

The current OS designs are primarily optimized for common hardware configurations. Modern systems have multiple processor cores with diverse architectural features such as memory hierarchies, interconnects, instruction sets, and IO configurations. It's no longer practical to optimize a general-purpose OS for a specific hardware model, as hardware configurations vary widely and become obsolete quickly but the dynamic nature of workloads and hardware variations makes it difficult to statically optimize an operating system for all scenarios.

The paper argues that the difficulties in OS engineering arise from the traditional structure of a shared-memory kernel with data structures protected by locks. The three design principles of Multikernel is:

1. Make all inter-core communication explicit.
2. Make the OS structure hardware-neutral.
3. View state as replicated instead of shared.

Difference in cores:

Cores within a single machine can differ and vary in terms of their performance characteristics. This diversity is a trend, with some cores optimized for parallelized programs and others for sequential tasks. Modern machines are trending towards using a mix of different cores. Core heterogeneity means that cores cannot share a single OS kernel instance this could be due to varying performance tradeoffs or differences in ISAs. A processor with large cores may be inefficient for parallelized programs, while one with only small cores may perform poorly on sequential tasks. Additionally, there may be cores with different ISAs for specialized functions.

Message sharing advantages:

For shared-memory systems, direct updates on a small set of memory locations are performed by threads pinned to each core, with cache coherence managing data migration between caches. Latency increases linearly with the number of threads and modified cache lines. As more cores update the same data, the update process becomes significantly slower due to cache misses, resulting in less productive work. In contrast, message-passing involves client threads issuing lightweight remote procedure calls to a server process. This incurs relatively constant costs regardless of the number of modified cache lines. While queuing delays may occur at the server, for updates involving multiple cache lines, the RPC latency is still lower than shared memory access. Additionally, with asynchronous or pipelined RPC, client processors can continue working without stalling on cache misses. Latency refers to the amount of time it takes for a remote procedure call (RPC) to be initiated, transmitted over a network or inter-process communication channel, processed by a remote server, and for the result (if any) to be returned to the caller. Lower RPC latency indicates that the communication between the client and server is more efficient and faster. The inherent lack of scalability in the shared memory model, combined with ongoing hardware innovation, is anticipated to pose increasingly challenging software engineering problems for OS kernels. With an increase in the number of processor cores and the intricacy of interconnections, maintaining hardware cache coherence becomes more resource-intensive. So, using message-based communication enables greater throughput, reduces interconnect utilization, and naturally accommodates heterogeneous hardware.

Concern with message sharing:

inability to access shared data in a message-based system. Message parsing system can also lead to complex control flow known as stack ripping.

In a multikernel OS, all communication between cores is conducted using explicit messages. No memory is shared between cores, except for the messaging channels. Explicit communication patterns provide transparency regarding which parts of shared state are accessed and by whom. Treating the machine as a network enables the OS to utilize well-known networking optimizations, such as pipelining and batching, to make more efficient use of the interconnect. Message passing allows for split-phase operations, where a request is sent and processing continues while awaiting a reply. This can lead to more efficient resource utilization, including power-saving strategies, as the requesting core can perform other tasks while waiting. A system built on explicit communication exhibits a naturally modular structure. This makes it easier to evolve, refine, and ensure robustness against faults.

In a multikernel system, only two aspects are tailored to specific machine architectures: the messaging transport mechanisms and the hardware interface (CPUs and devices). This reduces the need for extensive, cross-cutting code changes when adapting the OS to hardware with different performance characteristics. A multikernel system can support various messaging implementations, such as a user-level RPC protocol using shared memory (This means that processes (or components) in the system can communicate with each other by invoking procedures or functions that reside in separate address spaces, potentially even on different processors. The shared memory facilitates this communication, allowing processes to exchange data efficiently.) or a hardware-based channel to a programmable peripheral (This implies that the system has the capability to use specialized hardware channels to communicate with external devices) This flexibility is crucial, especially as some hardware platforms lack cache coherence or shared memory.

Explicit communication between cores naturally leads to a model where the OS state is replicated across cores. This replication approach is different from the traditional shared-memory kernel design. Replicating data structures can improve system scalability by reducing interconnect load, contention for memory, and synchronization overhead. It brings data closer to the processing cores, reducing access latencies. Making replication integral to the multikernel design helps in preserving OS structure and algorithms as underlying hardware evolves. It also provides a framework for supporting changes in the set of running cores, whether through hotplugging processors or shutting down hardware subsystems for power savings.

Here we are studying the Barrelfish implementation of multikernel system. In the case of Barrelfish, this structure is divided into two main components: CPU drivers and Monitors.

CPU drivers:

Since the CPU driver doesn't share state with other cores, it operates in an event-driven, single-threaded, and non-preemptable manner. It sequentially handles events (Easier to develop and debug due to this). Its compact size also allows it to reside in core local memory. The CPU driver efficiently manages local messaging between processes on the core via a lightweight, asynchronous interprocess communication mechanism within the same core. It is responsible for delivering hardware interrupts to user-space drivers and locally time-slicing user-space processes. The CPU driver is invoked through standard system call instructions.

Monitor:

They are user-space processes that operate on a single core, making them schedulable. Monitors are adept at managing message queues and overseeing lengthy remote operations. On each core, shared data structures like memory allocation tables and address space mappings are maintained in a globally consistent state through an agreement protocol managed by the monitors. When application requests involve global state, monitors act as intermediaries, facilitating access to remote replicas of the state. They also handle the setup of interprocess communication and wake up blocked local processes in response to messages from other cores. Furthermore, a monitor can put the core into an idle state to conserve power if no other processes on that core are ready to run. Core sleep can be achieved by waiting for an inter-processor interrupt (IPI) or, if supported, by using the MONITOR and MWAIT instructions.

Additionally, the remainder of Barrelfish includes:

Device Drivers and System Services:

These components, similar to a microkernel, operate as user-level processes. They encompass various functionalities like network stacks, memory allocators, and more.

Interrupt Handling:

Device interrupts are directed by hardware to the appropriate core. The CPU driver of that core manages the demultiplexing of interrupts and forwards them to the relevant driver process using messages.

A process in Barrelfish is represented by a set of dispatcher objects (The dispatcher in Barrelfish serves as an abstraction layer for scheduling and executing tasks or threads. It provides an interface between the operating system and the hardware, allowing the operating system to manage and allocate resources effectively.), with one dispatcher for each core where it might execute. Communication in Barrelfish doesn't occur directly between processes, but rather between dispatchers (and consequently, between cores). Dispatchers on a core are managed by the local CPU driver, which triggers an upcall interface provided by each dispatcher.

Barrelfish currently utilizes a variant of user-level RPC (URPC) between cores. This involves using a region of shared memory as a channel for transferring cache-line-sized messages directly between a single writer core and a single reader core. As an optimization, the throughput of pipelined URPC messages can be enhanced at the cost of single-message latency by employing cache prefetching instructions. This can be selected at the time of channel setup for workloads likely to benefit from it. Receiving URPC messages is achieved through memory polling, which is efficient because the line remains in the cache until invalidated. A dispatcher awaiting messages on URPC channels will poll those channels for a short period before blocking and sending a request to its local monitor to be notified when messages arrive. (as it will not wait indefinitely).

A name service is utilized to locate other services in the system, mapping service names and properties to a service reference which can be used to establish a channel to the service. Channel setup is conducted by the monitors.

In a multikernel OS, managing global resources like physical memory consistently across cores is crucial. Barrelfish employs a capability system. In this approach, all memory management operations are explicitly performed through system calls that manipulate capabilities (Barrelfish follows the principle of least privilege, which means that a program or component is granted only the minimum level of access or permissions necessary to perform its function). Barrelfish employs an object-capability model, where a capability is an unforgeable token that grants its possessor specific access rights to a resource (such as a memory region, a network interface, or other system components). Capabilities are used as a basis for resource management and access control. These capabilities act as user-level references to kernel objects or regions of physical memory. (Notably, dynamic memory allocation is removed from the CPU driver's responsibilities, which is solely tasked with verifying the correctness of operations that manipulate capabilities related to memory regions, it is handled by user level code.) A benefit of the capability-based approach is its uniformity. Many operations in Barrelfish that require global coordination (meaning they involve multiple components or entities that need to synchronize their actions) can be framed as instances of capability copying or retyping, allowing for the use of generic consistency mechanisms in the monitors (For example, if two components need to coordinate access to a shared resource, they may do so by copying or retyping the necessary capabilities. This allows them to establish a consistent view of the resource and its state). These operations are not tied exclusively to capabilities, and would need to be supported regardless of the chosen accounting scheme. However, in hindsight, the use of capabilities for memory management in Barrelfish is considered overly complex and not significantly more efficient than the per-processor memory managers found in conventional OSes. Capability retyping, including revocation, presents a more intricate challenge. This involves changing the usage of a memory area and requires global coordination, as reassigning the same capability in different ways on different cores can lead to an inconsistent system. To address this, all cores must agree on a single ordering of the operations to maintain safety. In this scenario, the monitors initiate a two-phase commit protocol to guarantee that all changes to memory usage are consistently ordered across the processors. The monitors in Barrelfish provide a mechanism for sending capabilities between cores, ensuring that the capability is not pending revocation and is of a type that can be transferred. User-level libraries responsible for capability manipulation invoke the monitor as needed to maintain a consistent capability space across cores. Cross-core thread management is also handled in user space. The thread schedulers on each dispatcher communicate via messages to create and unblock threads, as well as to migrate threads between dispatchers (and therefore cores). Barrelfish is primarily responsible for multiplexing the dispatchers on each core through the CPU driver scheduler. It also coordinates the CPU drivers to perform tasks like gang scheduling or co-scheduling of dispatchers. This flexibility allows for the application of various spatio-temporal scheduling policies based on OS policy requirements.

Barrelfish supports the traditional process model where threads share a single virtual address space across multiple dispatchers (and hence cores). This is achieved through coordination of runtime libraries on each dispatcher. This coordination involves three key OS components: virtual address space, capabilities, and thread management. It serves as an example of how traditional OS functionality can be provided over a multikernel architecture.

In Barrelfish, addressing the challenges of heterogeneous hardware and making informed decisions about system mechanisms is crucial. To achieve this, Barrelfish employs a component known as the System Knowledge Base (SKB). The SKB maintains knowledge about the underlying hardware in a subset of first-order logic., the SKB enables concise expression of optimization queries. For instance, it can be used to allocate device drivers to cores in a topology-aware manner, select appropriate

message transports for inter-core communication, and provide the necessary topology information for NUMA-aware (Non-Uniform Memory Access) memory allocation.