

Unified Compute Architecture: A Roadmap for Post-GPU Processing

Abstract

Brief overview of the vision: unify CPU and GPU responsibilities through architectural changes that emphasize shared memory, flexible execution models, and display independence.

1. Introduction

1.1 Motivation: Why Current CPU/GPU Separation is Inefficient

For decades, CPUs and GPUs have existed as largely separate processing units, each with its own responsibilities, memory, and specialized hardware. This separation has led to a variety of inefficiencies that are starting to become more apparent as workloads evolve. Graphics, AI, and compute workloads have become more interdependent, and yet the hardware remains siloed.

- **Memory Duplication:** Both the CPU and GPU often require copies of the same data—leading to increased memory requirements, bandwidth bottlenecks, and wasted resources.
- **Draw Call Overhead:** The traditional graphics pipeline relies on repeated "draw calls," where the CPU passes data to the GPU, which then processes it. This back-and-forth can incur significant latency and power consumption, especially in more complex applications.
- **Rendering and Compute Segmentation:** The GPU is primarily tasked with rendering graphics, while the CPU focuses on computation, logic, and control. However, workloads today often require both types of processing simultaneously, with no clean way to harmonize them.

A more unified approach to CPU and GPU processing would address these issues, allowing for greater flexibility and efficiency.

1.2 Historical Context: Draw Calls, VRAM, and GPU Display Control

Historically, GPUs were designed exclusively for rendering graphics, with VRAM dedicated to holding frame buffers, textures, and other visual data. The CPU's role was focused on application logic, managing operating systems, and controlling hardware.

- **Draw Calls:** In the graphics pipeline, the CPU must send multiple draw calls to the GPU for each object it wishes to render. This creates a **bottleneck in performance** and introduces **latency**.
- **VRAM:** Video RAM (VRAM) is used exclusively for the GPU to store textures, framebuffers, and other graphical data. This means that the CPU and GPU each maintain separate memory pools that hold identical copies of the data.
- **GPU Display Control:** The GPU traditionally controls the final step in the display pipeline, rendering images and passing them to the display hardware. The CPU has limited control over this process, creating another point of inefficiency.

The current design has become increasingly outdated as the lines between CPU, GPU, and other accelerators blur. Newer workloads like **machine learning**, **3D simulations**, and **AI-driven graphics** call for a more flexible, unified approach.

1.3 Why Now? Industry Trends Driving the Change

The landscape of computing has dramatically shifted in the last decade. Mobile devices, AI, and cloud computing have all placed new demands on processing power, requiring a more unified approach to compute.

- **AI/ML and Data Parallelism:** Machine learning models, especially those in fields like natural language processing and computer vision, require both **compute-heavy** workloads and **graphical rendering**. The data parallelism involved in AI workloads makes the separation of CPU and GPU memory and processing inefficient.
- **Mobile-First and ARM SoCs:** ARM-based systems on chips (SoCs), like Apple's M1/M2 chips, have demonstrated that it's possible to combine CPU, GPU, and other accelerators (like neural engines) into a single package with shared memory. These designs have shown **impressive performance gains** due to **unified memory** and tightly coupled components.
- **Cloud and Distributed Computing:** The rise of cloud computing has also exposed the inefficiencies of traditional hardware. Multi-GPU systems in data centers are cumbersome, require specialized APIs, and often suffer from poor resource allocation and data copy overhead.

The next step is to take the lessons learned from ARM's approach and expand it to a broader, more flexible architecture capable of handling a wide variety of workloads. This whitepaper proposes a roadmap for achieving such a unified system.

2. The Problem

2.1 Fragmented Memory Models

One of the most significant issues with current CPU/GPU systems is the **fragmentation of memory** between the CPU's main memory (RAM) and the GPU's dedicated memory (VRAM). This separation leads to a number of inefficiencies:

- **Memory Duplication:** Both the CPU and GPU often maintain copies of the same data. For example, when rendering a scene, textures, buffers, and other resources must be loaded into the GPU's memory, even though the CPU has already processed or owns the same data in system memory.
- **Data Transfer Latency:** Every time data needs to be moved between the CPU and GPU (e.g., transferring textures, buffers, or model data), it incurs latency. This back-and-forth can be significant in time-sensitive applications like gaming or real-time simulation.
- **High Memory Overhead:** With separate memory pools, there is significant waste due to redundant data storage. Furthermore, the size of the GPU's memory (VRAM) often limits the complexity of models and textures that can be processed in real-time.

These inefficiencies demand that memory resources be optimized to be shared by both the CPU and GPU, which would eliminate the need for duplication and reduce latency.

2.2 GPU-Owned Display Path

In traditional architectures, the **GPU** controls the final output to the display. The CPU, in many cases, has no direct involvement in this process unless explicitly instructed to perform operations that affect rendering. There are several key issues with this model:

- **Limited Flexibility:** The current approach assumes that display output always comes from the GPU. However, certain applications (such as software rendering or when using remote display setups) could benefit from having the display controller more directly tied to the system-on-chip (SoC) and not solely dependent on the GPU.
- **Multi-GPU Complexity:** When multiple GPUs are used (e.g., for rendering different portions of a scene in parallel), coordinating the output across multiple GPUs and ensuring a seamless display experience becomes increasingly complex. The current architecture requires careful synchronization of framebuffers between GPUs and often leads to unnecessary overhead.

- **Display Overhead:** The GPU needs to handle both the graphical rendering and the actual display output. While modern GPUs are more than capable of handling this load, it also means that the CPU is entirely dependent on the GPU to handle display tasks. This can lead to unnecessary bottlenecks or wasted power, especially when the GPU isn't fully utilized.

A more unified display path—where the **display controller** and **framebuffer** are part of the system-on-chip (SoC) rather than tied exclusively to the GPU—would allow for more efficient resource management, better flexibility, and lower power consumption.

2.3 CPU-GPU Communication Bottlenecks

Despite the advances in GPU technology, **communication bottlenecks** between the CPU and GPU remain a significant challenge. Traditional CPU/GPU communication works through explicit, often costly, interfaces such as PCIe, which introduces several limitations:

- **Draw Calls and Data Transfers:** The CPU often needs to send data to the GPU in the form of draw calls or buffer transfers. Each of these transfers requires a context switch, further increasing latency and wasting CPU cycles. The more complex the scene or operation, the more frequent these calls become.
- **Lack of Seamless Integration:** Current systems lack a true hardware-level integration between CPU and GPU. The separation means that both units are required to work with entirely different memory models, data types, and operational contexts. This makes it difficult to process workloads in parallel, reducing overall system efficiency.
- **Synchronization Issues:** Both units often need to synchronize and manage data dependencies. For example, the CPU may need to wait for GPU computation to finish before continuing with application logic. This synchronization often leads to inefficiencies and idle time, particularly in workloads that could otherwise be processed concurrently.

A unified system could alleviate these bottlenecks, allowing CPU and GPU to share memory seamlessly and interact more efficiently, leading to a smoother and faster workflow for a variety of applications.

2.4 Underutilized SIMD in CPUs

Current CPUs support SIMD (Single Instruction, Multiple Data) instructions that allow for parallel processing of data. However, **SIMD lanes in modern CPUs are often underutilized**, especially for floating-point-heavy workloads, which are common in tasks like graphics rendering, scientific simulations, and AI.

- **AVX2 and AVX512 Limitations:** AVX2 supports 256-bit vectors, and AVX512 extends this to 512-bit, but both are heavily underutilized for floating-point operations like FP32 and FP64. For example, AVX2 can handle **eight 32-bit floats simultaneously**, but many workloads fail to fully utilize this power.
- **Single-Precision and Double-Precision Efficiency:** Many workloads rely heavily on FP32 (single-precision) operations, yet modern CPUs are still optimized for general-purpose computations rather than highly parallel, floating-point-heavy workloads. This limits the overall throughput of modern processors.
- **GPU-Specific Features Unavailable:** Many modern GPUs support specialized floating-point operations (e.g., FP16, FP32) and the ability to process large datasets in parallel. However, CPUs lack this type of specialized hardware, meaning that many workloads that would benefit from GPU-like acceleration are instead limited to slower, general-purpose processing.

Introducing specialized SIMD support that can be used in more general-purpose workloads would help bridge the performance gap between CPU and GPU, enabling more flexible and efficient processing of parallel workloads across both architectures.

3. Proposed Architecture

3.1 Unified Memory

The foundation of this proposed architecture lies in **Unified Memory**, a crucial design shift where all processing units—whether CPU, GPU, or specialized accelerators—share a single high-bandwidth system memory pool. This eliminates the need for dedicated VRAM (Video RAM) and allows for efficient, seamless access to memory across different processing units. Traditionally, CPUs and GPUs have operated in separate memory spaces, requiring costly and time-consuming data transfers between them. By using **memory-mapped buffers**, this architecture ensures that both the CPU and GPU can access the same memory directly, resulting in significantly reduced overhead and better data locality. The unified memory system optimizes performance, reducing latency and avoiding the need for duplicate copies of data, making it ideal for workloads that involve both CPU and GPU-driven computations.

This approach simplifies the system architecture and enhances flexibility. For instance, a rendering task could be processed by the GPU, while a physics or AI simulation could be handled by the CPU, with both processing units accessing the same memory resources. This allows for more efficient parallel execution and better resource utilization.

3.2 Unified Display Controller

In a conventional system, the **display controller** is often a separate module, located either in the GPU or a discrete hardware block. However, in this unified architecture, the display output is controlled directly by the **System-on-Chip (SoC)**, with the **framebuffer** stored within the same unified memory pool. This approach ensures that the SoC has full control over display output, while still allowing both CPU and GPU to contribute to rendering tasks.

By adopting a single framebuffer architecture, the system eliminates the need for multiple framebuffers or separate memory regions dedicated to different rendering tasks. This centralization allows any rendering unit—whether the CPU or GPU—to write to the same framebuffer, improving performance and simplifying multi-renderer workflows. More importantly, the display controller's ownership by the SoC provides energy efficiency benefits, particularly in scenarios where the GPU is idle or underutilized. The display system can operate independently of the GPU, drawing only on the CPU when the demand for graphics output is low.

3.3 True Simultaneous Multithreading

True Simultaneous Multithreading (SMT) takes a leap forward with the introduction of **SIMD-aware SMT**, enabling more efficient use of the CPU's vector execution resources. Traditional SMT architectures generally execute instructions in a more linear or scalar fashion, with limited parallelism in highly vectorized workloads. The proposed architecture introduces **SIMD-aware scheduling**, dynamically adjusting the number of logical threads based on the instruction set in use. For example, when running an AVX2 workload (using 256-bit SIMD instructions), the architecture could schedule 8 threads concurrently, whereas with AVX512 (which uses 512-bit SIMD instructions), 16 threads could be scheduled, effectively utilizing the wider data paths and maximizing throughput.

This sophisticated thread scheduling ensures that both scalar and vector workloads can be handled optimally, providing higher performance and better parallelism. It allows the architecture to scale with the computational complexity of different instruction types, making it well-suited for a variety of tasks, from high-performance computing to gaming and AI workloads. By leveraging **logical thread scheduling** based on the data granularity of operations, this architecture fully exploits the potential of modern SIMD instructions.

3.4 FP Hierarchy

The architecture introduces a **floating-point hierarchy** to support varying levels of precision for different workloads. With dedicated hardware support for **FP16 (half-precision)**, **FP32 (single-precision)**, and **FP64 (double-precision)** operations, the system is able to execute floating-point tasks more efficiently by selecting the appropriate precision level based on workload demands. This precision flexibility is key in applications

such as scientific computing, machine learning, and graphics, where the balance between accuracy and performance must be dynamically optimized.

In addition to offering hardware-level support for multiple floating-point formats, the architecture enables **dynamic scalar/vector mode switching**. Depending on the type of workload, the CPU can seamlessly switch between scalar execution (for tasks that don't benefit from parallelization) and vector execution (for tasks that benefit from massive parallelism, such as matrix operations or deep learning). This dynamic switching ensures that the system can deliver the best performance for any task, optimizing both throughput and energy efficiency without requiring manual intervention from developers.

3.5 GPU Feature Backports

To further blur the line between CPU and GPU responsibilities, this architecture **backports GPU features** to the CPU. **Texture sampling**, **mipmap generation**, and other GPU-specific tasks are brought to the CPU side, allowing the system to handle these tasks more effectively without needing to offload them to a discrete GPU. For example, texture sampling, traditionally the domain of GPUs, can be performed by the CPU in this architecture, enabling applications to handle complex image processing tasks without the need for a GPU.

Furthermore, **mipmap generation**—which helps optimize texture rendering by providing lower-resolution versions of textures for distant objects—can be performed by the CPU. This reduces the need for the GPU to handle preprocessing, freeing up GPU resources for real-time rendering tasks.

Lastly, the architecture introduces **shader-like abstractions** at the CPU level. This means that the CPU can execute tasks usually reserved for shaders in the GPU, such as parallel computations on large data sets. By backporting shader functionality to the CPU, developers gain more flexibility in choosing how to distribute compute workloads, reducing the need for explicit GPU involvement in some applications.

3.6 Scope-Based Shaders

In traditional graphics pipelines, **draw calls** are used to instruct the GPU to render objects one at a time. However, in this architecture, the **scope of shaders** extends far beyond simple rendering tasks. By merging **logic**, **rendering**, and **compute** into a single pipeline, the system eliminates the need for explicit draw calls. This approach allows the CPU and GPU to collaborate more efficiently, with each task—whether it involves physics simulations, AI processing, or visual rendering—being handled within a unified context.

This unification streamlines the entire pipeline, reducing bottlenecks that arise when switching between different hardware components. The absence of explicit draw calls also results in **lower latency**, as the system can process data continuously, without pausing to send commands between the CPU and GPU. Moreover, this architecture supports **GPU-like dataflows** within software and CPU contexts, allowing developers to seamlessly execute parallel tasks typically reserved for the GPU, directly within the CPU context. This reduces the reliance on a discrete GPU for complex computations, enabling flexible and efficient workloads across the entire system.

With **scope-based shaders**, the architecture opens up new possibilities for developers, enabling them to more easily control the flow of rendering and compute tasks in a unified, high-performance manner.

4. Advantages

4.1 Performance Gains

One of the most significant advantages of this unified architecture is the **performance gains** it offers, particularly in terms of **reduced copy overhead** and **improved cache locality**. In traditional systems, the CPU and GPU operate in separate memory spaces, requiring frequent data transfers between the two, which introduces latency and reduces performance. By adopting a unified memory model, where both the CPU and GPU share a single high-bandwidth memory pool, the need for these costly memory copies is eliminated.

Data can be accessed directly by either processing unit, resulting in faster data transfers and lower overall latency.

Additionally, the unified memory design improves **cache locality**, as both the CPU and GPU access the same memory pool. This reduces the likelihood of redundant copies of data and optimizes cache usage, leading to better performance when handling large, complex datasets or rendering workloads. Tasks like texture mapping, AI inference, and scientific simulations become more efficient, as the data needed by different units is more readily available, enhancing the overall throughput of the system.

4.2 Developer Simplicity

From a developer's perspective, the **simplicity** of the architecture cannot be overstated. Traditionally, developers need to work within distinct environments for CPU and GPU, each with its own API stack and development tools. The proposed unified architecture removes this distinction by providing a **unified API space**. Developers no longer need to learn and manage separate GPU-specific graphics stacks or worry about the intricacies of transferring data between CPU and GPU memory. Instead, they can write code that seamlessly utilizes both CPU and GPU resources, with all processing units accessing the same memory space.

This streamlined approach reduces the complexity of multi-threaded and multi-processing programming, making it easier to write efficient code. By eliminating the need for separate graphics APIs, such as Vulkan, OpenGL, or DirectX for GPU programming, the development process is significantly simplified. Developers can focus on the logic of their applications rather than worrying about managing separate hardware contexts, making cross-platform and cross-hardware development much more accessible.

4.3 Security and Reliability

Security and reliability are critical considerations in modern computing, and the proposed architecture provides significant improvements in this area. In traditional systems, GPUs often act as **black-box** components, with limited visibility or control over what the GPU is doing, especially when it comes to display output. This separation can create security concerns, as malicious code running on the GPU could potentially compromise system integrity without being easily detected.

By integrating the **display controller** within the SoC, the architecture removes the need for a discrete GPU to handle display output. This not only simplifies the system but also **improves security**. Since the CPU has direct control over the display pipeline, it's easier to monitor and enforce security protocols, preventing potential exploits that rely on the GPU's lack of transparency. Moreover, having the display controller within the SoC enhances **reliability**, as the CPU can directly manage system states and interactions, leading to more predictable and stable behavior, particularly in critical applications.

4.4 Scalability

This unified architecture is also highly **scalable**, opening new possibilities for a range of system configurations. On one hand, it makes **small and embedded systems** more viable by reducing the need for discrete, power-hungry GPUs, which are often impractical in low-power or resource-constrained environments. The unified design allows for efficient use of the limited resources available in embedded systems, making it easier to integrate high-performance processing and display capabilities without the need for separate components.

On the other hand, this architecture offers **cleaner multi-GPU setups**, particularly in larger systems. Traditional multi-GPU configurations often suffer from complications such as memory fragmentation, synchronization issues, and inefficiencies in data sharing. By utilizing a unified memory pool and display controller, the overhead associated with multi-GPU configurations is reduced, simplifying setup and improving performance scaling. Multiple GPUs can operate more harmoniously, as they share the same memory and can easily access common resources, leading to more effective scaling for tasks like rendering, machine learning, and computational simulations.

This scalability is a key advantage for industries that require high-performance computing, from gaming to scientific research, where systems can dynamically adjust the number of GPUs based on the workload, providing a flexible and efficient computing environment.

5. Potential Challenges & Solutions

5.1 OS/Kernel-Level Changes for Shared Framebuffer Access

One of the primary challenges in implementing the proposed unified architecture is the requirement for **OS and kernel-level changes** to support **shared framebuffer access**. In traditional systems, the GPU typically handles the framebuffer, and the OS must manage the GPU context to provide display output. However, in this unified architecture, where the display controller is part of the SoC, a more direct connection between the OS and all processing units is necessary to allow access to the framebuffer by both the CPU and GPU.

Solution: To address this challenge, a modification of the operating system's graphics stack is required. The OS must be updated to support shared memory spaces and synchronization between the CPU, GPU, and the display controller. This could involve integrating new kernel modules or drivers that allow both the CPU and GPU to interact with the framebuffer seamlessly, without the need for traditional GPU drivers or APIs.

Additionally, the OS must include robust memory management tools to ensure that both the CPU and GPU access the framebuffer efficiently and without contention.

Kernel-level enhancements would also need to ensure that resource allocation for display output is properly managed, preventing race conditions or inconsistencies in the displayed content. By enabling more efficient coordination between the CPU, GPU, and display controller, this solution would allow for smooth and direct framebuffer access by both rendering components.

5.2 Reimagining Shader Programming with Scope-Based Shaders

A significant conceptual challenge in this unified architecture is the shift from traditional **GPU shaders** to **Scope-Based Shaders**—a new way of writing graphical programs that fundamentally differs from conventional shader programming. Traditional shaders are typically small, highly specialized programs that operate within fixed function pipelines (e.g., vertex shaders, fragment shaders) and are governed by constraints such as limited access to data (e.g., via uniforms or draw calls) and a rigid pipeline architecture.

Solution: In this new architecture, **Scope-Based Shaders** eliminate the need for rigid pipeline stages like vertex, fragment, and compute shaders. Instead, developers can write more **flexible, general-purpose code** that seamlessly blends **rendering, compute, and logic**. These shaders are not constrained by the traditional limitations of GPU shaders, such as limited data access or the need to manage complex, stage-based workflows.

To make this new paradigm accessible, new **graphics APIs and frameworks** would need to be developed to support this approach. These tools would provide abstractions for working with shared memory, allowing developers to focus on logic without having to worry about separate rendering and computation stages. Additionally, the programming model would need to be adapted to take advantage of the **scope-based execution model**, where rendering operations can directly manipulate data and logic in a more fluid, integrated manner. Developers would also need new debugging and performance profiling tools to accommodate this departure from the traditional shader pipeline.

Not only does this shift allow for simpler and more efficient programming, but **Scope-Based Shaders** can also **increase graphical complexity** and bring it to an entirely new level. By removing the traditional constraints of shaders, developers are free to create more intricate, dynamic, and detailed visual effects without the overhead of managing separate rendering stages. This could lead to more lifelike simulations, intricate lighting models, and detailed world-building in games and simulations. The ability to directly manipulate data and rendering simultaneously, within a single, unified context, opens up new possibilities for creating richer and more immersive graphical experiences.

While this shift in programming style is a significant departure from current GPU programming practices, it has the potential to **simplify** graphical and compute programming, enabling more efficient and flexible applications, especially in systems where tight integration between graphics and computation is required.

5.3 Resistance from GPU Vendors to Decoupling Display and Shaders

A potential challenge to the widespread adoption of this architecture is the potential **resistance from GPU vendors** to decouple **display handling** from traditional GPU operations. Currently, GPUs are designed with display output as a core function, and many GPU manufacturers have invested heavily in the development of specialized display subsystems and interfaces. The proposed architecture, which places the display controller within the SoC and moves away from GPU-centric rendering pipelines, may be seen as a threat to these established models.

Solution: Overcoming this resistance will require **industry-wide collaboration** to demonstrate the benefits of the unified architecture. One solution could be to present a **hybrid model** where the GPU still retains control over display functions in certain configurations, but the architecture allows for flexibility in scenarios where the display controller is handled by the SoC. This compromise would allow GPU vendors to maintain some of their display-focused capabilities, while still enabling the larger benefits of the unified architecture.

Additionally, **open standards** could be developed for shared memory and display access that would be flexible enough to accommodate both traditional and emerging hardware architectures. By building on existing GPU vendor frameworks and gradually integrating the new display model, the industry could move toward a unified, more flexible approach to CPU-GPU integration. As more developers adopt the new model and demonstrate its advantages in terms of performance, scalability, and developer simplicity, resistance may diminish over time.

Finally, a successful demonstration of the architecture's **market demand** and its ability to **simplify the development process** could convince vendors that this change is not a threat, but rather an opportunity to provide more versatile hardware platforms for the future of computing.

6. Use Cases

6.1 Minimalist Linux Systems (No GPU Needed)

One of the most compelling use cases for this unified CPU/GPU architecture is its application in **minimalist Linux systems** or low-power environments where traditional discrete GPUs are not necessary. In these systems, the **integrated display controller** within the SoC can handle all display output duties, removing the need for a dedicated GPU. This is particularly useful for **embedded systems**, **low-cost computing**, and **headless setups** where graphics rendering is minimal or not required at all.

With a unified memory architecture, even resource-constrained environments can access the same memory pool for both display and computing tasks, allowing for efficient resource use without the overhead of separate GPU memory. Developers can optimize their applications without worrying about the complexities of external GPUs, making it easier to deploy lightweight, efficient systems. Whether for kiosks, thin clients, or IoT devices, the unified architecture can significantly reduce both hardware requirements and power consumption while still providing reliable display and processing capabilities.

6.2 Game Engines with Logic-Aware Rendering

Game engines stand to benefit significantly from this architecture, particularly in the area of **logic-aware rendering**. Traditionally, rendering in game engines has been limited to the GPU's ability to execute pre-defined shaders and render pipelines. However, with **Scope-Based Shaders**, game engines can fully integrate **game logic with rendering and computation** in a seamless flow. This allows for a more dynamic and responsive rendering pipeline that can adjust to in-game logic in real time, offering a new level of interactivity.

For example, complex lighting models or world-building algorithms could adapt to in-game decisions or real-time events, enabling **dynamic rendering** based on player actions or environmental changes. Game engines could render complex scenes on-the-fly, directly manipulating game world data and graphical output without the need for constant draw calls or redundant processing. This new paradigm also simplifies the integration of real-time simulation and rendering, allowing for **immersive, interactive experiences** that blur the line between gameplay logic and visual presentation.

6.3 AI Workloads Blending Compute + Graphics

The unified CPU/GPU architecture is particularly well-suited for **AI workloads** that require intensive **compute** and **graphics** processing. Many AI tasks, such as **machine learning** and **computer vision**, benefit from the parallelism offered by GPUs, but they also require heavy computation that involves processing large datasets or training complex models. With the proposed architecture, both **compute tasks** and **graphics tasks** can seamlessly coexist on the same system, with the CPU and GPU sharing a unified memory pool and resources.

This enables **AI inference** and **training** workloads to directly interact with rendered data, such as images or videos, without needing separate memory management or additional data transfers. For instance, an AI system that processes video feeds for object recognition could render frames and simultaneously run inference algorithms on the same memory, without having to copy data between different components. This reduces overhead, improves performance, and streamlines workflows, making the system highly adaptable for AI applications in real-time, such as robotics, autonomous vehicles, and advanced visual recognition tasks.

6.4 Headless GPUs with Software or Remote Rendering

The **headless GPU** concept is gaining traction, particularly in **remote rendering** and **software-based rendering environments**. In traditional systems, GPUs are often tightly coupled with local displays. However, in many high-performance computing (HPC) setups or remote workstations, the need for a physical display is eliminated, and rendering is performed remotely or in a server environment. With this unified architecture, the concept of **headless GPUs** becomes more flexible and efficient.

By decoupling the GPU from the display output, the system can perform **software rendering** or **remote rendering** without requiring a physical monitor. The display controller within the SoC can be utilized for rendering to remote clients or via a network, making it ideal for **cloud gaming**, **virtual desktops**, and **distributed computing**. Systems could render complex graphical data on remote servers and stream the results to clients over the internet, while using the same unified memory for compute tasks. This could allow for more efficient use of **remote resources** and **virtualized environments**, where the display output is just one of many tasks handled by a distributed processing system.

6.5 Multi-GPU Rendering

A major advantage of the unified memory and framebuffer architecture is its ability to significantly enhance **multi-GPU rendering** workflows. Traditional multi-GPU setups face challenges such as fragmented memory spaces, inefficient data sharing, and synchronization issues between GPUs. These bottlenecks can limit the effectiveness of using multiple GPUs for rendering tasks.

With **unified memory** and a **shared framebuffer**, multiple GPUs can directly access the same memory and framebuffer, removing the need for complicated memory transfers and synchronization between GPUs. This creates a more seamless multi-GPU environment where GPUs can more efficiently work together on complex rendering tasks. Data can be shared and processed between GPUs without the overhead of duplicating data in multiple memory spaces, leading to higher performance and reduced latency.

Additionally, this unified approach simplifies the **management of multi-GPU configurations**, making it easier for both developers and end-users to scale their rendering setups. As a result, large-scale rendering tasks, such as those needed for real-time ray tracing or high-quality video rendering, can benefit from **faster rendering times, higher frame rates, and better scalability** with minimal system overhead. This is

particularly useful in industries such as **gaming**, **film production**, and **scientific visualization**, where multi-GPU setups are often essential for meeting the demands of high-performance rendering.

7. Conclusion: The Future is Unified

In this whitepaper, we have outlined a visionary architecture that combines the power of **CPU** and **GPU** into a single, unified system. By decoupling traditional boundaries between processors and memory, and integrating both **compute** and **rendering** capabilities into a cohesive unit, this approach paves the way for more **efficient**, **flexible**, and **powerful** computing systems. From **simplifying development** to **enhancing scalability**, this unified architecture represents a paradigm shift that could transform industries ranging from gaming to AI, and even embedded systems.

The **unified memory** model removes the need for complex data transfers between the CPU and GPU, allowing for more direct, efficient access to memory for all processing units. The **unified display controller** eliminates the need for a separate GPU for rendering, leading to simpler, more reliable systems. Meanwhile, the advent of **Scope-Based Shaders** offers developers new possibilities for integrating game logic, rendering, and computation in novel ways, without the traditional limitations of graphics pipelines.

As with any ambitious new architecture, **collaboration** across the industry is essential for its success. From hardware manufacturers to software developers and open-source communities, widespread adoption of this approach will require aligning on common standards, developing new programming models, and optimizing tools for integration. We envision a future where the **CPU**, **GPU**, and **display subsystems** are no longer isolated entities, but instead work together in harmony, unlocking new performance levels and simplifying development in a way that benefits all users.

The future is unified, and with the collaboration of the industry, we can build the next generation of computing systems that are more **power-efficient**, **developer-friendly**, and **capable** than ever before. This is a call to action for all stakeholders to come together, share knowledge, and work toward a unified architecture that can push the boundaries of what is possible in both **compute** and **graphics**.