

The SawMill Multiserver Approach

Alain Gefflaut Trent Jaeger Yoonho Park
Jochen Liedtke* Kevin Elphinstone* Volkmar Uhlig*
Jonathon E. Tidswell Luke Deller† Lars Reuther‡

*IBM T. J. Watson Research Center
Hawthorne, NY 10532*

Email: {sawmill@watson.ibm.com}

1 Introduction

Multiserver systems, operating systems composed from a set of hardware-protected servers, initially generated significant interest in the early 1990's. If a monolithic operating system could be decomposed into a set of servers with well-defined interfaces and well-understood protection mechanisms, then the robustness and configurability of operating systems could be improved significantly. However, initial multiserver systems [4, 14] were hampered by poor performance and software engineering complexity. The Mach microkernel [2] base suffered from a number of performance problems (e.g., IPC), and a number of difficult problems must be solved to enable the construction of a system from orthogonal servers (e.g., unified buffer management, coherent security, flexible server interface design, etc.).

In the meantime, a number of important research results have been generated that lead us to believe that a re-evaluation of multiserver system architectures is necessary. First, microkernel technology has vastly improved since Mach. L4 [13] and Exokernel [8] are two recent microkernels upon which efficient servers have been constructed (i.e., L4Linux for L4 [11] and ExOS for Exokernel [12]). In these systems, the servers are independent OSes, but we are encouraged that the kernel and server overheads, in particular context switches overheads, are minimized. Second, we have seen marked improvements in memory management approaches that enable zero-copy protocols (e.g., fbufs [6] and emulated copy [3]). Other advances include, improved kernel modularity [9], component model services [7], multiserver security protocols, etc. Note that we are not the only researchers who believe it is time to re-examine multiservers, as a multiserver system is also being constructed on the Pebble kernel [10].

In addition, there is a greater need for multiserver architectures now. Consider the emergence of a variety of specialized, embedded systems. Traditionally, each embed-

ded system includes a specialized operating system. Given the expected proliferation of such systems, the number of operating systems that must be built will increase significantly. Tools for configuring operating systems from existing servers will become increasingly more valuable, and adequate protection among servers will be necessary to guard valuable information that may be stored on such systems (e.g., private keys). This is exactly the motivation for multiservers.

In this paper, we define the SawMill multiserver approach. This approach consists of: (1) an architecture upon which efficient and robust multiserver systems can be constructed and (2) a set of protocol design guidelines for solving key multiserver problems. First, the SawMill architecture consists of a set of user-level servers executing on the L4 microkernel and a set of services that enable these servers to obtain and manage resources locally. Second, the SawMill protocol design guidelines enable system designers to minimize the communication overheads introduced by protection boundaries between servers. We demonstrate the SawMill approach for two server systems derived from the Linux code base: (1) an Ext2 file system and (2) an IP network system.

The remainder of the paper is structured as follows. In Section 2, we define the problems that must be solved in converting a monolithic operating system into a multiserver operating system. In Section 3, we define the SawMill architecture, which defines the types of components from which a multiserver is to be constructed. In Section 4, we outline guidelines for designing protocols that satisfy the multiserver requirements while minimizing communication overheads. In Section 5, we demonstrate some of these guidelines in the file system and network system implementations.

2 Multiserver Issues

An effective multiserver system must: (1) protect its servers from errors or malice in other servers; (2) implement coherent system semantics; and (3) incur minimal performance overhead. We define these requirements more precisely in this section.

*Faculty of Informatics, Universität Karlsruhe, Germany

†School of Computer Science and Engineering, University of New South Wales, Australia

‡Department of Computer Science, Technische Universität Dresden, Germany

2.1 Protection

Multiserver systems must preserve the integrity of each server's execution (i.e., protect servers from one another) and protect the integrity and secrecy of the data processed by each server (i.e., ensure the protection of user data). Specifically, we list the following protection requirements:

- Protect the execution integrity of each server:
 - Prevent modification of another server's code
 - Prevent modification of another server's 'control data' (i.e., data that is interpreted for execution, such as the stack)
 - All code and control data must only be obtained from a trusted source
- Protect the secrecy and integrity of user data:
 - Prevent leakage of data to unauthorized subjects
 - Prevent modification of data by unauthorized subjects
 - Protect data from accidental modification by other servers

The first set of requirements prevents malicious or buggy servers from crashing other servers (i.e., protects server integrity). A malicious or buggy server can only cause another server to crash if it modifies the data interpreted by the server in its execution: code, stack, and 'control' data (i.e., memory references and meta-data). We distinguish between *control data*, the data interpreted by the server during execution, and *user data*, the data being transferred from the users to the devices or vice versa via the operating system. The server must both obtain its code, stack, and control data from a trusted source and limit modification of this data to verifiably correct protocols (e.g., Clark-Wilson's well-formed transactions for transformation procedures [5]).

The second set of requirements are designed to protect the user data being processed by the server. In general, there are two types of requirements embodied here. The first two requirements embody traditional access control, including control over overt and covert channels, in general. The third requirement implies that user data transferred among servers must be protected effectively from accidental server modification. That is, as user data is being transferred from one server to another, there must be some protocol to prevent the first server from modifying this user data.

2.2 Semantics

In general, the problem of decomposing a monolithic operating system into a multiserver operating system is shown in Figure 1. In a multiserver system, the monolithic system

is implemented by a set of user-level system tasks upon a microkernel (not shown). The combination of these tasks define the semantics seen by the user tasks. We define how a multiserver must preserve these semantics below.

- Each system call in the monolithic system must be supported by one or more servers in the multiserver system ¹.
- Each system call must be processed as restricted by its atomicity requirements.
- Any server must be able to obtain and enforce system policies (e.g., for security, resource management, etc.).

First, the decomposed system must be able to respond to all the same system calls as the monolithic system. Thus, the multiserver system has the same functionality as the monolithic system. Second, the atomicity requirements of system calls in the monolithic system must be enforced in the multiserver system. If system call data is distributed among multiple servers, concurrency control mechanisms must account for all those servers to prevent race conditions. Third, a coherent system composed of a set of servers requires that any server be able to obtain and enforce system-wide policies, such as access control and resource management.

2.3 Performance

The goal of the SawMill multiserver design is to achieve the protection and semantic requirements with no significant performance degradation. Maintaining performance in a multiserver system is non-trivial. Consider Figure 1, in which a single-threaded, monolithic operating system is broken into two single-threaded system tasks. The following additional operations may result due to this decomposition:

- **IPC Frequency:**
 - IPCs (i.e., context switches) replace procedure calls
 - Additional IPCs may be necessary to maintain the consistency of replicated data
 - Additional IPCs may be necessary to synchronize access to shared data
 - Additional IPCs may be necessary to negotiate resource allocation
 - Cache and TLB locality are reduced by additional context switches

- **IPC Overhead:**

¹In an embedded system, only a subset of the system calls may be necessary, so only these must be supported.

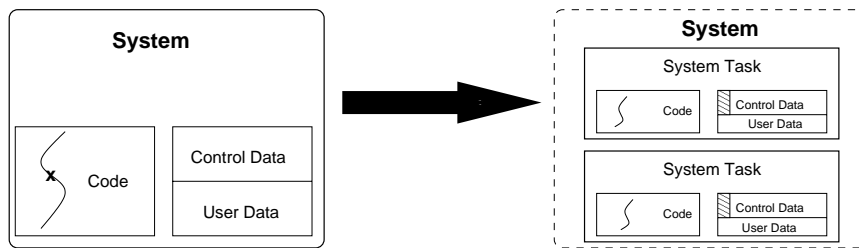


Figure 1: Decomposing a monolithic operating system involves distributing system functionality and meta-data into protection domains. Introducing protection domain boundaries may introduce communication, copying, mapping, and consistency management overheads which must be mitigated.

- Parameters must be marshalled and unmarshalled between servers
- Parameters need to be transferred between the servers

First, each procedure call that crosses a server boundary in now convert to an IPC, which consists of context switching, marshalling, unmarshalling, copying, and mapping overheads. Clearly, the number of IPCs should be minimized, and the overhead of each IPC and related functionality must be minimized. Also, depending on the protocol chosen, additional IPCs may be necessary to maintain data consistency, synchronize access, and obtain resources. Second, data may need to be transferred across these protection boundaries. While control data transfer can be limited to effective partitioning and caching, significant amounts of user data will be transferred through the servers. This data may be copied, mapped, or shared depending on the protection and protocol semantics.

3 SawMill Architecture

The SawMill architecture for multiserver operating systems is shown in Figure 2. The SawMill architecture consists of three types of components:

- **System servers:** These tasks provide the main functionality of the operating system (e.g., network systems, file systems, etc.)
- **Ubiquitous Services:** These components provide general functionality that may be of use to any system server (e.g., synchronization, access control, naming, communication, etc.).
- **Resource Servers:** These servers manage the core resources for distribution to the system servers (e.g., memory, IRQs, security policy, etc.).

System server code is augmented by libraries, called *ubiquitous services*, that provide multiserver-aware management of system data. In general, a ubiquitous service obtains resources from the appropriate *resource servers*

which control the distribution of resources among servers. Ubiquitous services manage these resources locally (e.g., caches) to limit communication overhead.

The basic protocol is as follows for memory management [1]. First, the server requests that its memory service library obtain access to its *dataspaces* (i.e., memory objects). The memory service library then opens the appropriate dataspaces (i.e., similar to memory objects) on the appropriate memory servers (i.e., resource server). When the server attempts to access this memory, the page fault is directed to memory service library which requests that appropriate memory server service the page fault. Notice that memory servers may be stacked, such that a memory server may need to obtain the page from its memory server, and so on. Similar *resource faulting* approaches are also used for obtaining access control data, names, tasks, and mount points. Other ubiquitous libraries provide multiserver-specific functionality, such as Flick’s cross-domain procedure call stubs [7] and synchronization.

4 SawMill Protocols

In order to build an efficient multiserver, system protocols must designed to minimize IPC frequency and overhead (see Section 2.3). We identify these principles for designing efficient multiserver protocols:

- Make system calls directly to the processing server
- Partition server-specific control data
- Share data as widely as possible

First, to reduce the number of IPCs that are necessary to implement a performance-critical system call, clients should communicate directly to the servers that process the system call, where possible. Second, server control data should be partitioned where possible. In cases, where the data is shared (i.e., cached read-only in multiple servers), few update messages should be necessary (i.e., few writes or weak consistency). Third, data sharing should be utilized as widely as protection requirements can allow. In

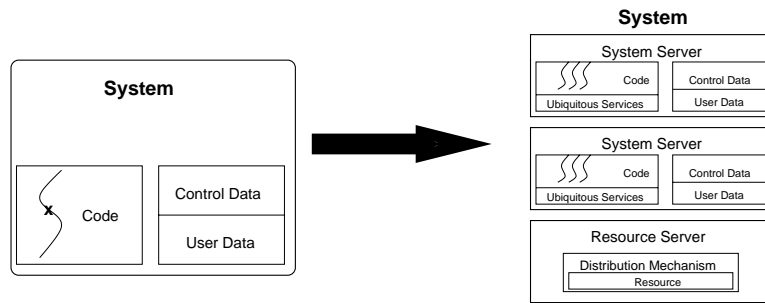


Figure 2: The SawMill approach enables the decomposition of monolithic operating systems into individual *system servers* where *ubiquitous services* manage system data locally in a multiserver environment. The ubiquitous services obtain the system data from *resource servers* which also control the distribution of such resources. Thus, a coherent system can be constructed from disjoint servers.

particular, the mapping of user data can be avoided in system servers, because the same buffers can be reused for the same user. Also, it may be desirable to share control data read/write between drivers and servers, as long as each can prevent themselves from crashing due to errors or malice.

5 Some Implementation Details

We now examine some interesting examples of applying these concepts to the implementation of the file system and network system. This implementation is based on Linux 2.2.1 code and preserves Linux semantics.

5.1 Protocol Options

First, each system protocol should be designed to require as few IPCs as possible. Consider the protocol options for a multiserver file system in Figure 3. The first protocol simply replaces the inter-server procedure calls with IPCs. In this case, the VFS is called on `open` to obtain a file handle. On subsequent `read/write` calls, the client calls the VFS which authorizes the handle and forwards authorized requests onto the PFS. In the second protocol, the VFS obtains a handle from the PFS that the client can use on the PFS directly. Therefore, only the `open` call is sent to the VFS, but there are potentially several IPCs between the VFS and PFS for name resolution. The third protocol option considers the VFS only as a store for mount point and access control information. This is typical of the resource faulting protocols described in the previous section. In this case, a client emulation library obtains the PFS-mount point mapping, so it only needs to call the VFS for each new mount point it uses². Also, the PFS caches access control data, so it only needs to call the VFS when a request does not hit in the access control data cache.

²An exception is on the processing of symbolic links between file systems, but these can be handled by having each PFS provide the client with the PFS to which its part of the path resolves.

In the normal case, the third option will result in fewer IPCs, but, due to the complexity in modifications necessary to implement it, the second option is being used currently. Overhead for read/write is the same for both the second and third case, however (see Section 6).

For the network system, an option analogous to the third option is implemented. In this option, a network manager performs the role of the VFS in naming stacks and devices and providing access control information. The network stack corresponds to a PFS.

5.2 Data Partitioning

Fortunately, a significant amount of control data can be partitioned between the servers in both the file system and the network system. However, user data must flow (not necessarily be copied) through the servers to reach its destination, so it cannot be partitioned. In the file system, the main control data are the superblocks and the inodes. Superblocks change very infrequently, so caching is not a problem. Inodes are provided by the driver and used read/write by both the PFS and the VFS. However, the inode data used and updated by the VFS and PFS are orthogonal. So, the *master copy* is stored by the PFS, and the VFS caches its copy and sends updates to the PFS. The PFS does not send updates to the VFS. If the third protocol option were used, little or no sharing of inode data would be necessary between the VFS and PFS.

A problem is that the buffer cache is shared among the PFSs in Linux. Thus, one PFS could disrupt the inode data of another PFS. Clearly, in order to prevent this, a trusted entity (either a buffer cache server or a driver) must partition the buffer cache among PFSs. Since the buffer cache usage of the different PFSs may vary over time, this server will need significant resource management abilities. Currently, our decomposed PFS and driver share their own separate buffer cache, but this is still an open problem.

The network system control data includes device structures, buffer lists, and `sk_buffs`. Since the device structure is used read-only by the network stacks, the drivers

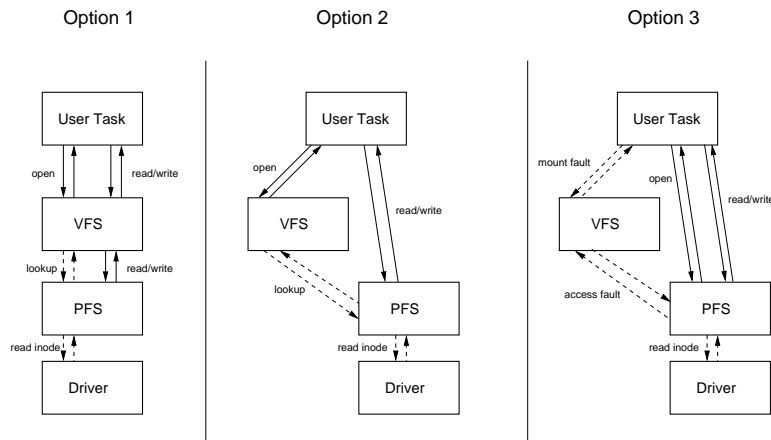


Figure 3: The communication requirements for three file system protocol options. Dashed IPCs are optional. For example, a PFS lookup may not be necessary because the appropriate inodes may be present in the directory cache.

maintain the master copy. The network stacks get a copy when they open the device, and the driver notifies them of any changes. There are function pointers in the structure, so these pointer are localized at load time. The buffer lists are partitioned between the driver and the stack. For example, when a stack wishes to send a packet, it sends the driver a reference to an `sk_buff` which the driver can queue on its own list. There is no sharing of buffer lists between the stack and driver. `sk_buffs` and packet data are shared as discussed below.

5.3 I/O Management

There are two important issues in I/O management: (1) the protection boundaries that the I/O data crosses (user-system or server-server within the system) and (2) the type of the data being transferred (user or control data). First, in a traditional UNIX system, such as Linux, data is typically *copied* between the user and the system. In a multi-server system, we have this type of data transfer and data transfer between servers in the system. System servers originally *shared* this data in the monolithic system, so sharing semantics are more natural. Second, we must consider the protection requirements of the type of data that is being transferred. Protection requirements permit user data to be shared among multiple servers, as long as sufficient copy semantics are enforced between the user and system (some form of *copy-on-write*), but servers must protect themselves from the generation of illegal control data.

First, user data transferred between the user and system is transferred with copy semantics. However, user data buffers are shared among the servers. In effect, the system servers are *originators* (in the `fbufs` sense [6]) of the data, so they maintain write access until the data is transferred to the user or to an untrusted server. Then, a decision must be made whether to downgrade the servers' access to the data to read-only (e.g., using *emulated copy* [3] or

fbuf semantics) and apply copy-on-write for the user, permit the servers to maintain read-write access (e.g., volatile buffers) with copy-on-write for the user, or copy on transfer. Unfortunately, neither emulated copy or `fbuf` semantics work for the Linux file system. Buffer cache buffers may not be aligned as necessary for emulated copy and Linux copy semantics does not enable the use of `fbufs`. At present, data is copied between user tasks and the SawMill Linux system servers, although further investigation is ongoing.

For control data, each server must trust the originator of the control data it uses to specify the necessary data, but it must ensure that the control data that it uses does not cause it to crash. For example, the network stack specifies `sk_buffs` for the network driver which the network driver must assume refer to the data that it is to send (if it is within the shared user data region). However, if the `sk_buff` is erroneous (i.e., refers to an illegal memory location), the network driver must catch this. In this case, we prefer an optimistic approach where the network stack and driver share `sk_buff` data and the driver uses the `sk_buff` directly. Any error in the `sk_buff` either results in the wrong data being sent (which the network stack could arrange regardless) or results in an illegal page fault. Since the driver's page fault handler is a local thread, it can recover the driver thread to a consistent state on such an illegal page fault.

6 Performance Summary

Initial performance results on the file system are shown in Figure 4. We compare SawMill Linux 2.2.1, L4Linux 2.2.1 (trampoline version), and Linux 2.2.1 systems using the IOZone read benchmark. All the data is from a 266 MHz, Pentium II with a 512KB L2 cache and a 32KB split L1 cache. The nature of the benchmark is that all the read operations hit in the file cache, so the time shown consist

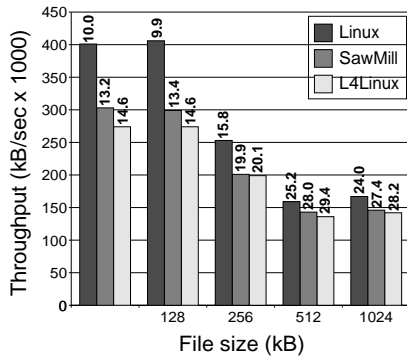


Figure 4: Initial performance comparison between monolithic Linux, SawMill Linux, and L4Linux (trampoline version)

of a round-trip communication and a copy.

For each 4KB record read, the SawMill transaction times are 3.2 μ s slower than monolithic Linux. Per record read, this amounts to about 800 cycles difference. Overall, for 64KB reads, the SawMill throughput is 305,000 KB/s while monolithic Linux throughput is 401,000 KB/s. Since the same copy routines are used in all cases, the difference is attributable to communication overhead differences. The performance of the L4Linux trampoline version is as expected: it should be slower than SawMill Linux because of the extra system call to implement the trampoline.

However, we know that SawMill performance can be improved in two ways: (1) our Flick stubs can be optimized by writing the parameters directly into the registers which we estimate will save 200-300 cycles based on a hand-coded example and (2) we can use the Pentium II-specific fast `sysenter/sysexit` instructions which will save 100-150 instructions³. Thus, we expect that we can save another 1.25 to 1.75 μ s on the multiserver communication time. Also, modifying Flick to use L4's scatter-gather will reduce the number of communications necessary to transfer a file. Thus, the cost of communication can be amortized for larger reads. We are working on scatter-gather currently and plan to have numbers for larger record sizes by the publication deadline.

7 Conclusions

In this paper, we describe the SawMill approach to constructing multiserver operating systems. This approach

³Note that a round-trip IPC requires two system calls vs. one system call for a monolithic Linux read.

consists of a set of services and design guidelines that enable the development of a multiserver from an existing code base. We demonstrate this approach on the Linux operating system by building multiserver file and network systems. The main complexity lies in the efficient and secure management of data over the disjoint servers, although we have a number of potentially promising approaches. Also, initial performance results indicate that the multiserver overhead is small on the file system and can be further optimized to reduce the impact of multiserver protection barriers. Soon, we will have the network system performance measurements.

References

- [1] M. Aron, J. Liedtke, Y. Park, K. Elphinstone, and T. Jaeger. Implementing and using VM diversity. Unpublished report available at URL <http://www.research.ibm.com/sawmill/>.
- [2] R. V. Baron *et al.* *Mach Kernel Interface Manual*, 1990. Unpublished manuscript from the School of Computer Science, Carnegie Mellon University.
- [3] J. C. Brustoloni and P. Steenkiste. Effects of buffering semantics on I/O performance. In *Proc. OSDI'96*, 1996.
- [4] T. Bushnell. Towards a new strategy for OS design, 1996. At URL <http://www.gnu.ai.mit.edu/software/hurd>.
- [5] D. D. Clark and D. R. Wilson. A comparison of commercial and military computer security policies. In *Proc. of IEEE Symp. on Security and Privacy*, 1987.
- [6] P. Druschel and L. Peterson. Fbufs: A high-bandwidth cross-domain transfer facility. In *Proc. SOSP'93*, 1993.
- [7] E. Eide, K. Frei, B. Ford, J. Lepreau, and G. Lindstrom. Flick: A flexible, optimizing IDL compiler. In *Proc. PLDI'97*, 1997.
- [8] D. Engler, F. Kaashoek, and J. O'Toole. Exokernel: An operating system architecture for application level resource management. In *Proc. SOSP'95*, 1995.
- [9] B. Ford *et al.* The Flux OSKit: A substrate for kernel and language research. In *Proc. SOSP'97*, 1997.
- [10] E. Gabber, C. Small, J. Bruno, J. Brustoloni, and A. Silberschatz. Building efficient operating systems from user-level components in Pebble. In *Proc. USENIX'99*, 1999.
- [11] H. Härtig, M. Hohmuth, J. Liedtke, S. Schönberg, and J. Wolter. The performance of μ -kernel-based systems. In *Proc. SOSP'97*, 1997.
- [12] M. F. Kaashoek *et al.* Application performance and flexibility on exokernel systems. In *Proc. SOSP'97*, 1997.
- [13] J. Liedtke. On μ -kernel construction. In *Proc. SOSP'95*, 1995.
- [14] J. M. Stevenson and D. P. Julin. Mach-US: UNIX on generic OS object servers. In *Proc. USENIX'95*, 1995.