# Research-Level Kernel Component Proposals for Advanced Systems

## Executive Summary: Bridging the Gap Between Hardware and the Operating System

Modern computing systems, from large-scale data centers to specialized robotics and devices, rely on a mix of powerful components like GPUs, specialized CPUs, and fast memory. The standard operating system (OS) kernel, designed to be general-purpose, is struggling to keep up with this complexity. This results in three major problems:

1. **Wasted Performance:** The kernel's generic policies lead to inefficient use of specialized hardware (like GPUs) and poor handling of high-speed data, slowing down mathematical and I/O-intensive tasks.
2. **Unpredictable Timing (Jitter):** Even when dedicated CPU cores are used, hidden activity in the kernel causes sudden, unpredictable delays, which is disastrous for real-time systems like circuit simulations or robotics.
3. **Security Vulnerabilities:** Traditional security relies on software layers that can be compromised. New hardware security features are available, but the kernel needs to be completely re-engineered to effectively use them.

This report outlines nine distinct, research-level kernel component projects designed to address these gaps. Each project focuses on adapting the OS kernel to the specific demands of modern, heterogeneous hardware in areas like **efficiency, determinism, and hardware-enforced security.**

---

# I. Boosting Performance and Efficiency in Heterogeneous Systems (GPU, CPU, Math)

These projects aim to make the kernel smarter about resource scheduling and memory usage

to maximize the performance of specialized workloads (like physics simulations and GPU tasks) and manage power consumption more intelligently.

## Project 1: Kernel Segment Monitor (K-SeM)

**Focus:** GPU Scheduling, Real-Time Performance

| Area | Description |
|---|---|
| **The Problem** | When a high-priority GPU task needs to run, the OS kernel typically dedicates the GPU entirely to that task. However, if the high-priority task spends a long time performing calculation on the **CPU** before sending its next command to the GPU, the **GPU sits idle**, wasting valuable resources.[1] |
| **Existing Solution Limits** | One approach allows the kernel to manage GPU access without modifying the application code, but it is prone to this underutilization problem. Other solutions require inserting special code "macros" into the application, which is impractical for large, pre-compiled software libraries.[1] |
| **Novel Kernel Component** | The **Kernel Segment Monitor (K-SeM)** is a smart extension to the GPU scheduler. It transparently monitors the high-priority task's execution on the CPU. If K-SeM detects that the task is in a long CPU-only phase, it **dynamically allows a low-priority task to briefly use the GPU** (a "dynamic filler").[1] This eliminates GPU underutilization without requiring any changes to the application's source code, maintaining system predictability while improving efficiency.[1] |

# Project 2: Communication-Avoiding Page Policy (CAPP)

**Focus:** Memory Management, Mathematical/Physics Simulation

| Area | Description |
|---|---|
| **The Problem** | Mathematical applications, especially those dealing with **sparse matrices** (like in certain physics or circuit simulations), are extremely sensitive to memory speed. Programmers often reorder the data to improve **locality** (keeping related data close to the CPU cache). However, when system memory is low, the generic OS kernel will blindly evict or swap out these critical, highly optimized data pages, forcing the CPU to fetch them repeatedly from slow memory—a costly process called "communication".[2] |
| **Existing Solution Limits** | Optimization occurs entirely at the application level; the kernel's memory management unit (MMU) is completely unaware of which data pages are critical for high performance.[4] |
| **Novel Kernel Component** | The **Communication-Avoiding Page Policy (CAPP)** modifies the kernel's Virtual Memory Management (VMM). The application uses a new system call to **"hint"** the kernel about which memory pages are part of the critical, reordered data segments. The kernel then **tags** these pages and gives them the lowest possible priority for eviction or swapping. This guarantees that the kernel actively enforces the programmer's locality |

| | optimization, directly translating application-level math optimization into system-level memory policy.[2] |
|---|---|

## Project 3: Reverse-Engineered Lightweight Controller (RLC)

**Focus:** Circuit/Hardware Efficiency, Power Management (DVFS)

| Area | Description |
|---|---|
| **The Problem** | Modern processors use complex, proprietary hardware to manage power consumption by adjusting voltage and frequency (Dynamic Voltage and Frequency Scaling or DVFS).[5] Linux's built-in DVFS governors are too generic and cannot achieve optimal energy savings for specific workloads, especially those with soft deadlines (like many embedded or IoT systems).[6] Furthermore, security vulnerabilities have been found where malicious drivers can exploit software control over DVFS to perform attacks (CLKSCREW).[7] |
| **Existing Solution Limits** | The internal logic of hardware-managed DVFS is hidden and proprietary, making it impossible for the OS to write a truly optimal policy.[5] |
| **Novel Kernel Component** | The **Reverse-Engineered Lightweight Controller (RLC)** is a custom kernel module that uses privileged instructions (Model-Specific Register, or MSR, access) and high-precision timers to **sample the processor's internal state at an extremely high frequency** (hundreds of |

| | |
|---|---|
| | KHz).[5] This effectively "steals" the proprietary hardware model. This data is then used to train and deploy a **lightweight, data-driven power-management policy** (like a machine-learned controller) back into the kernel.[6] The project also requires a mandatory security analysis of the controller to implement **mitigation strategies** against attacks that exploit power management.[7] |

# II. Eliminating Jitter and Boosting I/O Speed (Determinism and Zero-Copy)

These projects address the performance bottlenecks and unpredictability caused by the standard OS architecture when handling extremely high-speed data flow and real-time processing requirements.

### Project 4: Hybrid Kernel-Bypass Framework (HBF)

**Focus:** Zero-Copy I/O, Kernel Bypass

| Area | Description |
|---|---|
| **The Problem** | The process of copying data between application memory and the kernel memory (system calls) is a massive performance bottleneck for high-throughput I/O (e.g., networking or high-speed storage).[8] **Kernel bypass** techniques grant applications direct access to I/O devices to avoid these copies |

| | |
|---|---|
| | and context switches.[10] However, this direct access can introduce severe security vulnerabilities, such as memory unmapping exploits.[11] |
| **Existing Solution Limits** | Developers must choose between extreme speed (kernel bypass) and security (traditional kernel I/O).[10] |
| **Novel Kernel Component** | The **Hybrid Kernel-Bypass Framework (HBF)** aims for a middle ground. It uses **designated memory buffers** that are pre-registered with a minimal kernel module. Instead of a full kernel transition, a system call performs a single, fast check on the pointer within the kernel's control interface.[11] This greatly reduces the overhead of context switching while ensuring that the application can only access memory within the verified, designated buffers, mitigating the security risk of prior kernel bypass work.[11] |

## Project 5: Adaptive Core Cohesion Agent (ACCA)

**Focus:** Zero-Copy I/O, CPU Affinity

| Area | Description |
|---|---|
| **The Problem** | Even when using zero-copy techniques (avoiding memory copies), performance can be ruined if the data is processed by different CPU cores in the kernel and user space. This forces the cores to continuously synchronize their caches, creating massive, time-consuming cache coherence traffic.[12] This is a "semantic |

| | |
|---|---|
| | gap" between how the kernel handles I/O (networking) and how the application (like a storage engine such as SPDK) processes it. |
| **Existing Solution Limits** | Current solutions rely on static configuration (like hard-pinning processes to cores), which fails in dynamic, multi-threaded, or cloud environments.[12] |
| **Novel Kernel Component** | The **Adaptive Core Cohesion Agent (ACCA)** is a kernel agent that dynamically monitors cache traffic and core pressure associated with high-speed I/O. Based on real-time metrics, ACCA actively and intelligently adjusts two things: 1) The core affinity of the kernel's I/O processing threads, and 2) The memory allocation so that the buffer is physically and logically **cohesive** (on the same core) with the user application thread. This ensures that the data is processed on the same core throughout the entire path, guaranteeing sustained zero-copy performance.[12] |

## Project 6: Dynamic Cache Partitioning Daemon (DCPD)

**Focus:** Hardware Efficiency, Cache Management

| Area | Description |
|---|---|
| **The Problem** | In multi-threaded systems, different processes can interfere with each other by accidentally filling up the shared CPU cache (L3/L2) with their own data, pushing out data needed by another, more critical task. This causes performance variability.[13] |

| | |
|---|---|
| **Existing Solution Limits** | The **Page Coloring** technique can be used to isolate critical data structures into separate cache partitions, but this currently requires explicit, manual control by a programmer—a method that is not feasible for general-purpose or virtualized systems hosting mixed workloads.[2] |
| **Novel Kernel Component** | The **Dynamic Cache Partitioning Daemon (DCPD)** is a kernel process that automates and dynamically manages Page Coloring. Using specialized resource monitoring (like cgroup QoS or access pattern traces), the DCPD automatically assigns and **recolors** pages in real-time, effectively and dynamically allocating a guaranteed fraction of the cache to high-priority tasks. This extends the deterministic performance benefits of cache partitioning from highly specialized single-application HPC to dynamic, mixed-workload environments.[2] |

## Project 7: Adaptive Control Module (ACM)

**Focus:** Determinism, OS Noise (Circuit Simulation/Real-Time)

| Area | Description |
|---|---|
| **The Problem** | **OS Noise** (unpredictable delays or "jitter") is the primary enemy of real-time systems, like high-frequency trading or industrial robotics. Even when dedicated cores are statically isolated for critical tasks, jitter persists due to high-rate contention on **fundamental, shared kernel resources**, such as the global memory allocator lock |

| | |
|---|---|
| | (zone->lock).[14] |
| **Existing Solution Limits** | Current isolation methods rely on static configuration (e.g., using isolcpus at boot).[16] They minimize temporal noise but cannot eliminate contention on shared data structures.[15] |
| **Novel Kernel Component** | The **Adaptive Control Module (ACM)** implements a **dynamic, closed-loop feedback system** in the kernel. It continuously measures scheduler latency jitter (OS noise) using low-overhead tracing tools.[18] If the measured jitter exceeds a safe threshold, the ACM triggers immediate, dynamic corrective actions. These actions could include temporarily prioritizing the real-time task, throttling non-critical kernel work, or redirecting memory allocation requests from the isolated core to a pre-allocated, lockless pool, thus completely bypassing the contended global lock (zone->lock) until stability is restored.[14] |

# III. Next-Generation Security: Hardware-Enforced Trust and Safety

These projects focus on re-engineering the kernel to effectively utilize and combine new processor security features (like hardware memory tags and trusted execution environments) to protect sensitive memory and data flow.

## Project 8: Tagged Page Manager (TPM)

**Focus:** Hardware Security (MTE), Memory Allocation

| Area | Description |
|------|-------------|
| **The Problem** | Newer ARM processors support the **Memory Tagging Extension (MTE)**, a hardware feature that tags memory and pointers to detect memory safety bugs (like buffer overflows).[19] To be effective, the kernel's memory allocator must assign tags randomly for security. This conflicts with the allocator's main goal: reducing memory fragmentation to utilize large, fast pages (**hugepages**).[21] The allocator must choose between performance and security. |
| **Existing Solution Limits** | Standard allocators are not "tag-aware" and cannot intelligently balance the need for contiguous pages (performance) with the need for highly randomized tags (security).[20] |
| **Novel Kernel Component** | The **Tagged Page Manager (TPM)** is a new MTE-aware kernel memory allocator. It manages physical memory based on both contiguity (for hugepages) and **tag value** (for security). The TPM implements a **tag-aware compaction strategy** that only coalesces memory blocks if they have compatible tags.[19] This minimizes the risk of tag collisions while maximizing the availability of large, contiguous blocks, ensuring both robust security and high-performance backing for memory allocations. |

# Project 9: Secure I/O Mediator (SIM)

**Focus:** Confidential Computing, TEE Architecture

| Area | Description |
|---|---|
| **The Problem** | **Trusted Execution Environments (TEEs)** like Intel SGX are excellent at protecting highly sensitive code and data **inside an application** against a malicious OS.[22] However, SGX **lacks support for generic trusted I/O paths**, meaning data coming into or leaving the secure area (e.g., from a network card or sensor) is vulnerable to interception by a malicious OS or hypervisor.[23] |
| **Existing Solution Limits** | The other major TEE, ARM TrustZone, provides system-wide security and **hardware support for trusted I/O paths**, but its isolation is broader and more system-level.[22] |
| **Novel Kernel Component** | The **Secure I/O Mediator (SIM)** proposes a hybrid architecture to combine the best features of both TEEs. A minimal, highly verified kernel component (the SIM) is deployed within the **TrustZone Secure World**. This SIM handles only the trusted device drivers and I/O channels.[23] It then establishes an attested, secure channel (using shared, protected memory) to the application's secure computation running in an **SGX enclave** (in the Normal World).[22] This ensures end-to-end confidentiality: SGX protects the computation, and TrustZone protects the data path. |

**Works cited**

1. Unleashing the Power of Preemptive Priority-based ... - arXiv, accessed on November 19, 2025, https://arxiv.org/abs/2401.16529
2. SMOReS: Sparse Matrix Omens of Reordering Success, accessed on November

19, 2025,
https://scholarship.tricolib.brynmawr.edu/bitstream/handle/10066/7572/2011Woo
dS_thesis.pdf?sequence=1

3. A Novel Cache-Blocked MPI-parallel Matrix Power Kernel: - NHR@FAU, accessed on November 19, 2025,
https://hpc.fau.de/files/2024/10/Dane-Lacey-Masterarbeit-finished.pdf

4. Memory Management - Michele Zanella, accessed on November 19, 2025,
https://zanella.faculty.polimi.it/wp-content/uploads/Es5.pdf

5. Reverse Engineering DVFS Mechanisms - Tanvir Ahmed Khan, accessed on November 19, 2025,
https://takhandipu.github.io/papers/piersma-reverse-2025.pdf

6. Department of Computer Science Studying Practical Kernel-Level DVFS Techniques for Deadline- Constrained Periodic Applications - St. Francis Xavier University, accessed on November 19, 2025,
https://www.stfx.ca/sites/default/files/documents/Computer_Science-ZhouTi_Aug
92022_proposal.pdf

7. CLKSCREW: Exposing the Perils of Security-Oblivious Energy Management, accessed on November 19, 2025,
https://www.cs.columbia.edu/~simha/preprint_USENIX17_clkscrew.pdf

8. Zero-copy - Wikipedia, accessed on November 20, 2025,
https://en.wikipedia.org/wiki/Zero-copy

9. Zero-copy: Principle and Implementation | by Zhenyuan (Zane) Zhang | Medium, accessed on November 20, 2025,
https://medium.com/@kaixin667689/zero-copy-principle-and-implementation-9a
5220a62ffd

10. Safe Sharing of Fast Kernel-Bypass I/O Among Nontrusting Applications - arXiv, accessed on November 20, 2025, https://arxiv.org/html/2509.02899v1

11. Safe Sharing of Fast Kernel-Bypass I/O Among Nontrusting Applications - arXiv, accessed on November 20, 2025, https://arxiv.org/pdf/2509.02899

12. SPDK NVMe over TCP Optimization on Arm, accessed on November 20, 2025,
https://developer.arm.com/community/arm-community-blogs/b/servers-and-clou
d-computing-blog/posts/spdk-nvme-over-tcp-optimization-on-arm

13. (PDF) Controlling Cache Utilization of HPC Applications, accessed on November 20, 2025,
https://www.researchgate.net/publication/221236170_Controlling_Cache_Utilizatio
n_of_HPC_Applications

14. Meet osnoise, a better tool for fine-tuning to reduce operating system noise in the Linux kernel | Red Hat Research, accessed on November 20, 2025,
https://research.redhat.com/blog/article/osnoise-for-fine-tuning-operating-syste
m-noise-in-linux-kernel/

15. Locked In, Leaked Out: Measuring Isolation via Kernel Locks - arXiv, accessed on November 20, 2025, https://arxiv.org/html/2507.21248v1

16. real time - Realtime OS: PREEMPT_RT Linux vs QNX and other - Stack Overflow, accessed on November 20, 2025,
https://stackoverflow.com/questions/77618622/realtime-os-preempt-rt-linux-vs-q

[nx-and-other](nx-and-other)

17. Real-Time | DPsim, accessed on November 19, 2025, https://dpsim.fein-aachen.org/docs/getting-started/real-time/
18. Effectively Measure and Reduce Kernel Latencies for Real-time Constraints - Chung-Fan Yang - YouTube, accessed on November 20, 2025, https://www.youtube.com/watch?v=epcPeMlBJW0
19. Arm memory tagging extension | Android Open Source Project, accessed on November 20, 2025, https://source.android.com/docs/security/test/memory-safety/arm-mte
20. Armv8.5-A Memory Tagging Extension - Arm Developer, accessed on November 20, 2025, https://developer.arm.com/-/media/Arm%20Developer%20Community/PDF/Arm_Memory_Tagging_Extension_Whitepaper.pdf
21. Beyond malloc efficiency to fleet efficiency: a hugepage-aware memory allocator - USENIX, accessed on November 20, 2025, https://www.usenix.org/system/files/osdi21-hunter.pdf
22. What are the key differences between Intel SGX and ARM TrustZone in terms of TEEs?, accessed on November 20, 2025, https://massedcompute.com/faq-answers/?question=What%20are%20the%20key%20differences%20between%20Intel%20SGX%20and%20ARM%20TrustZone%20in%20terms%20of%20TEEs?
23. SGXIO: Generic Trusted I/O Path for Intel SGX, accessed on November 20, 2025, http://library.usc.edu.ph/ACM/SIGSAC%202017/codaspy/p261.pdf