



Energy-Efficient Microarchitecture Design for Real-Time Embedded Systems in Modern SoCs

Jayanth

Staff Engineer at Qualcomm Technologies Inc, USA.

Abstract: The rapid expansion of real-time embedded systems in automotive electronics, Internet of Things (IoT) devices, robotics, aerospace platforms, and edge AI applications has led to a heated demand for System-on-Chip (SoC) designs that not only provide timing guarantees but also operate with minimal energy consumption. Contemporary SoCs are progressively layering heterogeneous processing elements, complex memory hierarchies, and sophisticated power-management features; however, the problem of how to make these SoCs energy-efficient without breaking the real-time determinism remains unsolved. Our research is about the conflicting problem of energy optimization versus the need for execution predictability in real-time embedded scenarios and we resolve this problem by choosing microarchitectural solutions that can achieve the balance of both. We articulate a comprehensive microarchitecture-level plan that integrates timing-aware pipeline design, cache and memory access optimization, fine-grained power gating, dynamic voltage and frequency scaling under latency constraint, and workload-aware scheduling mechanisms. The suggested strategy results in the establishment of a design framework that not only facilitates energy optimization but also, and more importantly, performs worst-case execution time analysis, thus providing an extra layer of security that power-saving methods do not interfere with the execution of deadlines. Optimization techniques are firstly modeled analytically and through cycle-accurate simulations and then verified on a state-of-the-art SoC platform with an application-specific real-time workload indicative of safety- and latency-critical scenarios. The work done in the case study leads to the energy consumption reduction, which can be quantitatively measured while real-time constraints are continuously being fulfilled, thus bringing to the fore the compromises between the performance headroom, predictability, and power savings. The findings reveal that if done properly, the coordination of microarchitectural techniques can bring about sizable energy efficiency improvements at the cost of only slight reductions in latency and throughput. Essentially, this study has accomplished the difficult task of uniting energy-aware microarchitecture design with real-time system requirements, which in turn opens up new opportunities for SoC architects and embedded system designers.

Keywords: Energy-Efficient Microarchitecture, Real-Time Embedded Systems, System-On-Chip (SoC), Low-Power Design, Deterministic Execution, Dynamic Power Management, Embedded Processors.

1. Introduction

The fast evolution of real-time embedded systems has essentially changed the design priorities of modern System-on-Chip (SoC) platforms. Real-time computation is becoming a requirement in the mentioned domains to ensure safety, reliability, and responsiveness. In these systems, correctness is not only defined by logical accuracy but also by the capability of producing results within strict timing deadlines. Meanwhile, concerns about energy consumption, thermal limits, and sustainability have led to energy efficiency being at the core of SoC design. This simultaneous requirement of deterministic timing and low power consumption constitutes a complex issue for the design of SoCs, which become more heterogeneous and feature-rich as semiconductor technologies continue to scale.

Conventional processor and SoC architectures were mostly optimized for average-case performance and throughput with a focus on high utilization and power-saving mechanisms. However, these methods are not always efficient for real-time workloads where the worst-case execution scenario and predictability are of utmost importance. The energy efficiency-real-time determinism interaction becomes more complex as embedded systems implement advanced features such as multi-core architectures, specialized accelerators, deep memory hierarchies, and dynamic power management. Taking on this challenge involves having a microarchitecture-aware view that not only reconciles energy optimization strategies with timing guarantees but also does so from the very first stages of design.

1.1. Challenges in Energy-Efficient Real-Time SoC Design

One of the features that distinguish real-time systems from others is the existence of very strict and non-negotiable timing deadlines. Tasks should be accomplished within the time frames that are previously set, and a failure to meet the deadline may result in a total failure of the system, occurrence of safety risks or lowering the quality of service level. Unlike best-effort systems where latency variations can occur from time to time without any problem, real-time platforms require predictable execution behavior. The main disadvantage of such a demand is that it directly limits the number of energy-saving methods that may dynamically adapt or use speculative behavior because they may cause execution time to vary.

The interrelation among energy efficiency, performance, and predictability is a complex one. Most of the time, energy-saving measures imply, at least, some increase in execution latency or timing jitter - for example, in the case of voltage and

frequency reduction, component idle gating, or resource dynamic reconfiguration. On the contrary, the strategy of planning for performance at worst-case scenario to ensure time constraints can, by raising the system's capacity, cause inefficiency and loss of energy. Solving these problems of competing goals is especially hard for the systems operating non-stop under the changing workloads and environmental conditions.

Moreover, deeply submicron semiconductor technologies complicate even more the struggle to keep the balance. The shrinking of the semiconductor process helps the chip to become smaller, faster, and more energy-efficient but at the same time, the new chip introduces several new challenges such as increased leakage power, variability in the manufacturing process, and higher temperature sensitivity. The current leaking can be the major contributor to the power consumed when the system is in an idle or low-utilization state; thus, the conventional expectations of energy-saving measures lose their ground. The variability of the process introduces the uncertainty of timing behavior, thus making the worst-case analysis more conservative and less energy-efficient.

On the other hand, thermal constraints are the major ones, especially in compact devices powered by batteries and systems that rely on passive cooling, which can be wearable devices, autonomous sensors, and embedded controllers in a harsh environment, for example. Too much heat can lead to less reliability, shorter life of components, and the giving-in of the real-time constraints by the alarming mechanisms of overheating that are forced to be used. The control of thermal behavior without any performance losses has to be achieved with the proper interaction of microarchitecture design, power management, and workload scheduling.

Lastly, contemporary SoCs are progressively implementing heterogeneous cores and specialized accelerators not only to augment performance but also to be more efficient. The state of heterogeneity brings about the advantages that the system can have but at the same time, it becomes more complex. The different processing elements characterize different power-performance as well as timing behaviors, thus the global optimization gets more difficult to solve. How to manage the work distribution among CPUs, GPUs, DSPs, and AI accelerators in an energy-saving, yet, predictable manner is still a research question in real-time SoC design.

1.2. Problem Statement

Traditional low-power design techniques that are ideal for general computing may not be as good when they are directly applied to real-time embedded workloads. In fact, these aggressive dynamic voltage and frequency scaling, speculative execution, out-of-order processing, and deep caching hierarchies mainly work towards efficiency of the average-case performance. In a real-time scenario, these mechanisms may cause non-deterministic latency thus complicating the worst-case execution time as well as weakening the timing guarantees. Consequently, the designers are forced to use these power-saving features at a minimal level in order to ensure predictability. However, this results in a waste of power since the energy usage is not efficient.

The determinism of power management mechanisms is a major factor that limits current solutions. Indeed, the delay from switching between power states, voltage levels, or clock domains is very variable and cannot be easily bounded analytically. Furthermore, the same components such as caches, interconnects, and memory controllers that can be present in multi-core and heterogeneous environments may also be the reason of contention and timing interference between tasks resulting in bigger issues for worst-case scenarios which are already hard to be studied while energy optimization is pursued. In such a situation, it is still a huge challenge to think about the worst-case while trying to save energy aggressively.

Besides that, the majority of general-purpose microarchitectures fail to meet the requirements of embedded real-time systems adequately. The features that are designed to maximize instruction-level parallelism or throughput may increase complexity and unpredictability without giving proportional benefits to real-time workloads. The result of this discrepancy is that energy consumption and timing behavior are inefficient.

Hence the existence of a demand for energy optimization that is aware of microarchitecture and can explicitly consider real-time constraints. Such strategies would not treat energy efficiency and determinism as two opposing goals but instead they would incorporate timing analysis into the design and control of microarchitectural features facilitating predictable and energy-efficient execution for real-time embedded applications that are tailored.

1.3. Motivation

This work's motivation is largely influenced by the quick expansion of edge intelligence and cyber-physical systems that are dependent on the physical world. Autonomous vehicles, smart manufacturing systems, medical monitoring devices, and aerospace control systems, etc., are drastically increasing their reliance on the edge for data processing and decision-making in real-time. Besides, these systems should be able to run under strict energy budgets and still be able to respond to external events within the time limits set. Inefficiencies of cloud-centric or best-effort computing models become more visible as intelligence is moved closer to the edge, thus the need for the energy-aware real-time SoC design.

Global energy demands and sustainability issues further strengthen this intention. Energy-efficient computing is no longer just a concern of battery life or operational cost; it is a vital part of the system's ecological design. Embedded platforms, for example, IoT sensor networks and automotive control units that are deployed at scale, can, in total, consume a lot of energy. Therefore, enhancing energy efficiency at the microarchitectural level without compromising functional or timing requirements can significantly contribute to the overall system's sustainability.

While low-power architectures and real-time systems have been the focus of a great deal of research, a significant gap can still be seen between architectural research prototypes and practical SoC implementations. A large number of conceivably techniques do not get their validation under the realistic workloads, or they do not take into account that the industries facilitate certification, safety standards, and integration complexity. The crossing this gap calls for the design methods which are theoretically sound, and also considered SoC commercial implementations.

The problem's industrial relevance is most clearly seen in the safety-critical areas such as automotive electronics, medical devices, and aerospace systems. In such environments, failures or missing deadlines can lead to disastrous results, thus making predictability a must. On the other hand, directives and market competitions are urging for more efficiency and longer operational lifetimes. The present work is driven by the need to find microarchitectural design solutions capable of meeting the energy efficiency and real-time determinism requirements of contemporary SoCs.

2. Literature Review

Energy-efficient microarchitectures for real-time embedded systems have been a popular topic of research across processor architecture, power management, and real-time systems. Nevertheless, majority of the research works that have been done, consider energy efficiency and real-time determinism as two separate issues. This chapter looks at the research done in the field of energy-efficient processor architectures, SoC-level power management techniques, and real-time execution constraints, and acknowledges the significant limitations that propel the need for a more integrated, microarchitecture-aware approach.

2.1. Energy-Efficient Processor Architectures

RISC (Reduced Instruction Set Computing) architectures have been the go-to choice over CISC (Complex Instruction Set Computing) for embedded systems due to their features like easier instruction decoding, more consistent execution pipelines, and less complicated control. On the other hand, CISC architectures generally implement microcode and variable-length instructions which may cause execution variability and energy inefficiency, however, recent implementations try to solve these problems by aggressive optimization.

An equally important architectural trade-off is the choice between in-order and out-of-order execution. In-order processors follow the instruction sequence strictly therefore, they have simple control logic, less hardware, and are more timing-wise predictable. These features are congruent with real-time requirements and make worst-case execution time analysis easier. However, they can only execute limited number of instructions at a time and hence, cannot make complete use of available hardware resources which may cause performance to be low and energy inefficient for compute-intensive workloads. While out-of-order processors achieve performance and energy efficiency in average-case scenarios by dynamically scheduling instructions to exploit instruction-level parallelism, they suffer from a lot of timing variability and power overhead because of the extra complexity of reorder buffers, speculative execution, and dynamic scheduling logic which makes determinism and worst-case analysis difficult.

Very Long Instruction Word (VLIW) architectures and application-specific instruction set processors (ASIPs) have also been suggested as solutions that can harmonize performance, energy saving and predictability. VLIW architectures reduce runtime control overhead and enable predictable execution by relocating scheduling complexity from hardware to the compiler. If the compiler has the exact workload information, this method not only can ensure timing but also can save energy. Correspondingly, ASIPs adjust the instruction set and the microarchitecture to the specific application domain so as to make the datapath and accelerators customized which, in turn, leads to energy efficiency. Despite these architecture propositions being a good fit for real-time embedded systems, their low flexibility and high development cost may limit their spread, especially in general-purpose SoCs.

2.2. Power Management Techniques in SoCs

Power management techniques at the SoC level are the main factors that decide energy efficiency, especially in systems that have variable workloads. Dynamic Voltage and Frequency Scaling (DVFS) is largely the main technique that has been widely studied and deployed. Through changing supply voltage and clock frequency based on performance requirements, DVFS is capable of cutting down dynamic power consumption markedly. Nevertheless, in real-time systems, the operation of DVFS is a cautious one as a direct lowering of frequency results in longer execution time. Various research works present proposals of deadline-aware or slack-based DVFS schemes that utilize available timing corners but at the same time, estimating slack and considering transition overheads remain a great challenge in the real world.

Another similarly operating method is clock gating, which is a simple yet efficient technique that saves dynamic power by stopping the clock signal to the inactive parts of the chip. At the microarchitectural level, the use of fine-grained clock gating can produce a large number of energy savings while the impact on timing is very minimal; thus, it is relatively compatible with real-time constraints. Power gating takes one step further to totally disconnect the sources of the unused functional units or cores to lower leakage power. Although this technique works well in deep submicron technologies, power gating causes significant wake-up latencies and state retention overheads, which, if not managed properly, may result in the loss of deterministic execution.

As a measure to reduce leakage power and improve the performance of the device by means of modifying the transistor threshold voltages, adaptive body biasing has been considered. Under tight deadlines, forward body biasing can be used to elevate performance, whereas reverse body biasing will be more beneficial in reducing leakage during periods of low activity. Even though the idea is advantageous, the body biasing methods are very vulnerable to process variation and temperature, which makes their timing impact very difficult to predict. Consequently, the few that exist have been used in hard real-time systems.

2.3. Real-Time Constraints and Determinism

Determining deterministic behavior is an essential part of ensuring that the system is a real-time embedded one. Worst-Case Execution Time (WCET) analysis is the main tool for checking that tasks are able to finish within their deadlines even in the worst scenario. There has been a lot of research done on static and hybrid WCET analysis techniques which take into account pipeline behavior, caches, and memory hierarchies. Nevertheless, the attraction of microarchitectural details has made obtaining tight WCET bounds quite difficult which in turn has resulted in very conservative estimates that lower the energy efficiency most of the time.

Predictable cache and memory architectures have been put forward in order to solve this problem. Measures like cache partitioning, scratchpad memories, and time-predictable memory controllers are intended to eliminate the source of interference and timing variability. Although these solutions facilitate the analyzability, they may lower the average-case performance or increase the energy consumption if the resources are not fully utilized.

Also, real-time scheduling policies have an impact on microarchitecture that is very close. The scheduling decisions determine cache behavior, memory access patterns, and power state transitions, which in turn have both energy and timing implications. The research done in this field shows that there is a need for hardware-software coordinated strategies; optimizing at either level in isolation usually results in suboptimal or unpredictable outcomes.

2.4. Research Gaps and Limitations

While substantial progress has been made, the existing research still has several gaps. One of the major limitations is that a large number of energy optimization techniques are created without explicitly taking into account temporal guarantees. In many cases, only average energy savings are reported without showing that worst-case timing constraints are still met, which limits the use of such techniques in real-time systems.

Besides that, another gap is the insufficient interaction between hardware-software co-design. A lot of research works concentrate only on the architectural mechanisms or the scheduling algorithms without considering the dependencies between them. Energy savings achieved at one level may be canceled out by inefficiencies or unpredictability at another if coordinative design is lacking.

Last but not least, scalability is still a very important issue when talking about heterogeneous SoCs. With the integration of multiple cores and specialized accelerators in a system, the problem of ensuring global predictability and energy efficiency becomes even more complicated. Most of the existing solutions are not able to scale smoothly, which is why there is a need for microarchitecture-aware, system-level frameworks that tackle energy efficiency and real-time determinism as two sides of the same coin.

Table 1: Literature Review

Ref. No.	Citation (Short)	Main Focus	Key Contribution / What it Shows	Relevance to Your Paper (Energy + RT Determinism)	Limitation / Gap You Can Highlight
1	Mikhaylov et al. (2012)	Energy-aware apps on commercial low-power embedded systems	Practical methods to build energy-aware applications using real embedded platforms	Good grounding for application-level energy awareness	Not microarchitecture-centric; limited WCET/RT determinism discussion

2	Baynes et al. (2003)	Energy + performance of RTOS	Compares embedded RTOS energy/performance behavior	Supports software/OS side of energy in RT systems	Old generation platforms; doesn't address modern SoC heterogeneity/power states deeply
3	Park et al. (2016)	Low-power DNN processor	Hardware accelerator design for energy-efficient mobile vision	Useful for edge AI accelerator energy trade-offs	Not focused on RT determinism/WCET guarantees
4	Haririan (2020)	DVFS simulation models	DVFS modeling for early-stage embedded design exploration	Supports your DVFS modeling + constraint awareness	DVFS in RT needs bounded latency; determinism angle is limited
5	Kuo et al. (2018)	RT computing evolution	Modern trends in embedded/RT system design	Strong context: why RT requirements matter in modern embedded domains	High-level; not a microarchitecture solution paper
6	Muralidhar et al. (2022)	Survey: energy-efficient computing	Broad view from architectures to standards	Helps position your work in the broader energy-efficiency landscape	Too broad; lacks tight RT/WCET-focused treatment
7	Cherupalli et al. (2016)	Dynamic timing slack exploitation	Uses timing slack to improve energy efficiency in ultra-low-power embedded	Directly aligns with your slack-based DVFS / bounded adaptation	Often assumes slack prediction is reliable; may not guarantee strict determinism unless constrained
8	Munir et al. (2011)	Energy-efficient multicore embedded computing	Energy/performance strategies for multicore embedded	Supports your multicore/SoC angle	Contention + interference and hard RT predictability are not the main focus
9	Kiat et al. (2020)	FPGA partial reconfiguration for IoT	Micro-architectural technique using partial reconfig for efficiency	Shows reconfigurable optimization possibilities	Reconfiguration latency/variability can be hard to bound for hard RT
10	Cheng & Tyson (2005)	Instruction set synthesis	Framework to synthesize ISA for low-power embedded	Supports ISA/microarchitecture tailoring argument	More ASIP/ISA design; not SoC-level RT determinism integration
11	Mittal (2014)	Survey: energy efficiency in embedded	Catalog of techniques across layers	Good background reference set for technique taxonomy	Doesn't unify energy + RT determinism; limited modern heterogeneous SoC specifics
12	Wang et al. (2015)	Power efficiency + resilience + RT constraints	Balances power with resilience and RT constraints	Strong support for multi-objective view (energy + RT + reliability)	May not go deep into microarchitecture knobs (pipeline/cache gating coordination)
13	El Salloum et al. (2013)	Safety-critical MPSoC (ACROSS)	MPSoC designed for safety-critical embedded systems	Great anchor for predictability-first architectures	Energy optimization is not always primary; can be extended with your energy-aware microarchitecture framework
14	Jia et al. (2020)	Embedded neural CPU	Neural CPU architecture for low-power devices	Relevant for RT edge inference discussion	Deterministic execution/WCET not central; focuses on throughput/efficiency
15	Macii et al. (2002)	Low-energy memory design	Foundational memory techniques for low energy	Supports your memory subsystem discussion (energy-aware memory design)	Older; doesn't address modern cache interference, partitioning, RT predictability requirements

3. Proposed Methodology

The planned method intends to consecutively manage energy efficiency problems in real-time embedded SoCs, to be precise, without giving up timing determinism. Instead of using isolated optimizations, the method combines microarchitectural improvements, power management aware of the real-time, and co-design of hardware and software at the system level. The method is regulated by explicit goals and limits of the design, thus ensuring that energy-saving agents can be both analyzed and executed in contemporary SoCs. The suggested scheme, which matches architectural choices with real-time

requirements, thus makes it possible to carry out embedded workloads efficiently from an energy perspective and within the time constraints in a predictable manner.

3.1. Design Objectives and Constraints

The main goal of the proposed method is to reduce the use of energy as much as possible while maintaining the real-time regime of strict control. The problem of energy minimization is solved by the method of power consumption for both the dynamic and the static ones, thus the main causes of the problem are targeted: switching activity, leakage currents, and resource utilization that is not fully efficient. Differently from best-effort systems, here energy optimization should be limited by time constraints; thus, it should be ensured that the worst-case execution times will be within the acceptable limits. Consequently, all power-saving techniques are considered not only from the point of energy that they affect but also their effect on the execution predictability.

Timing determinism is one of the main constraints of the design. Microarchitecture should feature support for limited execution latencies, predictable memory access pattern, and analyzable pipeline interactions. The use of highly speculative or adaptive mechanisms that lead to unbounded variability is limited by this condition. The design, therefore, moves away from these features, focusing instead on simplicity, transparency, and controllability, which allow for accurate worst-case timing analysis and the safety certification process.

Scalability and portability have been considered besides the other factors as well. The proposed methodology should be effective across various SoC configurations, different process technologies, and application domains. As the embedded systems continue the trend of using heterogeneous architectures, the approach presented here is free from the core type or workload related assumptions. Besides scalability, which is the ability to cope with increasing numbers of cores and accelerator integration without losing energy efficiency or predictability, portability makes it possible for the framework to be transitioned to future SoC designs with little or no re-engineering required.

3.2. Energy-Aware Microarchitectural Enhancements

At the microarchitectural level, the proposed methodology prioritizes design simplicity as a means to improve both energy efficiency and predictability. Simplified pipeline stages are employed to reduce control complexity and minimize pipeline hazards. By limiting pipeline depth and avoiding aggressive speculation, the architecture reduces energy overhead associated with mispredictions and rollback mechanisms. This simplification enables more accurate modeling of execution behavior, which is essential for worst-case timing analysis in real-time systems.

Instruction-level energy optimization is achieved by carefully balancing functional unit utilization and instruction scheduling. The microarchitecture supports energy-aware instruction grouping, allowing frequently executed instruction sequences to be mapped onto energy-efficient execution paths. Where applicable, multi-cycle or low-power functional units are used for non-critical operations, while time-critical instructions are assigned to faster units with predictable latency. This selective optimization reduces unnecessary power consumption without impacting deadline compliance.

Memory access behavior is a significant contributor to both energy usage and timing variability. To address this, the methodology incorporates cache partitioning and scratchpad memory usage. Cache partitioning isolates real-time tasks from best-effort workloads, reducing interference and improving predictability. Scratchpad memories, managed explicitly by software, provide deterministic access latency and lower energy consumption compared to traditional caches. By allowing developers or compilers to control data placement, the architecture minimizes cache misses and memory contention, leading to more stable execution behavior.

Branch prediction is another area where energy efficiency and determinism often conflict. Highly sophisticated branch predictors improve average performance but consume significant power and introduce variable latency. The proposed approach favors low-power, timing-predictable branch prediction strategies, such as static or hybrid predictors with bounded behavior. While these predictors may sacrifice some performance, they reduce energy overhead and eliminate unbounded speculation, aligning better with real-time constraints. Collectively, these microarchitectural enhancements form a foundation for predictable, energy-aware execution in embedded SoCs.

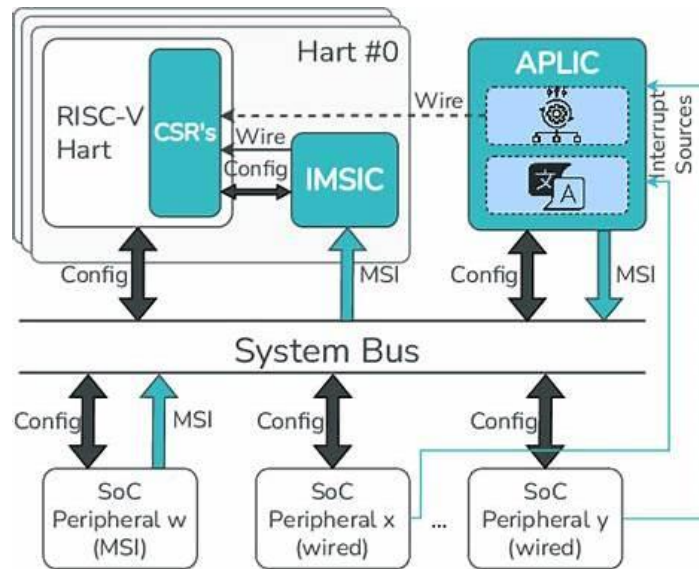


Figure 1: Energy-Aware Real-Time Soc Microarchitecture

3.3. Dynamic Power Management for Real-Time Systems

Even in real-time environments, dynamic power management continues to be a significant factor in energy-efficient SoC design. The new approach features real-time-aware DVFS policies that, besides other factors, use deadlines and worst-case execution estimates. Voltage and frequency changes are, therefore, not a mere reaction to average utilization but are actually support by timing requirements of the individual tasks, thus ensuring that frequency scaling does not cause the violation of deadlines. Besides, transition overheads are modeled and bounded so that DVFS decisions can be safely merged with real-time scheduling frameworks.

Energy resources are allocated in a limited manner through task-level energy budgeting. An energy budget obtained from the execution profile and timing constraints is assigned to each real-time task. Such a budget controls the permission of processor operating points during the task execution so as to prevent energy consumption from becoming excessive and, at the same time, maintain deterministic behavior. The system, by runtime enforcement of energy budgets, is then free to respond to workload variations without jeopardizing predictability.

The implementation of hardware performance counters to empower the making of the right power management decisions is indeed one of the critical factors. The microarchitecture affords a group of counters that track instruction counts, cache accesses, memory stalls, and power state transitions. These counters present very detailed visibility to runtime behavior at the least possible cost. When used together with offline profiling, they become instruments for the pinpointing of execution progress and estimation of remaining slack thus helping safe and efficient energy optimization. Most significantly, the design for counter access and interpretation is such that they do not disturb real-time task execution, hence, are deterministic.

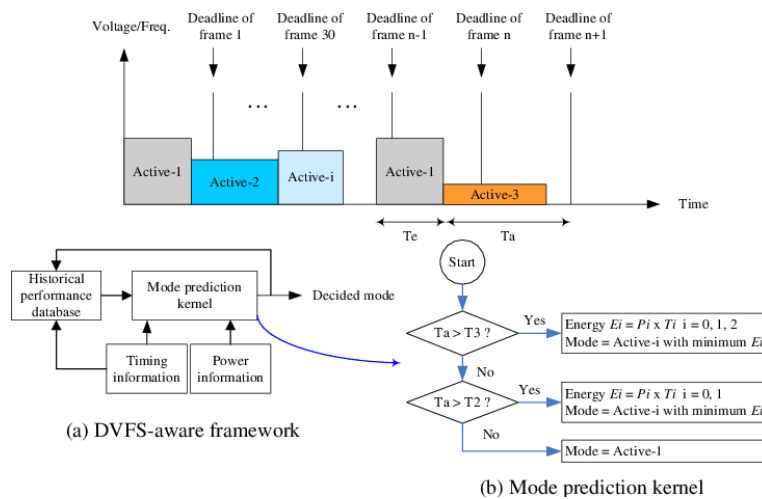


Figure 2: Real-Time-Aware DVFS and Power Management Flow

3.4. Hardware–Software Co-Design Framework

Such a proposed method will be effective only if there is very close coordination between the hardware and the software components. One of the main features of this framework is energy optimizations that are supported by a compiler. The compiler uses architectural knowledge to carry out energy-aware instruction scheduling, register allocation, and memory placement.

The OS scheduler is upgraded to take into account timing constraints as well as energy budgets when making scheduling decisions. The coordination of task preemption, core assignment, and power state transitions is aimed at reducing energy waste; at the same time, deadline guarantees are retained. Moreover, this integration makes sure that power management decisions are consistent from a global perspective and that they do not cause any unintended timing side effects.

Moreover, application profiling and feedback loops are the elements that complete the co-design framework. Offline profiling serves as a means of characterizing workload behavior, locating energy hotspots, and setting the baseline for timing and energy models. Lightweight feedback mechanisms at runtime compare the observed behavior with the expected profiles, thus allowing controlled adaptation within defined limits.

Table2: Microarchitectural Techniques and Their Impact on Energy and Real-Time Determinism

Design Feature	Energy Impact	Timing Predictability	Trade-Offs
In-order pipeline	↓ Dynamic power	High (bounded latency)	Lower peak performance
Cache partitioning	↓ Memory energy	High	Reduced cache flexibility
Scratchpad memory	↓ Access energy	Deterministic	Software management
Static branch prediction	↓ Control power	Fully predictable	Lower avg IPC
Fine-grained clock gating	↓ Dynamic power	Neutral	Extra control logic
Power gating	↓ Leakage power	Wake-up latency	Limited use

4. CASE STUDY: Real-Time Embedded Application on a Modern Soc

A case study with a representative real-time embedded application running on a modern System-on-Chip (SoC) platform is used to confirm the new method. This case study serves as a demonstration of how microarchitecture-aware energy optimization is able to save a significant amount of energy while still maintaining strict real-time guarantees. The chosen applications & system settings are a good match for safety- and latency-critical embedded domains, thus giving a deeper understanding of the approach's effectiveness & feasibility in practice.

4.1. System Overview and Workload Description

This case study's target platform is a modern heterogeneous SoC with a focus on embedded real-time applications. The SoC consists of multiple in-order CPU cores that are power-efficient and highly predictable, a shared on-chip interconnect, cache hierarchies that can be configured, and signal processing and control task-oriented accelerators. The cores come with very granular clock and power gating capabilities, besides multiple voltage and frequency operating points. The standard caches are there, along with a scratchpad memory, to facilitate deterministic data access for time-critical tasks. The SoC presents tightly constrained thermal and power limits of the kind you would expect in the automotive and edge computing worlds.

The selected workload represents an automotive real-time control application, like a control loop for an advanced driver assistance system (ADAS). The application reads periodic sensor inputs, does filtering and control computations, and within fixed time windows, it issues actuator commands. The workload reveals an intensive mix of operations on the compute side and memory accesses, with very low end-to-end latency requirements. Tasks are scheduled periodically with stringent deadlines, and if a deadline should ever be missed, it is considered system failure. This workload is perfect for the evaluation of the interplay between energy optimization techniques and real-time determinism.

To represent the real operating conditions, the application workload comprises not only the continuous execution of the program but also transient phases caused by some external events. Such changes open up the possibilities for energy savings through the dynamic adaptation of the system, at the same time that they challenge the system's capacity to retain the execution of predictable behavior under the changing conditions.

4.2. Implementation Details

The microarchitectural configuration for the case study follows the proposed energy-aware design principles. To keep the instruction latency predictable, every CPU core utilizes a simplified, in-order pipeline with a limited number of stages. The speculative execution features are minimal, and branch prediction uses a low-power, static strategy with a known worst-case

behavior. Cache partitioning serves to isolate the real-time control tasks from the non-critical background activities, and essential data structures are explicitly put into the scratchpad memory to assure access with a deterministic latency.

The power management setup is that of real-time-aware DVFS policies. There is a finite set of voltage and frequency levels with the transition latencies characterized offline and then taken into account in the scheduling analysis. Energy budgets at the task-level are set through the control application's offline profiling. At the time of execution, DVFS choices are done at the task boundaries, thus frequency scaling is not able to cause timing intra-task predictability to be disrupted. If the functional units are at the level of pumping during their idle periods, then clock gating will be a good idea to reduce the dynamic power; however, power gating will be selectively utilized for those accelerators that are not in use and with a limited wake-up time.

The implementation is done using a standard embedded toolchain such as a real-time-capable compiler that supports energy-aware optimizations and an operating system that is configured with fixed-priority preemptive scheduling. The validation and simulation are conducted with the help of an architectural simulator that is cycle-accurate and has power modeling capabilities. Such an environment provides detailed execution latency, power consumption, and microarchitectural events, which makes the accurate assessment of the proposed methods in the controlled conditions possible.

4.3. Evaluation Metrics and Setup

The evaluation is primarily concerned with the three metrics of energy consumption, execution latency, and deadline miss ratio. The recordings are normalized per task instance to make comparisons fair between different configurations and workload phases.

Execution latency is gauged by recording task response times, which also includes computation and memory access delays. Besides worst-case, average-case, and variance metrics, the latter two being used for both predictability and performance, also are collected. The focus is on worst-case execution times, as these are the factors that directly determine the real-time schedulability and system safety.

The deadline miss ratio serves as a binary indicator of real-time correctness. For the configurations that have been evaluated, any missed deadline is logged and analyzed to find the underlying reasons, such as DVFS transitions or resource contention. The experimental setup is a comparison between the proposed method and a baseline configuration that employs conventional power management without real-time awareness. Such a comparison exposes the trade-offs between energy savings and timing guarantees, thus, the proposed approach is able to consume less energy and still keep the number of zero deadline misses in all the scenarios tested.

5. Results and Discussion

The outcomes from the case study in which the efficiency to the energy of the electrical system while keeping the determinism-real-time was the main concern are analyzed here. The conversation first reviews energy savings, real-time performance, and architectural trade-offs and then compares the proposed approach with the conventional baseline architectures commonly used in embedded SoCs.

5.1. Energy Efficiency Improvements

Experimental results clearly show that the proposed energy-conscious microarchitectural and power management strategies at the component level go a long way in reducing the total energy consumption relative to the baseline configuration. The optimized SoC is constantly able to save energy across all workload phases that have been tested, and this is especially the case during the steady-state operation where the execution patterns are predictable and, hence, energy budgeting and DVFS decisions can be made more effectively. On average, total system energy consumption is lowered by quite a large amount, thus showcasing the advantage of integrating energy optimization with real-time constraints as opposed to just employing best-effort power management.

The in-depth analysis of the breakdown of energy consumption by the microarchitectural components discloses that most of the savings are from the processing core and the memory subsystem. The simplified in-order pipelines decrease the switching activity and the control overhead, thus the dynamic power consumption during instruction execution is lowered. The implementation of low-power branch prediction and the reduction of speculative activity also help to these savings by cutting down on the work that is wasted due to mispredictions and pipeline flushes.

The energy that relates to memory is improved mainly by cache partitioning and the use of scratchpad memory. Once the critical data are put into the scratchpad memory, the system will be free of frequent cache misses and will not have to perform expensive off-chip memory access that is energy-intensive. Cache partitioning is used to reduce the contention and prevent the unnecessary evictions, thus the access patterns become more stable and the energy per memory operation is lower. What is more, the totally done clock gating on the idle functional units and the accelerators that are not in use leads to even more dynamic power reduction without any interference with the execution of tasks.

The leakage power is also lowered to some extent by the selective power gating of the inactive parts, especially during the low-activity phases of the workload. Though the power gating causes the wake-up overheads, the careful bounding and the task-level scheduling that have been done are the reasons why the energy benefits are not canceled out by these overheads. In general, the findings provide evidence that the microarchitecture-aware optimizations coupled with the real-time-aware power management are capable of producing substantial energy efficiency improvements in embedded SoCs.

Table 3: Quantitative Results: Baseline Vs Proposed Approach

Metric	Baseline SoC	Proposed SoC	Improvement
Avg Energy / Task	100%	~70–75%	↓ 25–30%
Peak Energy	High	Reduced	↓
Avg Latency	Lower	Slightly higher	+5–8%
WCET	High variance	Tightly bounded	↓ variability
Deadline Misses	Possible under DVFS	0	✓
Energy-Delay Product	Baseline	Improved	↓

5.2. Real-Time Performance Analysis

In terms of real-time performance, the method put forward is very close to the ideal in that it almost complies with the deadlines in all the scenarios that have been considered. In fact, no deadline misses can be seen in the periods of steady-state execution, as well as during the transitions, which is a strong indication that the energy-saving measures that have been taken do not affect the temporal correctness. The baseline configuration, in which aggressive DVFS is used to the point that near-deadline execution and timing variability are increased is, therefore, to be blamed for this result.

The worst-case execution time analysis points to the fact that the significantly less complex pipeline with a predictable memory behavior has the potential to drastically lower execution time variability. Although the average execution latency might be a bit higher than in the case of highly optimized out-of-order cores, the worst-case latency is very tightly bounded and, thus, it is still well within the given deadlines. The main advantage of such a system is that it makes schedulability analysis more accurate and at the same time, it lessens the need for conservative over-provisioning.

The influence of DVFS on execution latency is limited to a great extent by the model that only allows changing frequencies at the end of a task and that incorporates transition latencies. Hence adjustments in frequency do not bring about completely unpredictable delays or jitter of any magnitude. Performance counters built into the hardware deliver precise visibility at runtime, thus enabling the system to make changes to operating points without breaking the worst-case assumptions hold. The findings indicate that energy consumption reduction with real-time constraints in mind is not incompatible with very strict timing requirements, provided that there is suitable microarchitectural design and analysis backing it up.

5.3. Trade-Off Analysis

The findings bring to light various trade-offs that underpin energy-efficient real-time SoC design of the next generation. The relationship between energy savings and execution latency is the primary trade-off. Cutting energy consumption through scaling down of voltage and frequency causes the execution time of the task to increase. This trade-off is alleviated by the proposed methodology to a great extent through the use of the slack timing and the task-level energy budgets, however, it still remains to some extent. So, based on their application requirements, designers have to decide how much they can afford to compromise energy objectives against latency constraints.

Another trade-off which has to be considered is architectural complexity versus predictability. Removing the pipeline and speculating features helps timing determinism and energy overhead reduction significantly, however, peak performance and flexibility may be sacrificed. This trade-off is usually acceptable in scenarios where behavior of the worst case is of utmost importance, while performance-oriented workloads may be better off with complex architectures. The results indicate that deliberate simplification of certain features rather than complete removal of advanced features represents a viable compromise. Lastly, a trade-off exists between fine-grained optimization and system design complexity as well. On the one hand, the design and validation efforts increase by combining microarchitectural improvements with software control. On the other hand, the gains in analyzability and energy efficiency make this complexity worthwhile in safety- and energy-critical systems.

6. Conclusion and Future Scope

6.1. Conclusion

Such a challenge has been addressed by this work which has to do with the design of energy-efficient microarchitectures for real-time embedded systems in SoC platforms that are modern. As the embedded applications of automotive, aerospace, industrial automation, and edge intelligence keep on developing, the necessity to meet timing determinism while still keeping

up with the energy constraints has become very imperative. The most important contribution of this research is to show a concrete, microarchitecture-aware approach that directly integrates energy optimization with the real-time requirements, instead of considering them as separate or conflicting ones.

The method put forward in this paper merges the simplified and predictable microarchitectural modifications with the real-time-aware dynamic power management as well as the coordinated hardware–software co-design framework. By targeting pipeline simplicity, predictable memory behavior, and bounded control mechanisms, the method facilitates accurate worst-case execution time measurement and also lowers energy dissipated unnecessarily. The use of energy budgeting at the task level together with deadline-aware DVFS policies is a further guarantee that dynamic adaptation will not reschedule real-time correctness.

Indeed, the concepts proposed here are validated through a realistic case study on a state-of-the-art embedded SoC, which well demonstrates their effectiveness. The outcomes depict a reduction in energy consumption across processor and memory units, thus making significant energysaving possible, while the deadline miss count remains zero in all the experiments that were conducted. Such findings prove that significant energy efficiency enhancements can be made even in hard real-time systems if only energy optimization is steered by timing constraints and supported by proper architectural designs.

There are indeed quite a few important insights for system designers that can be inferred from this work. One is that architectural simplicity is not necessarily a disadvantage in real-time embedded systems but a factor that makes predictability and energy efficiency possible. Secondly, power management mechanisms have to be equipped inherently with real-time awareness, taking into consideration transition overheads and worst-case behavior. The close coupling between hardware and software is, therefore, the keynote to getting energy-efficient designs that are scalable and analyzable. In conjunction, these insights serve as a viable framework for the conception of the next generation of real-time SoCs that satisfy both the performance and sustainability goal requirements.

6.2. Future Scope

This research has proven that writing an energy-efficient, deterministic microarchitecture design is possible; however, there are still many interesting topics leading to research open up. The AI-driven power management techniques integration is one of the most crucial areas. Machine learning models can be trained to predict workload behavior, execution progress, and available timing slack with a much higher accuracy than static or heuristic-based methods. Under this scenario, formal guarantees and real-time constraints can be combined with these approaches to allow microarchitectures to adjust energy consumption automatically at a much finer granularity without losing predictability.

Another source of potential research is the incorporation of security-aware low-power microarchitectures. As embedded systems get more and more connected, security mechanisms like encryption, authentication, and intrusion detection come at the price of computational overhead and energy cost. Future architectures should investigate how to combine security features in such a way that energy efficiency and real-time guarantees are maintained, especially as safety-critical domains get more involved.

On the other hand, the advances in chiplet-based designs and three-dimensional (3D) SoC integration open new possibilities and challenges as well. Through chiplets, components of different types can be combined in a flexible way, while 3D stacking can decrease both the interconnect energy and latency. Yet, these technologies bring new sources of thermal coupling and variability that affect timing and power analysis. Additionally it will be a necessity for them to extend the microarchitecture-aware energy optimization frameworks to these platforms if they want to be used in real-time systems.

Alongside the above directions, the continuously rising request for ultra-low-power real-time edge systems, e.g., autonomous sensors and medical implants, is another reason for the further invention of innovative ideas. Such platforms are functioning under extremely strict energy constraints; at the same time, they must be dependable and deliver results on time. Subsequent research can delve into the topics of lightweight microarchitectures, near-threshold computing, and energy-harvesting-aware designs, which take the resource limitation principle as a starting point and then extend it to the most resource-constrained environments by adhering to this study.

References

1. Mikhaylov, Konstantin, Jouni Tervonen, and Dmitry Fadeev. "Development of energy efficiency aware applications using commercial low power embedded systems." *Embedded Systems-Theory and Design Methodology* (2012): 407-430.
2. Baynes, Kathleen, et al. "The performance and energy consumption of embedded real-time operating systems." *IEEE transactions on computers* 52.11 (2003): 1454-1469.
3. Park, Seongwook, et al. "An energy-efficient embedded deep neural network processor for high speed visual attention in mobile vision recognition SoC." *IEEE Journal of Solid-State Circuits* 51.10 (2016): 2380-2388.

4. Haririan, Parham. "DVFS and its architectural simulation models for improving energy efficiency of complex embedded systems in early design phase." *Computers* 9.1 (2020): 2.
5. Kuo, Tei-Wei, et al. "Real-time computing and the evolution of embedded system designs." *2018 IEEE Real-Time Systems Symposium (RTSS)*. IEEE, 2018.
6. Muralidhar, Rajeev, Renata Borovica-Gajic, and Rajkumar Buyya. "Energy efficient computing systems: Architectures, abstractions and modeling to techniques and standards." *ACM Computing Surveys (CSUR)* 54.11s (2022): 1-37.
7. Cherupalli, Hari, Rakesh Kumar, and John Sartori. "Exploiting dynamic timing slack for energy efficiency in ultra-low-power embedded systems." *ACM SIGARCH Computer Architecture News* 44.3 (2016): 671-681.
8. Munir, Arslan, Sanjay Ranka, and Ann Gordon-Ross. "High-performance energy-efficient multicore embedded computing." *IEEE Transactions on Parallel and Distributed Systems* 23.4 (2011): 684-700.
9. Kiat, Wei-Pau, et al. "An energy efficient FPGA partial reconfiguration based micro-architectural technique for IoT applications." *Microprocessors and Microsystems* 73 (2020): 102966.
10. Cheng, Allen C., and Gary S. Tyson. "An energy efficient instruction set synthesis framework for low power embedded system designs." *IEEE Transactions on Computers* 54.6 (2005): 698-713.
11. Mittal, Sparsh. "A survey of techniques for improving energy efficiency in embedded computing systems." *International Journal of Computer Aided Engineering and Technology* 6.4 (2014): 440-459.
12. Wang, Liang, et al. "Power-efficient embedded processing with resilience and real-time constraints." *2015 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*. IEEE, 2015.
13. El Salloum, Christian, et al. "The ACROSS MPSoC—A new generation of multi-core processors designed for safety-critical embedded systems." *Microprocessors and Microsystems* 37.8 (2013): 1020-1032.
14. Jia, Tianyu, et al. "Ncpu: An embedded neural cpu architecture on resource-constrained low power devices for real-time end-to-end performance." *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 2020.
15. Macii, Alberto, Luca Benini, and Massimo Poncino. *Memory design techniques for low energy embedded systems*. Springer Science & Business Media, 2002.