

Building a Bare-Metal Operating System in Lateralus

bad-antics | March 2024 | Technical Report

Abstract

This paper documents the construction of a bare-metal operating system written in Lateralus, targeting x86-64 hardware. The OS implements process management, virtual memory, interrupt handling, device drivers, a filesystem, and a user-space shell, demonstrating Lateralus's capability for low-level systems programming.

1 Introduction

Building a bare-metal operating system is one of the most challenging and rewarding projects in systems programming. This paper documents the construction of a minimal operating system written primarily in Lateralus with critical bootstrapping components in assembly. The OS targets x86-64 hardware and implements process management, memory management, device drivers, and a simple filesystem.

The goal is not to compete with production operating systems but to demonstrate that Lateralus's low-level capabilities, pipeline-native design, and ownership model are well-suited for systems programming at the hardware level. Every abstraction from interrupt handling to memory allocation is implemented from scratch.

This paper covers the boot process, protected mode transition, memory management, interrupt handling, device drivers, process scheduling, filesystem implementation, and the shell interface. Each section includes design rationale, implementation details, and the Lateralus-specific techniques that simplify OS development.

2 Boot Process

The boot process begins with UEFI firmware loading a Lateralus bootloader from the EFI System Partition. The bootloader uses UEFI Boot Services to query the memory map, locate the kernel image, and configure the framebuffer. Once the kernel is loaded into memory, the bootloader exits Boot Services and transfers control to the kernel entry point.

The bootloader is compiled as a UEFI application using Lateralus's freestanding mode, which disables standard library dependencies. UEFI protocols are accessed through a foreign function interface that maps UEFI calling conventions to Lateralus function signatures. Error handling uses Lateralus's Result type to propagate UEFI status codes.

```
// UEFI bootloader entry point in Lateralus
@entry fn efi_main(handle: EfiHandle, table: *SystemTable) -> EfiStatus {
    let bs = table.boot_services;
    let gop = bs.locate_protocol::()?;
```

```
gop.set_mode(preferred_mode(gop))?;  
let kernel = load_file(bs, "\\kernel.elf"?);  
let mmap = bs.get_memory_map()?;  
bs.exit_boot_services(handle, mmap.key)?;  
kernel.entry(mmap, gop.framebuffer)  
}
```

After exiting Boot Services, the kernel receives the memory map and framebuffer information from the bootloader. The kernel performs early initialization in a small assembly stub that sets up the stack, enables the required CPU features, and jumps to the Lateralus kernel main function.

3 Protected Mode and Paging

The kernel initializes the Global Descriptor Table with segments for kernel code, kernel data, user code, and user data. The GDT is loaded using the lgdt instruction through an inline assembly block. Segment selectors are loaded into the appropriate registers to establish the flat memory model used by long mode.

Page tables are constructed using a four-level hierarchy: PML4, PDPT, PD, and PT. The kernel identity-maps the first 4 GB of physical memory and creates a higher-half mapping at 0xFFFF800000000000 for the kernel address space. Page table entries use Lateralus's bitfield types for readable flag manipulation.

```
// Page table entry using Lateralus bitfields  
type PageEntry = bitfield u64 {  
    present: 1,  
    writable: 1,  
    user: 1,  
    write_through: 1,  
    cache_disable: 1,  
    accessed: 1,  
    dirty: 1,  
    huge_page: 1,  
    global: 1,  
    _reserved: 3,  
    frame: 40,  
    _reserved2: 11,  
    no_execute: 1,  
}
```

The Translation Lookaside Buffer is flushed after page table modifications using `invlpg` for single-page invalidations and a full CR3 reload for bulk changes. The kernel tracks TLB state per-CPU to minimize unnecessary flushes on multiprocessor systems.

4 Physical Memory Management

Physical memory is managed using a buddy allocator that handles allocation and deallocation of

page frames. The allocator organizes free memory into power-of-two sized blocks from 4 KB (single page) to 2 MB (huge page). Splitting and coalescing operations maintain the buddy invariant, preventing external fragmentation.

The buddy allocator is initialized from the UEFI memory map, which describes available, reserved, and firmware-occupied memory regions. Only regions marked as conventional memory are added to the allocator. The initialization process coalesces adjacent free regions before inserting them into the buddy tree.

Allocation uses a bitmap to track the status of each block at each order level. The bitmap is stored in a reserved memory region and is sized proportionally to total physical memory. For a system with 4 GB of RAM, the bitmap requires approximately 128 KB of storage.

```
// Buddy allocator interface
pub fn alloc_frame() -> Result<PhysAddr, AllocError> {
    alloc_order(0) // single 4KB page
}

pub fn alloc_order(order: u8) -> Result<PhysAddr, AllocError> {
    if let Some(block) = free_list[order].pop() {
        return Ok(block);
    }
    // Split a larger block
    let parent = alloc_order(order + 1)?;
    let buddy = parent + (PAGE_SIZE << order);
    free_list[order].push(buddy);
    Ok(parent)
}
```

5 Virtual Memory and Heap

The kernel heap uses a slab allocator built on top of the buddy allocator. Slab caches are created for commonly allocated sizes (32, 64, 128, 256, 512, 1024, 2048 bytes) and for specific kernel structures (process descriptors, file handles, network buffers). Each slab cache maintains a free list of pre-allocated objects.

Virtual memory areas are tracked using a red-black tree of VMA descriptors. Each VMA records the virtual address range, permissions, backing source (anonymous, file-mapped, device), and page fault handler. The VMA tree supports efficient lookup, insertion, and merging of adjacent regions with compatible attributes.

Demand paging defers physical memory allocation until a page is actually accessed. Page faults trigger the VMA's fault handler, which allocates a physical frame, initializes it, and installs the page table mapping. This lazy approach reduces memory consumption for sparse allocations and speeds up process creation.

Copy-on-write sharing allows forked processes to share page frames until one process modifies a

shared page. The write fault handler allocates a new frame, copies the page contents, and updates the faulting process's page table. Reference counting on page frames tracks the number of sharing processes.

6 Interrupt Handling

The Interrupt Descriptor Table maps interrupt vectors to handler functions. Each IDT entry specifies the handler address, privilege level, and stack switching behavior. The kernel uses a separate interrupt stack for each CPU to prevent stack overflow when interrupts nest.

Hardware interrupts are managed through the Advanced Programmable Interrupt Controller (APIC). The kernel configures the Local APIC for timer interrupts and inter-processor interrupts, and the I/O APIC for device interrupts. Interrupt routing is configured to distribute device interrupts across available CPUs.

```
// Interrupt handler registration
pub fn register_handler(vector: u8, handler: fn(&InterruptFrame)) {
    IDT[vector] = IdtEntry::new(
        handler as u64,
        KERNEL_CS,
        IdtFlags::PRESENT | IdtFlags::INTERRUPT_GATE,
    );
}

// Timer interrupt handler
fn timer_handler(frame: &InterruptFrame) {
    TICKS.fetch_add(1, Ordering::Relaxed);
    scheduler::reschedule();
    apic::eoi();
}
```

Exception handlers provide detailed diagnostics for CPU exceptions including page faults, general protection faults, and double faults. The page fault handler distinguishes between valid demand-paging faults and invalid accesses, passing valid faults to the memory manager and terminating processes that access invalid addresses.

System calls use the syscall instruction for fast user-to-kernel transitions. The syscall handler saves user-space registers, switches to the kernel stack, and dispatches to the appropriate system call implementation based on the syscall number in the RAX register.

7 Device Drivers

Device drivers follow a common interface defined by the Driver trait. Each driver implements probe (device detection), init (initialization), read, write, and ioctl operations. The driver model supports hot-plugging through a device tree that is updated when hardware changes are detected.

The keyboard driver handles PS/2 keyboard input through IRQ 1. Scan codes are translated to key events using a configurable keymap. The driver supports key repeat, modifier tracking, and a circular buffer for storing keystrokes until they are consumed by the active process.

The framebuffer driver provides pixel-level access to the display using the UEFI-configured framebuffer. A software text renderer draws characters using an embedded bitmap font. The driver supports scrolling, cursor positioning, and ANSI escape codes for terminal emulation.

The serial port driver provides UART communication for debugging and headless operation. The driver configures baud rate, parity, and flow control through the UART's register interface. Output is buffered and transmitted using interrupt-driven I/O to prevent busy-waiting.

The storage driver implements AHCI for SATA devices, providing block-level read and write operations. Direct Memory Access transfers data between disk and memory without CPU involvement. The driver supports command queuing for improved throughput on modern SATA devices.

8 Process Management

Processes are represented by a `Process` structure containing the process ID, page table root, register state, open file descriptors, and scheduling parameters. Process creation allocates a new address space, loads the executable image, and initializes the stack with command-line arguments and environment variables.

The scheduler uses a multi-level feedback queue with four priority levels. New processes start at the highest priority and are demoted when they exhaust their time quantum. I/O-bound processes that voluntarily yield are promoted, ensuring interactive processes maintain responsiveness.

```
// Process creation
pub fn spawn(path: &str, args: &[&str]) -> Result<Pid, SpawnError> {
    let proc = Process::new()?;
    let elf = fs::read(path)?;
    proc.load_elf(&elf)?;
    proc.setup_stack(args)?;
    proc.state = ProcessState::Ready;
    SCHEDULER.enqueue(proc.pid);
    Ok(proc.pid)
}
```

Context switching saves the current process's register state to its process structure and restores the next process's state. The switch also updates the page table root in CR3 and flushes the TLB. The context switch path is carefully optimized to minimize latency, as it directly impacts system responsiveness.

Inter-process communication uses message passing through kernel-managed channels. Each channel has a fixed-size buffer for asynchronous message delivery. Processes can block on channel receive, enabling synchronous communication patterns. The channel implementation uses lock-free

algorithms for minimal overhead.

9 Filesystem

The filesystem implements a simple Unix-like hierarchy with directories, regular files, and symbolic links. The on-disk format uses an inode-based structure with a superblock, inode table, data block bitmap, and data blocks. The filesystem supports files up to 4 GB using direct, indirect, and double-indirect block pointers.

The block cache maintains recently accessed disk blocks in memory. Cache eviction uses the clock algorithm to approximate LRU behavior with minimal overhead. Write-back caching defers disk writes until the cache is flushed, a block is evicted, or a sync operation is requested.

Directory entries are stored as variable-length records within directory data blocks. Each entry contains the file name, inode number, and entry type. Directory lookup uses linear search within blocks, which is adequate for the moderate directory sizes expected in an educational OS.

The Virtual Filesystem Switch layer provides a unified interface above physical filesystems. VFS operations are dispatched to the appropriate filesystem driver based on mount point. This abstraction allows the kernel to support multiple filesystem types simultaneously.

10 Shell and User Space

The shell provides an interactive command-line interface for system interaction. Built-in commands include `cd`, `ls`, `cat`, `echo`, `mkdir`, `rm`, `ps`, and `kill`. External commands are loaded from the filesystem and executed as child processes. The shell supports command pipelines using the pipe system call.

User-space programs are compiled as ELF executables using Lateralus's freestanding standard library. The library provides system call wrappers, memory allocation (via `brk` and `mmap` system calls), string manipulation, and formatted output. The library is statically linked into each executable.

The `init` process is the first user-space process launched by the kernel after boot. `init` reads a configuration file specifying services to start and their dependencies. Services are started in dependency order, and `init` monitors them for crashes, restarting failed services automatically.

11 Conclusion

This bare-metal OS demonstrates Lateralus's capability for low-level systems programming. The ownership model prevents common kernel bugs including use-after-free and double-free errors. Pipeline-native data processing simplifies interrupt handling and device driver implementations. The complete system boots on real hardware and provides a functional interactive environment.

4.1 Memory Management Unit Configuration

The MMU configuration enables the No-Execute bit in the EFER MSR, allowing page table entries to mark pages as non-executable. This hardware-enforced W^X policy prevents code execution from data pages, providing a fundamental defense against buffer overflow exploits in kernel and user space.

Huge page support uses 2 MB pages for kernel code and data mappings, reducing TLB pressure for large contiguous allocations. The buddy allocator tracks huge-page-aligned blocks separately, ensuring that huge page allocations do not require splitting from the 4 KB pool.

The physical memory allocator uses per-CPU free lists to reduce contention on multiprocessor systems. Each CPU maintains a small cache of free pages that can be allocated without acquiring the global allocator lock. The per-CPU caches are periodically rebalanced to prevent memory hoarding.

Memory zones categorize physical memory by its capabilities: DMA zone (below 16 MB) for legacy device compatibility, Normal zone (16 MB to 4 GB) for general allocations, and High zone (above 4 GB) for user-space mappings. Zone-aware allocation ensures device drivers receive memory accessible by their hardware.

Kernel stack allocation uses a dedicated slab cache that provides guard pages above and below each stack. Guard pages are marked as not-present in the page table, causing an immediate page fault on stack overflow instead of silently corrupting adjacent memory.

The page frame database stores metadata for each physical page including reference count, mapping count, flags, and the owning slab cache. The database is indexed by page frame number and stored in a contiguous array mapped during early boot initialization.

Reverse mapping tracks which page table entries point to each physical page. This information is essential for page reclamation, where the kernel must update or invalidate all mappings before freeing a shared page. The reverse map uses a linked list anchored in the page frame database.

Memory compaction relocates movable pages to consolidate free memory into larger contiguous blocks. Compaction is triggered when a large allocation fails despite sufficient total free memory. The compaction algorithm scans from both ends of a memory zone, migrating pages toward the center.

NUMA awareness distributes allocations across memory nodes based on CPU affinity. Processes preferentially allocate from the memory node local to their running CPU, reducing memory access latency. The NUMA topology is discovered from ACPI tables during boot.

Memory pressure notifications alert subsystems when free memory drops below configurable thresholds. The notification system enables proactive cache trimming and buffer reduction before the allocator runs out of memory, preventing allocation failures during load spikes.

6.1 Advanced Interrupt Management

Interrupt affinity assigns device interrupts to specific CPUs for cache locality. The assignment algorithm considers interrupt frequency, CPU load, and NUMA proximity. Affinity settings are stored

in the interrupt descriptor and can be modified at runtime through the `sysctl` interface.

Message Signaled Interrupts bypass the I/O APIC entirely, delivering interrupt messages directly to the target CPU's Local APIC. MSI support reduces interrupt latency and eliminates the sharing and routing limitations of legacy pin-based interrupts. The kernel prefers MSI when supported by the device.

Interrupt coalescing reduces overhead for high-frequency devices by batching multiple events into a single interrupt. The coalescing timer delays the interrupt delivery for a configurable period, allowing multiple events to accumulate. This technique is particularly effective for network and storage devices.

Software interrupts provide deferred processing for work that is triggered by hardware interrupts but does not require immediate execution. The `softirq` mechanism runs handlers at a lower priority than hardware interrupts, allowing preemption by time-critical device handlers.

The interrupt latency monitor measures the time between interrupt assertion and handler entry. Measurements are collected per-vector and reported through a debug interface. High-latency interrupts indicate performance problems that may require interrupt routing or handler optimization.

Nested interrupt handling allows high-priority interrupts to preempt lower-priority handlers. The kernel maintains a per-CPU interrupt nesting depth counter and switches to a dedicated interrupt stack on the first nesting level. Maximum nesting depth is limited to prevent stack overflow.

Interrupt storm detection identifies vectors that fire at abnormally high rates, typically indicating hardware malfunction or driver bugs. When a storm is detected, the offending vector is temporarily masked and a diagnostic message is logged. The vector is re-enabled after a cooldown period.

Inter-processor interrupts are used for TLB shootdowns, scheduler notifications, and panic broadcasts. The IPI mechanism uses the Local APIC's ICR register to target specific CPUs or broadcast to all processors. IPI handlers are designed for minimal latency to reduce cross-CPU synchronization overhead.

The interrupt controller abstraction supports multiple backend implementations including legacy 8259 PIC, APIC, and x2APIC. The active backend is selected during boot based on hardware detection. The abstraction provides a uniform interface for interrupt masking, acknowledgment, and configuration.

Watchdog timers use the Local APIC timer to detect kernel lockups. The watchdog fires if a CPU does not service the timer interrupt within a configurable deadline. On expiration, the watchdog handler captures CPU state, logs a diagnostic dump, and optionally triggers a panic for post-mortem analysis.

7.1 PCI Express Enumeration

PCI Express enumeration discovers all devices connected to the system bus. The kernel scans the configuration space of each possible bus, device, and function number, reading vendor and device

IDs to identify present devices. Discovered devices are registered in the device tree.

Base Address Registers are programmed to assign MMIO and I/O port ranges to each device. The kernel's resource allocator assigns non-overlapping address ranges from the available memory windows. BAR assignments are stored in the device descriptor for driver use.

The PCI Express capability list is parsed to discover device capabilities including MSI/MSI-X support, power management, advanced error reporting, and SR-IOV. Capabilities are recorded in the device descriptor and used by drivers to configure device features.

Configuration space access uses memory-mapped I/O through the MCFG ACPI table for PCIe devices. The kernel maps the configuration space into its virtual address space during initialization. MMIO-based configuration is significantly faster than the legacy port-based mechanism.

Hot-plug support monitors PCI Express link status changes to detect device insertion and removal. When a new device is detected, the kernel enumerates it, assigns resources, and probes matching drivers. Removed devices trigger driver cleanup and resource deallocation.

Power management transitions PCI Express devices through D0 (active), D1, D2, and D3 (suspended) states. The kernel coordinates device power states with the system power manager, ensuring devices are active when needed and suspended when idle to conserve power.

Error handling uses the Advanced Error Reporting capability to detect and recover from PCI Express errors. Correctable errors are logged and counted. Uncorrectable errors trigger device reset and driver reinitialization. Fatal errors disable the affected link.

DMA remapping through the IOMMU provides memory protection for device-initiated transfers. Each device is assigned an I/O page table that restricts its DMA access to authorized memory regions. This prevents buggy or malicious devices from accessing arbitrary memory.

Vendor-specific driver matching uses a priority system to select the best driver for each device. Drivers register supported vendor/device ID pairs and a priority level. When multiple drivers match a device, the highest-priority driver is selected.

The device tree provides a hierarchical view of all discovered hardware. User-space utilities can query the device tree to list devices, view their configuration, and monitor status changes. The device tree is exported through a virtual filesystem mounted at `/sys/devices`.

8.1 Thread Implementation

Kernel threads provide concurrent execution within the kernel address space. Each kernel thread has its own stack but shares the kernel's page tables and global data structures. Kernel threads are used for background tasks including memory reclamation, journal flushing, and device polling.

Thread-local storage is implemented using the GS segment base register on x86-64. Each thread's TLS area contains per-thread variables including the current thread pointer, error number, and thread-specific allocator state. The TLS area is initialized during thread creation and reclaimed on

thread exit.

Futex-based synchronization provides efficient user-space locking with kernel support for contended cases. Uncontended lock acquisition and release occur entirely in user space without system calls. When contention is detected, the futex system call blocks the waiting thread and manages the wait queue.

Thread groups implement POSIX-style process behavior where multiple threads share an address space, file descriptors, and signal handlers. The thread group leader's PID serves as the process ID. Thread group management handles exit synchronization and zombie reaping.

Priority inheritance prevents priority inversion by temporarily boosting the priority of a thread holding a lock needed by a higher-priority thread. The boosted priority is reverted when the lock is released. The inheritance chain handles multiple levels of nested locking.

Read-copy-update provides a synchronization mechanism for read-heavy data structures. RCU readers access shared data without locks, while writers create new versions and defer reclamation of old versions until all readers have finished. The grace period detection uses per-CPU counters and quiescent state tracking.

Thread migration between CPUs is handled by the load balancer, which periodically checks per-CPU run queue lengths and migrates threads from overloaded to underloaded CPUs. The migration algorithm considers cache warmth, NUMA locality, and migration cost.

Deadline scheduling supports real-time threads with explicit deadline and period parameters. The deadline scheduler uses the Earliest Deadline First algorithm, preempting other threads when a deadline thread becomes runnable. Admission control prevents overcommitting CPU capacity.

Thread debugging support provides per-thread breakpoints, watchpoints, and single-stepping through the ptrace interface. The debug infrastructure uses hardware debug registers where available and falls back to software breakpoints when hardware registers are exhausted.

Wait queue management provides efficient blocking for threads waiting on events. Wait queues support both exclusive and non-exclusive waiters. Exclusive wake-up ensures that only one waiter is activated per event, preventing thundering herd problems for mutual exclusion scenarios.

9.1 Filesystem Journaling

The filesystem journal provides crash consistency by recording metadata changes before they are committed to the main filesystem. The journal uses a write-ahead protocol where journal entries are written and flushed to disk before the corresponding metadata blocks are modified.

Journal transactions group related metadata changes into atomic units. A typical file creation transaction includes inode allocation, directory entry addition, and bitmap updates. The transaction is either fully applied or fully rolled back during recovery, preventing partial updates.

Recovery after an unclean shutdown replays committed journal transactions that were not yet applied

to the main filesystem. The recovery process scans the journal for committed transactions, verifies their checksums, and applies the recorded changes. Recovery completes in seconds regardless of filesystem size.

Journal checkpointing periodically advances the journal head by verifying that old transactions have been fully applied to the main filesystem. Checkpointing frees journal space for new transactions and bounds the recovery time by limiting the number of transactions that need replay.

Ordered data mode ensures that file data blocks are written to disk before the metadata transaction that references them is committed. This ordering prevents newly allocated data blocks from containing stale data from previously deleted files, providing data integrity without full data journaling.

Journal space management reserves a fixed-size region on disk for the circular journal. When the journal fills, new transactions must wait for checkpointing to free space. The journal size is chosen during filesystem creation to balance between recovery speed and transaction throughput.

Asynchronous journal commit batches multiple pending transactions into a single disk write operation. The batching interval is configurable and trades write latency for throughput. Short intervals provide durability at the cost of more disk writes.

Extended attributes are journaled alongside regular metadata, ensuring that security labels, access control lists, and custom metadata maintain consistency across crashes. Extended attribute storage uses both inline (within the inode) and external (separate blocks) formats.

Journal checksums protect against journal corruption from hardware errors. Each journal block includes a CRC32 checksum computed over its contents. During recovery, blocks with invalid checksums are skipped, preventing corrupted journal entries from being applied.

Filesystem quotas track per-user and per-group disk usage in the journal. Quota updates are included in metadata transactions, ensuring accurate accounting even after crashes. The quota system supports both hard limits (allocation fails) and soft limits (warnings and grace periods).

5.1 Kernel Heap Debugging

The kernel heap provides allocation debugging through red zones, guard patterns, and allocation tracking. Red zones are extra bytes added before and after each allocation that are filled with a known pattern. Corruption of the red zone pattern indicates a buffer overflow or underflow.

Use-after-free detection poisons freed memory with a distinctive byte pattern. Any access to the poisoned pattern triggers an immediate assertion failure with diagnostic information including the allocation site and the freeing site. This catches dangling pointer bugs at the point of use.

Double-free detection maintains a free list integrity check. Each freed allocation records a canary value that is verified on subsequent free operations. Attempting to free an already-freed allocation triggers a panic with the original allocation and both free call stacks.

Allocation tracking records the call stack for every allocation and deallocation. The tracking database

is queryable through a debug interface, allowing developers to identify memory leaks by listing allocations without corresponding frees. The overhead is significant, so tracking is a compile-time option.

Slab cache statistics track allocation counts, free counts, cache hit rates, and fragmentation per cache. High fragmentation indicates that object lifetimes are mismatched with slab size, suggesting that the cache configuration should be adjusted.

Out-of-memory handling uses a priority-based reclamation strategy. When allocation fails, the kernel first attempts to reclaim cached filesystem data, then evicts user-space pages, and finally invokes the OOM killer to terminate the process with the largest memory footprint.

Memory leak detection runs periodically as a background kernel thread. The detector scans all kernel memory for pointers to allocated objects, flagging objects that are not referenced from any known root. Unreferenced objects are reported as potential leaks.

Stack usage monitoring tracks the maximum stack depth for each kernel thread. The monitor uses stack coloring during thread creation and checks the highest uncolored address periodically. Threads approaching their stack limit trigger warnings and optional stack expansion.

Allocation size histograms profile the distribution of allocation sizes for each slab cache. The histograms inform cache tuning decisions by revealing whether the fixed cache sizes match actual allocation patterns. Mismatched sizes waste memory through internal fragmentation.

Memory sanitizer integration provides compile-time instrumentation that detects uninitialized memory reads. The sanitizer shadows each byte of kernel memory with a validity bit. Reading an uninitialized byte triggers an error report with the allocation context and the reading location.

10.1 Testing the Kernel

Unit tests verify individual kernel subsystems in isolation. Each subsystem exports a test suite that exercises its public interface with known inputs and expected outputs. Tests run in a special mode where kernel services are mocked, providing deterministic test execution.

Integration tests boot the kernel in QEMU and exercise cross-subsystem interactions. Test scenarios include process creation, file I/O, memory pressure, and device operations. The test harness communicates with the kernel through the serial port, sending commands and collecting results.

Stress tests exercise the kernel under extreme conditions: thousands of processes, memory exhaustion, high interrupt rates, and rapid device insertion/removal. Stress tests expose race conditions, resource leaks, and performance degradation that are invisible under normal load.

Fuzzing uses coverage-guided fuzzing to test system call handling with random inputs. The fuzzer generates sequences of system calls with random arguments and monitors the kernel for crashes, hangs, and assertion failures. Found issues are minimized to the smallest reproducing sequence.

Boot testing verifies that the kernel boots successfully on a variety of hardware configurations. The

test suite uses QEMU with different CPU models, memory sizes, and device configurations. Boot tests catch configuration-specific issues that would otherwise require physical hardware testing.

Performance benchmarking measures throughput and latency for critical operations: context switch time, system call overhead, page fault handling, and filesystem I/O. Benchmarks are run automatically on each commit to detect performance regressions early.

Code coverage analysis identifies untested kernel code paths. Coverage reports highlight functions and branches that are not exercised by any test, guiding test development priorities. The target is comprehensive coverage of all error handling paths.

Static analysis uses the Lateralus compiler's built-in checks plus external analyzers to detect potential bugs without executing the code. Analysis targets include null pointer dereference, integer overflow, lock ordering violations, and unsafe memory access patterns.

Continuous integration builds and tests the kernel on every commit. The CI pipeline compiles the kernel, runs unit and integration tests, checks for compiler warnings, and verifies boot on multiple QEMU configurations. Failed builds block merging.

Crash dump analysis uses a custom tool that decodes kernel crash dumps from QEMU. The tool reconstructs the call stack, register state, and memory layout at the time of the crash. This information accelerates debugging by providing complete context for post-mortem analysis.

2.1 Kernel Image Format

The kernel is compiled as an ELF executable with position-independent code. The ELF format allows the bootloader to parse section headers and load segments at the addresses specified by program headers. Debug information is included in a separate section for debugging with GDB.

The kernel's text section is mapped as read-only and executable, while data and BSS sections are mapped as read-write and non-executable. This separation enforces W^X at the kernel level, preventing injected code from executing even if a buffer overflow vulnerability exists.

Symbol information is preserved in the kernel image for runtime stack trace generation. When a kernel panic occurs, the panic handler walks the call stack and resolves return addresses to function names using the symbol table. This provides meaningful crash diagnostics.

Kernel modules are loaded as relocatable ELF objects. The module loader resolves symbol references against the kernel's exported symbol table, allocates memory for module sections, and applies relocations. Loaded modules become part of the kernel and have full access to kernel interfaces.

Compressed kernel images use LZ4 decompression during boot. The bootloader loads the compressed image and the decompressor stub into memory. The stub decompresses the kernel in place before transferring control. LZ4 is chosen for its fast decompression speed.

The kernel command line is passed from the bootloader through a shared data structure. Command

line parameters configure debugging output, memory limits, driver options, and root filesystem location. The parameter parser supports typed parameters with default values.

Build system integration generates the kernel image, initramfs, and bootloader in a single build command. The build system tracks dependencies between assembly sources, Lateralus sources, and generated headers. Incremental builds recompile only changed components.

Kernel versioning embeds the version string, build timestamp, and compiler version in the kernel image. This metadata is displayed during boot and accessible through the `sysinfo` system call. Module loading verifies version compatibility between the module and the running kernel.

The initramfs is a compressed cpio archive embedded in the kernel image. It contains the init program, essential device drivers as loadable modules, and the filesystem utilities needed to mount the real root filesystem. The kernel extracts the initramfs into a `tmpfs` during early boot.

Link-time optimization produces a single compilation unit from all kernel sources, enabling cross-module inlining and dead code elimination. LTO significantly reduces kernel image size and improves runtime performance at the cost of longer build times.

11.1 Lessons and Design Reflections

Lateralus's ownership model proved invaluable for preventing use-after-free bugs in the memory manager. Page frame lifetimes are tracked by the type system, and the compiler rejects code that accesses freed frames. This catches an entire class of bugs that plague C-based kernels.

Pipeline-native processing simplified interrupt handling by expressing handler chains as pipelines. Device interrupts flow through acknowledgment, dispatch, and completion stages with each stage clearly separated. Error handling at any stage is expressed naturally using pipeline error propagation.

The algebraic type system enables exhaustive matching on hardware status registers. Interrupt vectors, error codes, and device states are modeled as enum variants. The compiler ensures every possible value is handled, preventing the silent ignoring of unexpected hardware conditions.

Foreign function interface support for assembly was essential for hardware interaction. While most kernel code is written in Lateralus, critical low-level operations (GDT loading, context switching, special register access) require inline assembly. The FFI provides safe wrappers around these operations.

Compile-time evaluation generates lookup tables, permission bitmasks, and configuration constants without runtime overhead. The kernel extensively uses compile-time functions to precompute values that would otherwise require runtime initialization, simplifying the boot process.

The experience highlighted areas for Lateralus improvement: better volatile memory access semantics, more flexible inline assembly integration, and support for custom allocators in the standard library. These findings inform the language's continued development.

Testing infrastructure benefited from Lateralus's freestanding library, which provides formatted output

and assertions without operating system dependencies. Unit tests run on the host system using a mock hardware abstraction layer, enabling rapid test iteration.

Build time for the complete kernel is under thirty seconds on modern hardware, thanks to Lateralus's fast compilation and incremental build support. This rapid feedback loop enables an experimental development style where changes can be tested immediately.

The kernel achieves a boot time of approximately two seconds from UEFI handoff to shell prompt on modern hardware. This performance validates the design choices and demonstrates that Lateralus-based systems can compete with hand-optimized C implementations.

Documentation is generated from source code comments using Lateralus's doc system. The kernel documentation includes API references, design documents, and tutorials for extending the kernel. Hyperlinked cross-references connect related concepts across subsystems.

References

- [1] OSDev Wiki. <https://wiki.osdev.org>, 2024.
- [2] Intel 64 and IA-32 Architectures Software Developer Manuals. Intel Corporation, 2024.
- [3] UEFI Specification, Version 2.10. UEFI Forum, 2023.
- [4] Tanenbaum, A.S. Modern Operating Systems, 4th Edition. Pearson, 2014.
- [5] Love, R. Linux Kernel Development, 3rd Edition. Addison-Wesley, 2010.
- [6] Arpaci-Dusseau, R. and Arpaci-Dusseau, A. Operating Systems: Three Easy Pieces. 2018.