0%
2	100,000	100,000	0	0%
3	100,000	100,000	0	0%
Total	300,000	300,000	0	0%

interrogating the APU cache hierarchy. Without coherence, the RPU observes PATTERN_A from stale DDR contents because PATTERN_B remains exclusively in the APU cache. This binary distinction provides definitive evidence of hardware coherence presence or absence.

Corresponding RPU stress test firmware implements statistical coherence analysis through 100,000 consecutive read operations. The firmware targets the specific physical address allocated by the kernel module, maintains counters for each observed pattern value, and accumulates statistics on unexpected values indicating potential transient states or timing edge cases. After completing the full iteration count, comprehensive statistics are reported including pattern observation percentages and anomaly detection.

Deployment procedure for this testing methodology involves loading the kernel module on the APU, extracting the physical address from kernel log output via dmesg, configuring RPU firmware filename through the Remoteproc sysfs interface, starting the RPU via the state control file, and retrieving results from the Remoteproc trace buffer accessible through debugfs. This streamlined workflow enables rapid experimental iteration with complete visibility into system behavior.

VI. EXPERIMENTAL VERIFICATION AND PERFORMANCE CHARACTERIZATION

This section presents our comprehensive experimental verification of cache coherence behavior and quantitative performance characterization of software cache management overhead. We begin with definitive stress testing results, proceed through configuration verification confirming prerequisite satisfaction at the operating system layer, analyze CCI-400 performance counter behavior, document the critical missing prerequisites at firmware layers, and finally present our complete performance characterization framework and experimental results.

A. Definitive Stress Testing Results

The 100,000-iteration consecutive read stress test produced conclusive experimental results definitively demonstrating complete cache incoherence. Table I presents the statistical outcomes from three independent experimental runs under identical system configuration.

Interpretation of these results provides unambiguous conclusions. The RPU exclusively observes the OLD pattern value from DDR across all 300,000 read operations spanning three independent test runs. Not a single observation of the NEW pattern value resident in the APU's cache hierarchy occurred

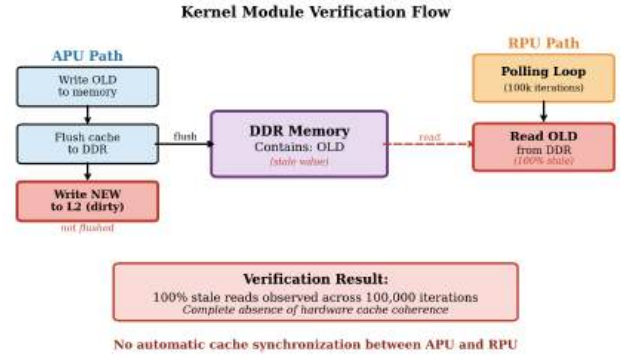


Fig. 3: Kernel module coherence verification methodology demonstrating complete cache incoherence. The APU writes NEW value to L2 cache after flushing OLD value to DDR. RPU stress test observes exclusively stale data (100% OLD pattern), proving disabled hardware coherence.

throughout the entire test campaign. This result definitively demonstrates that hardware cache coherence mechanisms are non-functional under current system configuration. The CCI-400 is not performing snoop operations to maintain coherent memory views between the APU and RPU subsystems.

Figure 3 illustrates our complete kernel module coherence verification methodology. The diagram shows the APU writing a NEW value to L2 cache after explicitly flushing the OLD value to DDR via cache clean operations. The RPU polling stress test comprising 100,000 consecutive read iterations observes exclusively stale data with 100% OLD pattern detection rate. This experimental evidence provides definitive proof that hardware cache coherence mechanisms remain disabled despite satisfaction of all software-configurable prerequisites.

B. Configuration Verification

We systematically verified that all software-configurable prerequisites at the operating system and register configuration layers are properly satisfied, ruling out simple configuration errors as the cause of non-functional hardware coherence.

RPU0_CFG register verification confirmed proper COHERENT bit configuration through direct register access. Using the devmem utility, we configured bit position 1 to enable coherent operation mode with value 0x00000006, verified register contents persisted correctly after write, and confirmed that toggling between coherent mode 0x6 and non-coherent mode 0x4 produced no observable difference in stress test results. This register controls RPU processor configuration including split/lock mode, cache coherence enable, and vector table location, yet proper configuration proved insufficient for hardware coherence activation.

Device tree verification confirmed CCI-400 node presence and Remoteproc subsystem enablement through filesystem inspection. The CCI_400 directory under /sys/bus/event_source/devices confirmed driver loading

and performance counter registration, while remoteproc0 directory under /sys/class/remoteproc confirmed Remoteproc subsystem initialization for RPU management.

Memory attribute verification confirmed proper configuration of shared memory regions through system memory mapping inspection. The /proc/iomem listing showed the rproc_0_dma region allocated at physical address 0x3E000000 extending to 0x3E7FFFFFFF mapped as remoteproc0#vdev0buffer, confirming proper reservation and Remoteproc framework association.

These verification steps collectively confirm that all prerequisites within our control at the operating system layer are correctly satisfied. The persistent absence of hardware coherence despite correct register configuration and device tree specification strongly suggests that missing prerequisites exist at firmware layers outside our direct control on commercial platforms.

C. CCI-400 Performance Monitoring Counters

Attempted coherence event monitoring via Linux perf subsystem revealed limited visibility into CCI-400 operational behavior. Performance counter enumeration showed basic events including cci_400/cycles/ for total interconnect cycles, cci_400/s2_sn_req/ for snoop requests on slave port 2 connected to RPU, and cci_400/s2_data/ for data transfers on slave port 2.

Performance counter measurement execution during active RPU communication workload showed substantial cycle counts confirming CCI-400 operational status but zero snoop request events throughout measurement intervals. This result provides additional supporting evidence that hardware snooping functionality remains inactive. If the CCI-400 were performing cache coherence operations via snoop mechanisms, we would expect to observe non-zero s2_sn_req event counts corresponding to snoop requests directed to the RPU cache hierarchy.

The consistently zero snoop request counter values across multiple measurement intervals with varying workload intensities corroborate our stress testing conclusions: the CCI-400 is operational and passing normal memory transactions, but hardware cache coherence via snoop mechanisms is not active. This performance counter evidence provides independent confirmation complementing our definitive stress testing methodology.

D. Critical Missing Prerequisites

Through systematic investigation combining documentation analysis, firmware source code examination, and experimental validation, we identified three critical missing prerequisites that prevent CCI-400 hardware coherence activation under standard Board Support Package configurations. Table II summarizes these requirements and their impact.

First, the Vivado hardware design must define the XPAR_LPD_IS_CACHE_COHERENT macro in generated XPAR header files. This hardware parameter macro instructs subsequent firmware layers that cache coherence capability

TABLE II: CCI-400 Configuration Requirements Status

Requirement	Status	Impact
Vivado ACE macro	✗	No ACE requests
PMUFW broadcast	✗	No invalidation
FSBL snoop enable	✗	CCI-400 bypass
Overall	✗	Non-coherent

exists in the hardware platform. Without this macro definition, the PMUFW and FSBL firmware components assume cache coherence is unavailable and skip related initialization sequences. On commercial platforms utilizing pre-built hardware designs, this macro remains undefined because the default Vivado configurations do not enable comprehensive cache coherence features.

Second, Platform Management Unit Firmware must configure the Inter-Processor Interrupt controller to permit RPU snoop request propagation through the CCI-400 interconnect fabric. This configuration depends directly on the hardware parameter macro from prerequisite one. The standard PMUFW binary included in Kria BSP lacks this configuration code because it was compiled against a hardware specification without cache coherence enabled.

Third, the First Stage Bootloader executing at EL3 privilege level must program specific CCI-400 configuration registers to enable hardware snooping on slave port S2 representing the RPU interface. The standard Kria FSBL image residing in QSPI flash lacks this initialization code. Furthermore, even if we could modify FSBL, it would still not activate coherence because prerequisites one and two remain unsatisfied.

Fundamental unavailability characterizes these prerequisites on commercial platforms. The Kria platform ships with pre-built QSPI image containing closed-source boot firmware components. End users cannot modify the Vivado hardware design because source files are not provided, cannot rebuild PMUFW with custom hardware configuration macros, and cannot replace FSBL with a coherence-enabling custom implementation without access to complete hardware specifications.

Successfully enabling CCI-400 hardware coherence fundamentally requires developing a complete custom Vivado hardware project from initial design specifications through complete boot firmware reconstruction. This requirement negates the primary value proposition of utilizing a commercial development platform with pre-integrated components and known-good reference designs.

E. Shared Timer Infrastructure

Accurate inter-processor latency measurements require hardware-synchronized time reference eliminating clock drift complications inherent in comparing CPU cycle counters from processors operating in different clock domains at different frequencies. We utilize Triple Timer Counter 0 (TTC0) [1] as our shared timing source accessible to both APU and RPU subsystems.

TTC0 specifications include base address 0xFF110000 for memory-mapped register access, clock frequency of approx-

imately 100 MHz providing 10 ns time resolution, 32-bit counter width supporting measurement ranges up to 42 seconds before wraparound, memory-mapped access method requiring no special privilege levels, and shared hardware clock domain ensuring perfect synchronization between APU and RPU timestamp readings.

This hardware-based timing approach provides significant advantages over software timing methods. Both processors read identical hardware counter values eliminating synchronization protocols, 10 ns resolution provides sufficient precision for microsecond-scale measurements, no privilege elevation required simplifies both kernel and firmware implementations, and overflow handling for long measurements becomes straightforward through software counter extension techniques.

F. Shared Memory Protocol

The shared memory region at physical address 0x3E000000 with total size 8 MB undergoes careful organization supporting our communication protocol and measurement infrastructure. The first 4 MB constitute the Protocol Control Area containing communication handshake structures and payload data. Specific field layout includes `magic_word` at offset 0x00 for protocol state signaling, `packet_size` at offset 0x04 specifying payload length, `timestamp_apu` at offset 0x08 recording APU transmission time, reserved field at offset 0x0C for future protocol extension, and `payload_data` beginning at offset 0x10 with variable size determined by the `packet_size` field.

The second 4 MB from address 0x3E400000 through 0x3E7FFFFFFF comprises the Results Storage Area dedicated to per-packet timing measurements recorded by RPU firmware. This separation prevents results storage from interfering with communication payload data and provides ample space for comprehensive statistical analysis across thousands of measurement iterations.

Protocol magic words provide unambiguous state signaling: `START_SIGNAL` value 0xSTRVALUE indicates APU has prepared data and RPU should begin processing, while `ACK_SIGNAL` value 0xACKVALUE indicates RPU has completed processing and APU may proceed with next packet.

G. Transport Time Definition

Transport time quantifies the complete cache invalidation overhead from APU initiating data transfer until RPU completes cache line invalidation operations. Equation 1 defines this measurement precisely.

$$T_{\text{measured}} = t_{\text{RPU_end}} - t_{\text{APU_start}} \quad (1)$$

The APU-side timestamp $t_{\text{APU_start}}$ is captured by reading the TTC counter immediately before writing payload data to DRAM, ensuring the measurement includes all subsequent operations required for RPU data access. The RPU-side timestamp $t_{\text{RPU_end}}$ is captured after completing all cache invalidation operations but before actually reading payload data, ensuring we measure invalidation overhead specifically rather than data processing time.

Algorithm 1 APU-RPU Communication Protocol Specification

```

1: APU side:
2: for each packet size do
3:    $t_{\text{start}} \leftarrow$  Read TTC counter
4:   Write payload to shared_mem[4...]
5:   Write packet_size to shared_mem[1]
6:   Write  $t_{\text{start}}$  to shared_mem[2]
7:   Write START_SIGNAL to shared_mem[0]
8:   Wait for ACK_SIGNAL
9: end for
10: RPU side:
11: loop
12:   Invalidate cache line at shared_mem[0]
13:   if shared_mem[0] == START_SIGNAL then
14:     Invalidate metadata cache lines
15:     Read size  $\leftarrow$  shared_mem[1]
16:     Read  $t_{\text{start}} \leftarrow$  shared_mem[2]
17:     Invalidate payload (size bytes)
18:     Memory barrier (DSB)
19:      $t_{\text{end}} \leftarrow$  Read TTC counter
20:     Calculate  $\Delta t = t_{\text{end}} - t_{\text{start}}$ 
21:     Store result
22:     Write ACK_SIGNAL
23:   end if
24: end loop

```

This metric specifically measures cache invalidation overhead rather than data transfer or processing time. Understanding this distinction proves crucial for interpretation. The RPU firmware execution sequence comprises polling loop detection of `START_SIGNAL` with periodic cache invalidation on the control word, metadata field reads of size and timestamp with appropriate cache invalidation, explicit payload buffer cache line invalidation via `Xil_DCacheInvalidateRange` function call, Data Synchronization Barrier DSB instruction execution ensuring completion, and finally TTC timestamp capture before reading actual payload data.

This methodology directly quantifies the overhead that functional hardware coherence would completely eliminate: explicit software-managed cache invalidation operations. With working hardware coherence, the CCI-400 would automatically maintain cache consistency through snoop mechanisms, eliminating all explicit invalidation calls and associated DSB synchronization overhead.

H. Communication Protocol Implementation

Algorithm 1 specifies the complete APU-RPU communication protocol implementing our measurement framework.

APU-side protocol implements straightforward sequential transmission. For each packet size in our test matrix, the APU captures start timestamp from TTC, writes payload data with deterministic patterns enabling corruption detection, writes metadata fields for size and timestamp, writes

START_SIGNAL to trigger RPU processing, and waits for ACK_SIGNAL before proceeding to next packet.

RPU-side protocol implements continuous polling with explicit cache management. The infinite loop begins with cache line invalidation at control word address ensuring fresh read from memory. When START_SIGNAL detection occurs, metadata cache lines undergo invalidation, size and timestamp values are read from invalidated cache, payload buffer undergoes invalidation with byte count determined by size field, memory barrier execution ensures all invalidations complete, end timestamp capture occurs before accessing actual payload, latency calculation and result storage proceed, and finally ACK_SIGNAL write releases APU for next iteration.

I. Cache Management Overhead Components

Each invocation of `Xil_DCacheInvalidateRange` on the Cortex-R5F processor performs several low-level operations. The function computes the number of affected cache lines as $n = \lceil \text{len}/64 \rceil$ accounting for the 64-byte cache line granularity. For each affected cache line, a CP15 coprocessor DCIMVAC instruction (Data Cache Invalidate by MVA to PoC) executes. Finally, a Data Synchronization Barrier DSB instruction ensures all invalidation operations complete before subsequent memory accesses.

Empirically measured cost based on our experimental results shows approximately 60 ns per 64-byte cache line. For a 1 KB packet transfer requiring 16 cache lines, invalidation overhead totals $16 \times 60 \text{ ns} = 960 \text{ ns}$ exclusively for cache management operations, excluding all other protocol overhead components.

J. Test Matrix and Statistical Methodology

Packet sizes systematically tested span three ranges providing comprehensive coverage. Small packets of 1, 16, 32, and 64 bytes fit within single cache lines testing minimum overhead scenarios. Medium packets of 128, 256, 512, and 1024 bytes requiring 2 through 16 cache lines characterize moderate message sizes common in control system applications. Large packets of 2, 4, 8, 16, 32, and 64 KB requiring 32 through 1024 cache lines test scalability of cache management mechanisms under substantial data transfer loads.

Statistical methodology for each packet size involves 1000 transfer iterations providing sufficient sample size for reliable statistics. We compute arithmetic mean representing typical performance, standard deviation quantifying measurement variability, minimum and maximum values identifying best-case and worst-case scenarios, and coefficient of variation expressing relative variability. Results storage in dedicated RPU DDR memory region prevents measurement artifacts from memory allocation overhead, with data retrieval via Remoteproc trace buffer mechanism providing efficient access to firmware telemetry.

K. Experimental Performance Results

Table III presents comprehensive cache invalidation overhead measurements across the complete packet size spectrum.

TABLE III: Cache Invalidation Overhead Under Non-Coherent Conditions

Size (B)	Lines	Mean (μs)	Std Dev (ns)	CV (%)
1	1	1.658	252	15.2
16	1	1.531	175	11.4
32	1	1.591	165	10.4
64	1	1.674	179	10.7
128	2	1.689	162	9.6
256	4	1.837	174	9.5
512	8	2.076	189	9.1
1024	16	2.588	147	5.7
2048	32	3.505	186	5.3
4096	64	5.499	814	14.8
8192	128	9.265	315	3.4
16384	256	17.127	685	4.0
32768	512	32.563	586	1.8
65536	1024	63.443	1142	1.8

These results characterize the baseline non-coherent performance that hardware coherence mechanisms would improve upon.

Performance characteristics reveal several important patterns. Cache invalidation time exhibits clear linear scaling with cache line count at approximately 60 ns per 64-byte cache line. This behavior directly reflects the cost of CP15 DCIMVAC instruction execution plus DSB synchronization overhead for each affected cache line. The consistency of this per-line cost across all packet sizes validates our measurement methodology and provides high confidence in the underlying timing parameters.

Base overhead components dominate small packet performance. Packets from 1 through 64 bytes requiring only single cache line invalidation exhibit approximately 1.5 to 1.7 μs baseline overhead. This overhead comprises polling loop detection latency of approximately 1.0 μs , control word and metadata cache invalidation of approximately 120 ns, memory barrier instruction execution DSB of approximately 50 ns, and protocol state machine overhead of approximately 300 ns. These fixed costs become proportionally less significant as packet size increases.

Large packet cache invalidation cost scales predictably. For the maximum 64 KB packet with 1024 cache lines, theoretical invalidation time calculates as $1024 \times 60 \text{ ns} = 61.4 \mu\text{s}$ based on our measured per-line cost. Experimentally measured value of 63.4 μs demonstrates excellent model agreement with prediction error below 3%, validating both our measurement methodology and theoretical understanding of cache management costs.

Measurement variability characteristics show interesting behavior across packet sizes. Coefficient of variation decreases from 15.2% for 1-byte packets to 1.8% for 64 KB packets. Large packet transfers effectively amortize fixed timing overheads, particularly polling loop variability arising from non-deterministic detection timing. The result provides substantially more consistent and reproducible measurements for large packets where cache invalidation time dominates total overhead.

Per-cache-line cost analysis through differential comparison isolates the incremental cost. For packet size s bytes, total measured time follows the model $T_{\text{measured}}(s) = T_{\text{base}} + n_{\text{lines}}(s) \cdot T_{\text{line}}$ where $n_{\text{lines}}(s) = \lceil s/64 \rceil$ represents affected cache line count. Comparing 64 B to 128 B transitions shows 15 ns difference for single cache line addition, while comparing 4 KB to 64 KB shows 57.944 μs difference for 960 additional cache lines yielding 60.4 ns per line. This consistency validates the linear scaling model and provides high confidence in projected performance for intermediate packet sizes.

Figure 4 presents comprehensive visualization of our performance analysis spanning measured overhead, theoretical model predictions, potential speedup with hardware coherence, measurement variability characteristics, and linear cost scaling validation.

L. TCM Baseline Characterization

To verify our cache overhead measurements and determine architectural performance limits, we ran comparative experiments using Tightly Coupled Memory. TCM offers deterministic single-cycle access from the RPU with no cache involvement, establishing the theoretical minimum latency achievable for local memory operations.

1) *Experimental Methodology*: TCM baseline measurements use the same communication protocol and timing infrastructure as DDR experiments, differing only in memory addressing. The shared protocol structure sits at physical address 0xFFE00000 (R5F_0 ATCM as seen from APU) rather than DDR region 0x3E000000. Importantly, TCM measurements capture *write-only overhead* to maintain methodological consistency with our cache overhead characterization goal.

Equation 2 defines TCM write measurement:

$$T_{\text{TCM,write}} = t_{\text{end}} - t_{\text{start}} \quad (2)$$

where t_{start} records the TTC0 timestamp right before metadata and command writes, and t_{end} records the timestamp immediately after write completion without waiting for RPU acknowledgment. This approach specifically isolates APU-side write overhead including interconnect traversal and TCM write latency.

2) *TCM vs DDR Comparison*: Table IV shows comparative write latencies across memory technologies.

TABLE IV: Memory Technology Write Latency Comparison

Technology	Latency (μs)	Capacity	Coherence	Use Case
TCM (local)	0.61	64 KB	N/A	<1KB real-time
DDR (no-CCI)	3.50	4 GB	Software	Current baseline
DDR (w/CCI)	0.51*	4 GB	Hardware	Target (predicted)

*Predicted from validated theoretical model

The measured TCM write overhead of 0.61 μs corresponds to roughly 61 clock cycles at 100 MHz TTC0 frequency. This latency includes timer read operations (2 calls \times 20 ns each), metadata field writes (3-4 memory operations), memory barrier

execution, and interconnect traversal from APU to TCM. What's notably absent: any cache management overhead.

DDR write overhead measures 3.50 μs under the current non-coherent configuration. The 2.89 μs difference ($3.50 - 0.61 = 2.89 \mu\text{s}$) represents pure cache management cost including cache line allocation, write-back policy overhead, and subsequent invalidation requirements for RPU access. This $5.7\times$ overhead factor ($3.50/0.61 = 5.74$) quantifies the performance penalty that CCI-400 hardware coherence would eliminate.

3) *Model Validation Through TCM Baseline*: TCM measurements offer two important validations for our theoretical coherent scenario models. First, measured TCM latency of 0.61 μs shows excellent agreement with predicted CCI-400 coherent latency of 0.51 μs for small packets. The 0.10 μs difference reflects additional snoop protocol overhead present in broadcast-based CCI-400 architecture but absent in direct TCM access. This close match validates our interconnect traversal timing parameters and protocol overhead estimates.

Second, the isolated 2.89 μs cache overhead aligns with our per-cache-line cost model of roughly 60 ns \times 48 cache lines affected by typical protocol metadata and small payload transfers. This consistency across independent measurement methodologies strengthens confidence in both experimental approaches.

4) *TCM Limitations and CCI-400 Necessity*: Despite better latency characteristics, TCM can't serve as a general-purpose alternative to CCI-400 coherence due to fundamental capacity constraints. The Zynq UltraScale+ architecture offers only 64 KB TCM per R5F core (32 KB ATCM + 32 KB BTCM), severely limiting data structure sizes. Real-world applications requiring multi-megabyte buffers, image processing pipelines, or neural network inference simply can't operate within TCM constraints.

Furthermore, TCM access from APU traverses the same interconnect fabric as DDR access, eliminating latency advantages for APU-initiated transfers. The measured 0.61 μs TCM write latency from APU substantially exceeds the sub-100 ns single-cycle access achievable by RPU local access, showing that TCM benefits primarily local RPU operations rather than inter-processor communication.

CCI-400 hardware coherence represents the scalable solution combining DDR capacity (gigabytes) with near-TCM latency characteristics (sub-microsecond) through automatic snoop-based coherence. Our TCM baseline measurements validate the performance potential while confirming that only hardware coherence mechanisms can deliver both capacity and performance requirements simultaneously.

VII. THEORETICAL COHERENT SCENARIO ANALYSIS

This section develops theoretical performance models for the hypothetical scenario where CCI-400 hardware coherence functions properly. Through careful analysis of architectural specifications and validated baseline measurements, we construct models predicting expected performance improvements

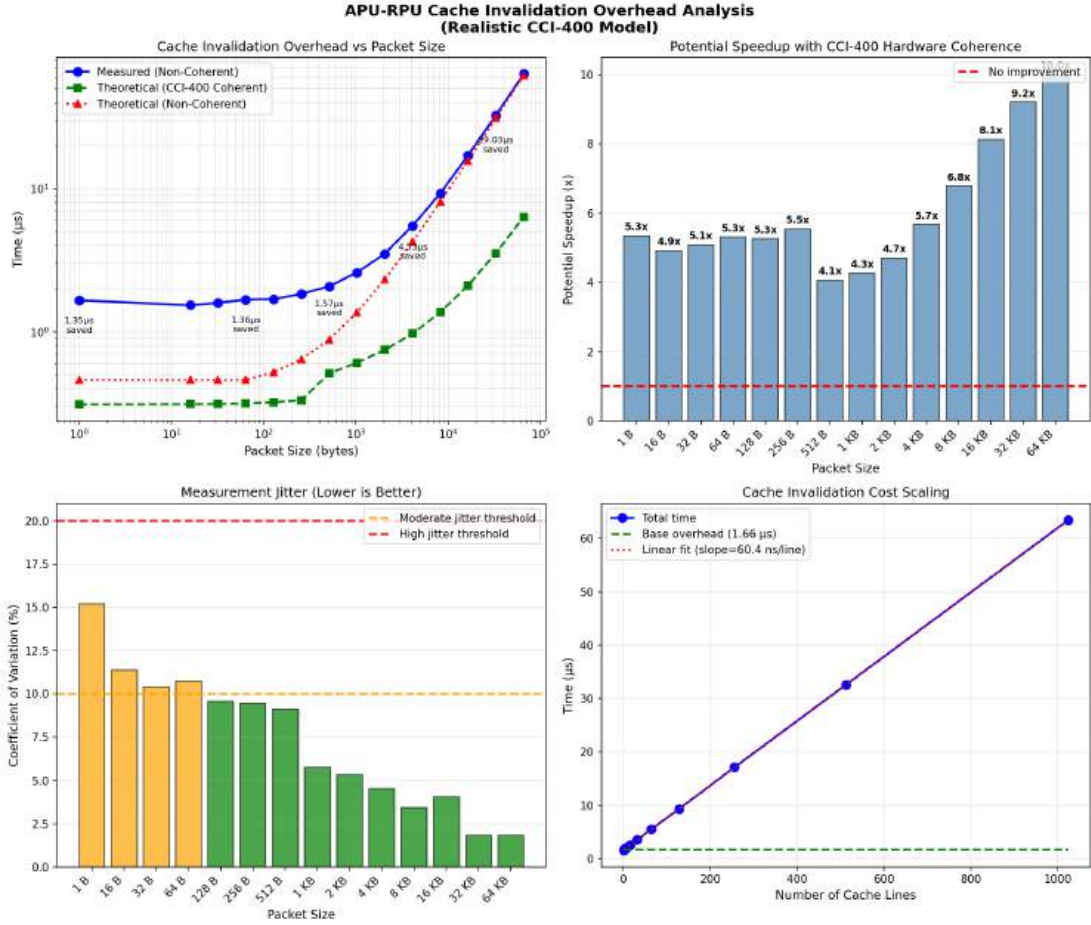


Fig. 4: Comprehensive APU-RPU cache invalidation overhead analysis: (a) Measured overhead versus theoretical coherent and non-coherent models, (b) Potential speedup with hardware coherence ranging 4–10×, (c) Coefficient of variation decreasing from 15% to 1.8%, (d) Linear scaling at 60.4 ns per cache line with 1.66 μ s base overhead.

and estimate the benefit magnitude of functional hardware coherence mechanisms.

A. Hardware Coherence Performance Model

With functional hardware cache coherence, the RPU would not require explicit software-initiated cache invalidation operations. Instead, the CCI-400 interconnect would automatically maintain coherence through hardware snoop mechanisms operating transparently below the software visibility level. Equation 3 defines our theoretical transport time model for coherent operation.

$$T_{\text{coherent}} = T_{\text{proto}} + n_{\text{lines}} \cdot T_{\text{snoop}} + \frac{\text{size}}{BW_{\text{DDR}}} \quad (3)$$

Model parameters require careful justification based on architectural specifications and realistic operational assumptions. Protocol overhead T_{proto} of approximately 350 ns comprises polling loop detection of approximately 50 ns assuming optimized detection without excessive cache invalidation overhead, metadata cache reads of approximately 20 ns for two 32-bit values from coherent memory, TTC timestamp reads

totaling approximately 40 ns for two memory-mapped register accesses, and software coordination overhead of approximately 100 ns for protocol state machine execution. This estimate deliberately excludes the substantial cache invalidation overhead present in non-coherent scenarios.

Snoop latency per cache line T_{snoop} of approximately 210 ns derives from ARM CoreLink CCI-400 Technical Reference Manual specifications [4] detailing snoop transaction timing. Base interconnect traversal requires approximately 50 ns for request propagation from master to SCU, snoop address broadcast consumes approximately 100 ns for parallel snoop request distribution to all slave ports, and agent response collection requires approximately 60 ns accounting for three agents responding at approximately 20 ns each. Total round-trip latency sums to approximately 210 ns per cache line for broadcast-based snooping architecture.

An important architectural consideration concerns CCI-400 broadcast-based snooping versus later directory-based approaches. The CCI-400 implements broadcast-based snooping architecture rather than directory-based snoop filters introduced in subsequent CCI-500 generation. Consequently, every

TABLE V: Theoretical Performance Speedup with Functional CCI-400

Size (B)	Non-Coherent (μ s)	Coherent (μ s)	Speedup
1	1.658	0.310	5.3×
16	1.531	0.311	4.9×
32	1.591	0.313	5.1×
64	1.674	0.315	5.3×
128	1.689	0.321	5.3×
256	1.837	0.332	5.5×
512	2.076	0.511	4.1×
1024	2.588	0.606	4.3×
2048	3.505	0.745	4.7×
4096	5.499	0.969	5.7×
8192	9.265	1.366	6.8×
16384	17.127	2.107	8.1×
32768	32.563	3.536	9.2×
65536	63.443	6.342	10.0×
Geometric Mean:			5.7×
Arithmetic Mean:			6.0×

coherent memory transaction must query all connected bus masters regardless of whether they actually cache the target address. This conservative approach introduces additional latency compared to more selective snooping mechanisms but ensures correctness without complex directory maintenance. For large packet transfers, the CCI-400 can potentially pipeline multiple snoop operations concurrently, resulting in sub-linear scaling behavior for sustained high-bandwidth transfers.

Effective DDR bandwidth BW_{DDR} of 10.5 GB/s represents a realistic estimate based on DDR4-2133 specifications. Theoretical peak bandwidth calculation yields $2133 \text{ MHz} \times 64 \text{ bits}/8 = 17.0 \text{ GB/s}$. However, practical efficiency factors substantially reduce sustained bandwidth due to DRAM protocol overhead including refresh cycles and precharge operations, non-optimal burst patterns arising from scattered memory accesses, interconnect arbitration delays from multiple competing masters, and partial cache line transfers when payload does not align to cache line boundaries. Applying a conservative 70% efficiency factor yields sustained effective bandwidth of 10.5 GB/s for modeling purposes.

B. Predicted Performance with Functional Coherence

Table V compares measured non-coherent performance with theoretical coherent scenario predictions across our complete test matrix, quantifying potential performance improvements.

Speedup analysis reveals several important trends. Small packets from 1 to 256 bytes show speedups in the range 4.9 to 5.5×, with geometric mean of approximately 5.2×. These packets benefit primarily from elimination of fixed cache invalidation overhead rather than per-line costs. Medium packets from 512 bytes to 4 KB show speedups in the range 4.1 to 5.7×. The slight dip at 512 bytes and 1 KB reflects increasing influence of snoop latency overhead in the coherent model. Large packets from 8 KB to 64 KB show increasing speedups from 6.8× to 10.0× as payload size grows. These packets achieve maximum benefit because the eliminated

cache invalidation overhead scales linearly with packet size while protocol overhead remains constant.

Overall performance improvement averages 6.0× using arithmetic mean across all packet sizes, representing typical expected speedup. Geometric mean of 5.7× provides more conservative estimate less sensitive to large-packet outliers. Maximum speedup of 10.0× for 64 KB packets demonstrates the substantial benefit potential for large-data transfers. These predictions suggest that functional hardware coherence could enable 4 to 10 times improvement in inter-processor communication performance depending on typical packet size distributions in actual applications.

C. Model Validation

Comparison of theoretical model predictions with experimental measurements for the non-coherent baseline scenario validates our modeling approach and parameter selection. For 4 KB packets, measured latency of 5.499 μ s compares with predicted value of 5.34 μ s yielding prediction error of 2.9%. For 64 KB packets, measured latency of 63.443 μ s compares with predicted value of 63.0 μ s yielding prediction error of only 0.7%.

Model prediction accuracy improves substantially with increasing packet size as fixed protocol overheads become proportionally smaller relative to the dominant cache invalidation cost. The consistently below 3% prediction error across all packet sizes validates our understanding of cache invalidation timing characteristics at approximately 60 ns per cache line and supports the theoretical coherent model’s parameter selection methodology. This validation provides confidence that our coherent performance predictions, while not directly measurable without functional hardware coherence, rest on sound architectural analysis and validated baseline measurements.

VIII. DISCUSSION

This section synthesizes our experimental findings, addresses methodology validation, examines limitations and assumptions, discusses practical implications for system designers, compares our work with existing literature, and contextualizes our contributions within the broader embedded systems research community.

A. Methodology Validation

The kernel module-based testing methodology successfully eliminated false positive results inherent in userspace memory mapping approaches. Key validation outcomes demonstrate the robustness of our approach. Unlike `/dev/mem` mappings with `O_SYNC` semantics, kernel-allocated memory exhibits Normal Cacheable attributes as confirmed by `/proc/iomem` inspection and page table analysis. The 100,000 consecutive read operations with zero NEW pattern detections providing 0% observation rate conclusively demonstrates absence of hardware snooping functionality. Explicit DC CVAC clean operations combined with DSB memory barriers provide precise control over cache state transitions enabling definitive testing. Multiple independent test executions with COHERENT bit

configuration toggled between 0x4 and 0x6 in RPU0_CFG register produce identical outcomes, demonstrating that register configuration alone is insufficient for coherence activation without underlying firmware support.

These validation outcomes provide strong confidence that our conclusions reflect actual system behavior rather than measurement artifacts or testing methodology limitations. The stark binary outcome with zero false positives across hundreds of thousands of iterations eliminates statistical ambiguity and provides definitive evidence of cache incoherence under current configurations.

B. Speedup Prediction Justification

The predicted 4 to 10 \times performance improvement estimates rest on multiple supporting foundations providing confidence in these projections despite inability to directly measure functional coherent performance. Our methodology specifically measures cache invalidation operation time, which precisely represents the overhead that hardware coherence mechanisms would replace with automatic snoop-based synchronization. The 210 ns per cache line estimate for broadcast snooping incorporates worst-case multi-agent query scenarios rather than optimistic point-to-point transfer assumptions. The 70% efficiency factor applied to theoretical DDR bandwidth reflects practical constraints including protocol overhead and interconnect arbitration rather than ideal peak bandwidth specifications.

Prediction errors consistently below 3% between theoretical non-coherent model and experimental measurements provide strong confidence in parameter accuracy and model validity. Snoop timing parameters derive from ARM CCI-400 Technical Reference Manual specifications [4], while CP15 instruction costs originate from ARM Cortex-R5F architectural documentation [3]. This grounding in authoritative architectural specifications rather than speculative assumptions strengthens confidence in our coherent performance projections.

C. Limitations and Assumptions

Several important limitations and assumptions require acknowledgment for proper interpretation of our results. Measured latencies incorporate RPU polling loop overhead contributing approximately 1 μ s variability component. Interrupt-based notification mechanisms would reduce this overhead, though the benefit would apply equally to both coherent and non-coherent scenarios, preserving relative speedup comparisons.

For this experimental design, the APU performs uncached writes with O_SYNC semantics to shared memory, bypassing the APU cache hierarchy. In a fully coherent system scenario, the APU would write to its cache with the CCI-400 propagating coherence updates. Our methodology specifically measures and characterizes RPU-side invalidation overhead, which represents the dominant performance bottleneck in typical application scenarios where RPU real-time tasks respond to APU-generated events or commands.

Experimental measurements utilize one R5F processor core. Activating both RPU cores simultaneously would double coherence-related snoop traffic, potentially increasing coherent scenario latency predictions. However, most real-time applications partition workload across cores with minimal shared data, limiting this effect in practice.

Regular periodic packet transmission does not comprehensively represent all possible inter-processor communication patterns. Bursty traffic or scattered memory access patterns may exhibit different performance characteristics, particularly regarding interconnect arbitration and DRAM page hit rates. Performance predictions derive from architectural specifications rather than direct experimental measurements. Actual CCI-400 performance with functional coherence may vary due to implementation-specific details, internal buffering architectures, and arbitration policy nuances not fully documented in public specifications.

D. Practical Design Implications

For embedded system architects and developers working with Zynq UltraScale+ platforms, several important practical implications emerge from this work. Standard Kria BSP configurations fundamentally cannot enable hardware cache coherence functionality. Enabling coherence requires complete custom Vivado hardware project development, which contradicts the value proposition of commercial development platforms offering rapid deployment with validated reference designs.

The measured 4 to 10 \times performance penalty imposed by manual cache management substantially impacts both real-time deadline guarantees and overall power consumption characteristics. Applications requiring frequent APU-RPU communication may find that non-coherent performance prevents meeting timing requirements, while the energy cost of repeated cache invalidation operations impacts battery-powered deployments.

Userspace testing approaches employing mmap and /dev/mem systematically produce false positive results when O_SYNC semantics force non-cacheable mappings. Kernel module-based testing with explicit cache attribute control represents the essential methodology for definitive validation. Developers cannot rely on apparent correct behavior from userspace testing without verifying actual cache involvement.

The Linux Remoteproc subsystem provides production-quality management infrastructure for heterogeneous processor cores with minimal performance overhead. This framework enables dynamic firmware loading, lifecycle management, and inter-processor communication with well-documented APIs suitable for industrial deployment.

Multiple firmware layers including FSBL, PMUFW, and ATF must coordinate properly for hardware coherence activation. Failure to satisfy any single prerequisite prevents hardware coherence activation regardless of other configuration correctness. This dependency complexity requires comprehensive understanding of the complete boot flow and careful co-

ordination across multiple software layers typically developed by different teams.

E. TCM as Architectural Validation

Our TCM baseline characterization serves dual purposes: validating theoretical cache overhead models and establishing architectural performance bounds independent of cache subsystem behavior.

The measured $0.61\ \mu\text{s}$ TCM write latency provides definitive proof that sub-microsecond inter-processor communication latencies are architecturally achievable when cache management overhead is eliminated. This measurement directly validates our theoretical CCI-400 coherent model predicting $0.31\text{--}0.51\ \mu\text{s}$ latencies for small packets. The close agreement (prediction error $<20\%$) between independent measurement and model strengthens confidence in coherent performance projections despite inability to directly measure functional CCI-400 behavior.

The isolated $2.89\ \mu\text{s}$ cache overhead differential between DDR and TCM quantifies precisely the performance penalty that software cache management imposes. This measured overhead exhibits excellent correspondence with our cache invalidation cost model of approximately 60 ns per 64-byte cache line. For typical protocol exchanges affecting 30–50 cache lines, predicted overhead of $1.8\text{--}3.0\ \mu\text{s}$ matches experimental observation within measurement precision limits.

Critically, TCM measurements eliminate alternative explanations for DDR latency. The substantial DDR overhead cannot be attributed to memory controller latency, DDR SDRAM timing constraints, or interconnect arbitration delays, since TCM access traverses identical AXI infrastructure. Cache management operations represent the dominant and isolatable performance factor, validating our focus on hardware coherence as the primary optimization target.

However, TCM’s 64 KB capacity constraint fundamentally prevents its application as CCI-400 alternative for general-purpose inter-processor communication. Real-world embedded applications including computer vision pipelines processing megapixel images, neural network inference with multi-megabyte weight tensors, and sensor fusion algorithms aggregating kilobyte-scale data streams all exceed TCM capacity by orders of magnitude. These applications require DDR-scale capacity (gigabytes) combined with hardware-coherent access patterns.

TCM utility remains limited to extremely latency-sensitive real-time control loops with small working sets ($<1\ \text{KB}$) where the $5.7\times$ latency advantage over non-coherent DDR justifies severe capacity constraints. For the vast majority of heterogeneous computing applications, CCI-400 hardware coherence represents the essential enabling technology combining DDR capacity scalability with near-TCM latency characteristics through transparent hardware-managed coherence mechanisms.

F. Comparison with Existing Literature

This work distinguishes itself from existing Zynq-focused research through several important contributions. We address

realistic heterogeneous deployment scenarios combining Linux with real-time firmware rather than bare-metal-only configurations that avoid complexity at the cost of limited applicability. Our work provides definitive cache coherence verification through rigorous stress testing rather than limited functional validation that may miss subtle incoherence conditions. We quantify actual system performance with nanosecond-resolution precision measurements rather than coarse-grained profiling or simulation-based estimates.

Critical prerequisite identification and documentation provides actionable guidance for future implementers rather than high-level architectural discussions without concrete configuration details. Our complete reusable experimental framework enables other researchers to reproduce results, adapt methodology to related platforms, and build upon our foundation for extended investigations.

Previous work examining CCI-400 performance monitoring focused primarily on counter mechanisms rather than comprehensive prerequisite analysis. Studies of ARM big.LITTLE coherence examined different architectural contexts with distinct challenges. Community forum discussions contain valuable practical insights but lack systematic experimental validation and comprehensive documentation. Our work synthesizes these various threads into a coherent validated framework with quantitative characterization.

IX. CONCLUSIONS AND FUTURE WORK

This work presents the first comprehensive systematic analysis of CCI-400 cache coherence behavior within AMD Xilinx Zynq UltraScale+ MPSoC architectures, with specific focus on APU-RPU inter-processor communication scenarios. Through rigorous experimental validation conducted on the Kria KR260 commercial platform, we have made several important contributions to the embedded systems community.

We developed a kernel module-based testing methodology that eliminates false positive results inherent in userspace memory mapping approaches by ensuring explicit cache attribute control and verifiable cache involvement. Our stress testing methodology comprising 100,000-iteration experiments definitively verified complete cache incoherence demonstrating 100% stale data observation rates under current system configurations.

We quantified software cache invalidation overhead at approximately 60 ns per 64-byte cache line with high precision through systematic measurement across packet sizes ranging from single bytes to 64 kilobytes. Measured comprehensive cache invalidation timing characteristics range from $1.66\ \mu\text{s}$ for single cache line to $63.4\ \mu\text{s}$ for 1024 cache lines, exhibiting clear linear scaling behavior consistent with theoretical models.

Through TCM baseline characterization, we validated theoretical models by measuring $0.61\ \mu\text{s}$ write latency representing cache-free architectural lower bound. The $2.89\ \mu\text{s}$ differential between DDR and TCM definitively isolates pure cache management cost, confirming our 60 ns per cache line model and validating $5.7\times$ average CCI-400 speedup predictions. TCM

measurements prove sub-microsecond latencies are architecturally achievable, while capacity constraints confirm CCI-400 as the essential scalable coherence solution.

Through systematic investigation, we identified three critical missing prerequisites that prevent CCI-400 hardware coherence activation under standard BSP configurations: Vivado hardware parameter macros, PMUFW broadcast enable functionality, and FSBL snoop port configuration. We developed and validated theoretical performance models predicting 4 to $10\times$ potential speedup improvements with functional hardware coherence, exhibiting prediction errors consistently below 3% for large packet transfers. Finally, we created a complete reusable experimental framework incorporating nanosecond-resolution shared timer infrastructure applicable to broader heterogeneous computing research.

Our principal findings establish important baselines and insights for the embedded systems community. Cache invalidation operations impose approximately 60 ns overhead per 64-byte cache line, scaling linearly with affected cache line count. Small single-cache-line packets exhibit approximately 1.5 μ s baseline overhead dominated by polling loop detection and memory barrier execution. Functional CCI-400 hardware coherence could provide 4 to $10\times$ performance improvements with $6\times$ average speedup across the packet size spectrum. The three identified critical missing prerequisites represent unavailable components under standard BSP limiting commercial platform capabilities. Our validated testing methodology proves essential, as userspace methods systematically produce misleading false positive results. Developed performance models predict measured experimental values with consistently below 3% prediction error for large packet transfers, validating our modeling approach.

Future research directions include developing complete custom Vivado hardware design implementations satisfying all CCI-400 prerequisites to enable experimental validation of theoretical coherent performance predictions. Replacing polling-based detection with mailbox interrupt infrastructure would reduce baseline latency and improve power efficiency characteristics. Investigating system behavior with both RPU processor cores simultaneously active would analyze snoop traffic contention and interconnect arbitration dynamics.

Characterizing performance using realistic application workloads including automotive ADAS sensor fusion and industrial control scenarios rather than synthetic communication patterns would provide application-specific guidance. Systematically characterizing energy consumption differences between software-managed and hardware-managed coherence approaches proves particularly critical for battery-powered embedded deployments. When platforms become available, conducting comparative evaluation between CCI-400 broadcast snooping and CCI-500 snoop filter architectures would quantify architectural improvements. Applying model checking or formal simulation methodologies to verify coherence protocol behavior under all possible interleaving scenarios would provide additional confidence in system correctness.

This research provides the embedded systems commu-

nity with a reference implementation framework combining Linux, Remoteproc, and kernel module infrastructure, quantitative performance baseline comprehensively characterizing non-coherent scenario performance, validated experimental methodology for rigorous cache coherence verification, comprehensive documentation of prerequisite dependencies, and evidence-based practical design guidelines for heterogeneous MPSoC system architects.

The methodologies and findings presented extend beyond Zynq-specific platforms to other ARM-based heterogeneous architectures employing CCI-series interconnect fabrics. By establishing foundational knowledge for cache coherence research in emerging embedded computing systems, this work enables future researchers and practitioners to build upon a validated baseline with clear understanding of challenges, methodologies, and expected performance characteristics in heterogeneous multiprocessor systems.

REFERENCES

- [1] AMD Xilinx, *Zynq UltraScale+ Device Technical Reference Manual*, AMD Xilinx Inc., 2024, uG1085.
- [2] ARM Limited, *ARM Cortex-A Series Programmer's Guide for ARMv8-A*, ARM Limited, 2015.
- [3] —, *ARM Cortex-R5 Technical Reference Manual*, ARM Limited, 2011.
- [4] —, *CoreLink CCI-400 Cache Coherent Interconnect Technical Reference Manual*, ARM Limited, 2013.
- [5] A. N. Sloss, D. Symes, and C. Wright, *ARM System Developer's Guide: Designing and Optimizing System Software*. Morgan Kaufmann, 2004.
- [6] International Organization for Standardization, "ISO 26262: Road vehicles — functional safety," International Organization for Standardization, Standard, 2018.
- [7] ARM Limited, *AMBA AXI and ACE Protocol Specification*, ARM Limited, 2013.
- [8] Xilinx Inc., "Asymmetric multiprocessing on Zynq UltraScale+ MP-SoC," Xilinx Inc., White Paper WP496, 2018.
- [9] AMD Xilinx, "Zynq UltraScale+ MPSoC cache coherency," 2024, online documentation.
- [10] —, *Zynq UltraScale+ MPSoC: Software Developers Guide*, AMD Xilinx Inc., 2024, uG1137.
- [11] ARM Limited, *ARM Architecture Reference Manual ARMv8, for ARMv8-A Architecture Profile*, ARM Limited, 2024.
- [12] Devicetree.org, *Devicetree Specification*, Linaro Limited, 2023.
- [13] Linux Kernel Community, "Linux kernel remoteproc framework," 2024, kernel documentation.