

Project Charter: Linux Kernel Module

What You Are Building

You are building a production-grade Linux character device driver—a loadable kernel module (`.ko`) that acts as a bridge between the Linux kernel and userspace. This isn't a simple "Hello World"; you will build a virtual device that manages a kernel-space buffer, exposes a `/dev/` file interface for data transfer, implements a custom `ioctl` control protocol, and provides a `/proc` entry for live system introspection. By the end, your driver will safely handle concurrent access from multiple processes using kernel synchronization primitives.

Why This Project Exists

Most developers interact with the operating system through high-level abstractions, treating the kernel as a "black box." Building a kernel module from scratch is the only way to truly understand the "everything is a file" philosophy. You will confront the hardware-enforced boundary between userspace and kernelspace, learning why direct pointer dereferencing is fatal and how the kernel manages memory, scheduling, and concurrency at the highest privilege level (Ring 0).

What You Will Be Able to Do When Done

- **Develop and Debug LKMs:** Compile, load, and unload out-of-tree modules using the Kbuild system and inspect logs via `dmesg`.
- **Bridge the Memory Divide:** Use `copy_to_user` and `copy_from_user` to safely move data across the kernel-userspace boundary.
- **Implement VFS Operations:** Write handlers for standard system calls like `open()`, `read()`, `write()`, and `close()`.
- **Design Binary Protocols:** Create custom control interfaces using `ioctl` with unique magic numbers and encoded command structures.
- **Master Kernel Concurrency:** Use `mutexes` to prevent race conditions and `wait queues` to implement blocking I/O and process sleeping.
- **Support Multiplexed I/O:** Implement the `.poll` operation to make your device compatible with `select()`, `poll()`, and `epoll()`.

Final Deliverable

A robust C source file (approx. 500 lines) and a Kbuild-compatible `Makefile`. The final module creates `/dev/mydevice` (for data) and `/proc/mydevice` (for stats). It must pass a concurrent stress test where 8

processes (4 readers, 4 writers) access the device simultaneously without data corruption, deadlocks, or a "Kernel Oops."

Is This Project For You?

You should start this if you:

- Are proficient in C (specifically pointers, structs, and manual memory management).
- Understand basic concurrency concepts like race conditions and locks.
- Are comfortable working in a Linux terminal and using `sudo` privileges.

Come back after you've learned:

- [C Programming Pointers & Structs](#)
- [Basic Linux Signal Handling](#) (If you don't know why `read()` might be interrupted).

Estimated Effort

Phase	Time
M1: Hello World & Metadata (Setup, Kbuild, and Params)	~4 hours
M2: Character Device Core (Read/Write & /dev creation)	~8 hours
M3: Control & Introspection (ioctl and /proc seq_file)	~8 hours
M4: Concurrency & Poll (Mutex, Wait Queues, and Poll)	~10 hours
Total	~30 hours

Definition of Done

The project is complete when:

- The module compiles with `-Werror` and loads without "tainting" the kernel (proper GPL license).
- A userspace `echo` can write data to `/dev/mydevice` and `cat` can read it back.
- Custom `ioctl` commands successfully resize the internal kernel buffer at runtime.
- `/proc/mydevice` displays real-time, accurate counts of total bytes read and written.
- The driver survives the "Thundering Herd" stress test (multiple concurrent readers/writers) without a kernel panic or data loss.



Before You Read This: Prerequisites & Further Reading

Read these first. The Atlas assumes you are familiar with the foundations below. Resources are ordered by when you should encounter them — some before you start, some at specific milestones.

1. The Kernel-Userspace Boundary

Topic: System Call Architecture & Memory Protection

- **Paper:** P. J. Braam, 2003, "The Linux Virtual File System."
- **Code:** `arch/x86/entry/entry_64.S` — The raw assembly handling the transition from ring 3 to ring 0.
- **Best Explanation:** *Linux Kernel Development* (3rd Ed) by Robert Love, Chapter 5: "System Calls."
- **Why:** It provides the most readable explanation of how the processor physically switches privilege levels and how the kernel validates the syscall table.
- **Timing:** Read **BEFORE Milestone 1** to understand the "physics" of the environment you are entering.
- **Spec:** POSIX.1-2017, Section 2.9.7: "Thread-Safety" and "Memory Sanity."
- **Code:** `arch/x86/lib/copy_user_64.S` — Look for the `_copy_to_user` implementation and the `.fixup` section.
- **Best Explanation:** LWN.net, "Accessing user-space memory" by Jonathan Corbet.
- **Why:** This article explains the "Trap and Recover" mechanism (exception tables) that makes `copy_from_user` safe without pre-checking pointers.
- **Timing:** Read **during Milestone 2** after your first kernel panic to understand why your direct pointer dereference failed.

2. Character Device Architecture

Topic: VFS & Device Registration

- **Code:** `fs/char_dev.c` — The core implementation of the major/minor number registry.
- **Best Explanation:** *Linux Device Drivers* (3rd Ed) by Corbet, Rubini, & Kroah-Hartman, Chapter 3: "Char Drivers."
- **Why:** Though written for older kernels, the architectural pattern of `struct cdev` and the VFS dispatch remains the definitive pedagogical map.

- **Timing:** Read **at the start of Milestone 2** to see the "contract" you are about to fulfill.
 - **Spec:** Linux Kernel Documentation: `admin-guide/devices.txt`.
 - **Code:** `include/linux/fs.h` — Focus on the `struct file_operations` definition.
 - **Best Explanation:** "The Linux VFS" by Neil Brown (LWN series).
 - **Why:** It explains how the VFS provides polymorphism in C, allowing a socket, a file, and your driver to share the same `read()` interface.
 - **Timing:** Read **after Milestone 2** to appreciate how your driver fits into the global Linux filesystem tree.
-

3. Memory Management

Topic: SLUB Allocator & Kernel Paging

- **Paper:** Jeff Bonwick, 1994, "The Slab Allocator: An Object-Caching Kernel Memory Allocator."
 - **Code:** `mm/slub.c` — Specifically the `slab_alloc_node` function.
 - **Best Explanation:** *Understanding the Linux Kernel* by Bovet & Cesati, Chapter 8: "Memory Management."
 - **Why:** It explains the "Bento Box" model of memory allocation and why `kmalloc` is faster than `malloc`.
 - **Timing:** Read **during Milestone 2** when you call `kzalloc()` for the first time.
-

4. Control Plane & Introspection

Topic: ioctl Design & Virtual Filesystems

- **Spec:** Linux Kernel Documentation: `userspace-api/ioctl/ioctl-number.rst`.
- **Code:** `include/uapi/asm-generic/ioctl.h` — The macros that define the 32-bit command encoding.
- **Best Explanation:** "The new way of doing ioctl()" (LWN) regarding the transition to `unlocked_ioctl`.
- **Why:** Essential for understanding why the Big Kernel Lock was removed and how modern ioctl dispatch works.
- **Timing:** Read **at the start of Milestone 3** to avoid defining "garbage" ioctl numbers that conflict with other drivers.
- **Code:** `fs/seq_file.c` — The implementation of the iterator pattern for `/proc`.
- **Best Explanation:** Linux Kernel Documentation: "The seq_file Interface."
- **Why:** It is the only resource that correctly explains why you don't need to worry about 4KB buffer limits when using `seq_printf`.

- **Timing:** Read **during Milestone 3** while implementing your `/proc` entry.
-

5. Concurrency & Synchronization

Topic: Mutexes, Wait Queues, and Scheduling

- **Code:** `kernel/locking/mutex.c` — Look for the "Fast path" vs "Slow path" logic.
 - **Best Explanation:** *Linux Kernel Development* by Robert Love, Chapter 9: "Kernel Synchronization Methods."
 - **Why:** Love differentiates between spinlocks and mutexes based on "context," which is the single most important concept for kernel stability.
 - **Timing:** Read **BEFORE starting Milestone 4**; it is required knowledge to prevent deadlocks.
 - **Paper:** S. Molloy, 2006, "The 'Thundering Herd' Problem in Network Servers."
 - **Code:** `kernel/sched/wait.c` — The `prepare_to_wait` and `finish_wait` functions.
 - **Best Explanation:** "Wait queues in the Linux kernel" (KernelNewbies Wiki).
 - **Why:** It provides a line-by-line breakdown of how `sleep()` actually informs the scheduler to stop running a task.
 - **Timing:** Read **during Milestone 4** while debugging why your reader process won't wake up.
-

6. Event Multiplexing

Topic: Poll, Select, and Epoll

- **Code:** `fs/select.c` — The `do_poll` function loop.
 - **Best Explanation:** "The Implementation of epoll (1)" by Marek Majkowski.
 - **Why:** It traces the `poll_wait` callback from your driver all the way up to the `epoll_wait` system call in userspace.
 - **Timing:** Read **at the end of Milestone 4** to understand how your driver enables high-performance servers (Nginx/Node.js).
-

7. Reliability & Debugging

Topic: Tainting & Oops Analysis

- **Best Explanation:** Linux Kernel Documentation: `admin-guide/tainted-kernels.rst`.
- **Why:** You need to know what the letters in `Tainted: P G W` mean when your kernel crashes.
- **Timing:** Read **the first time you see a "Kernel Oops"** in your dmesg logs.

- **Code:** `scripts/decode_stacktrace.sh` in the Linux source tree.
 - **Best Explanation:** "Kernel Address Space Layout Randomization (KASLR)" (LWN).
 - **Why:** It explains why your crash addresses look different every time you reboot.
 - **Timing:** Read **after Milestone 4** to transition from a student to a professional systems engineer.
-

Linux Kernel Module: Interactive Atlas

This project builds a fully functional Linux character device driver from scratch — starting with a minimal loadable kernel module and progressively adding file operations, ioctl control interfaces, /proc introspection, and concurrent access handling with wait queues and poll support. The journey traverses the most fundamental boundary in modern computing: the line between kernel space and user space. Every system call you've ever made — `read()`, `write()`, `open()` — terminates in code that looks exactly like what you'll write here.

By the end, you'll have a character device that multiple processes can open, read, write, and poll concurrently, with proper mutex synchronization, blocking/non-blocking I/O semantics, and runtime introspection via /proc. This isn't a toy: it's the same architecture used by `/dev/null`, `/dev/random`, `/dev/tty`, and every hardware driver in the Linux kernel.

The real learning isn't just the API — it's understanding WHY the API looks the way it does. Why does `copy_to_user()` exist instead of `memcpy()`? Why does `wait_event_interruptible()` return `-ERESTARTSYS`? Why must `poll()` both register a wait queue AND return a mask? Each answer reveals a deep design constraint that ripples across all of systems programming.

Milestone 1: Hello World Kernel Module

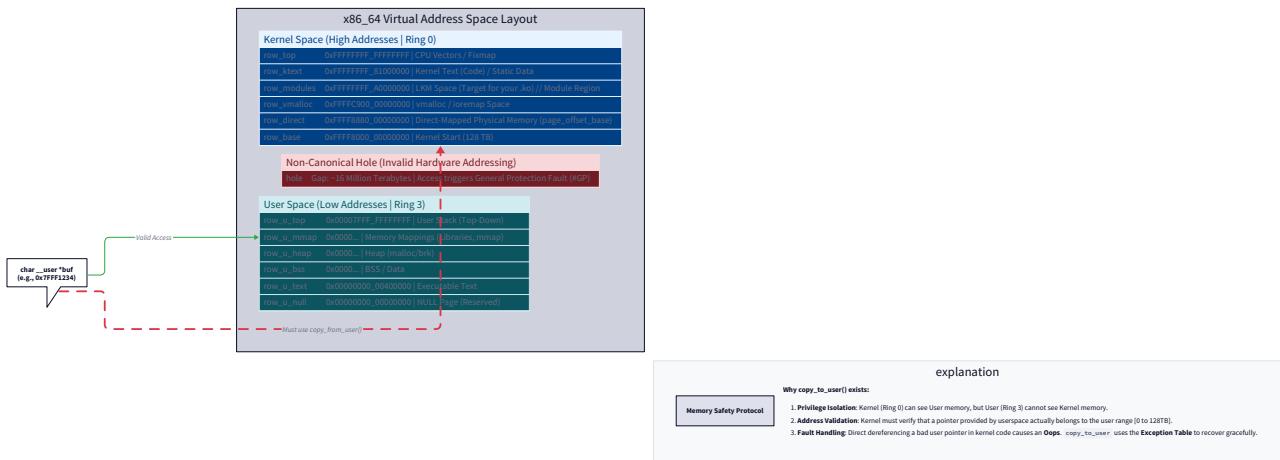
The Ground Beneath Your Feet

Before you write a single line of kernel code, you need to confront a misconception that trips up nearly every developer coming to kernel

programming for the first time. It's worth addressing head-on because getting it wrong doesn't just cause bugs — it causes *kernel panics*.

The Revelation: You Are Not Writing a Library

You probably know that Linux supports dynamically loaded shared libraries — `.so` files that a process loads at runtime via `dlopen()`. The mental model is clean: the library gets mapped into the process's address space, the dynamic linker resolves symbols, the library runs in the same process context as the caller, and if the library crashes, it crashes the process (not the whole system). You've worked with this model before. When you hear "loadable kernel module," your brain reaches for this analogy: *it's probably like a kernel-side .so file, loaded on demand, running in some contained kernel process, debuggable with gdb.* **Every part of that model is wrong.** A kernel module is not loaded *into* a process. It is injected *into the kernel itself* — the same 40-million-line entity that manages every process, every memory mapping, every hardware interrupt on your machine. There is no isolation. There is no "kernel process" running your module's code in a sandbox. When your module's `init` function runs, it runs with **ring 0 CPU privilege** — the highest privilege level the hardware supports. It can access any physical memory address. It can execute any hardware instruction. It can corrupt the kernel's own data structures. When your module's functions are called later (by processes making system calls), they run directly in the kernel's address space, on whatever CPU the calling process happened to be scheduled on, with no process boundary separating them from the kernel's internals. The consequence: a NULL pointer dereference in your module doesn't produce a segmentation fault. It produces a **kernel oops** — a diagnostic dump showing the CPU state, register values, and stack trace at the moment of the crash — and potentially a kernel panic that halts the machine. There is no signal handler to catch this, no garbage collector to save you, no memory-safe language to interpose. You are the kernel.



🔑 Foundation: Kernel vs userspace address spaces

- 1. What it IS** In a modern operating system, the CPU uses virtual memory to isolate processes. This virtual memory is split into two primary regions: **Userspace** and **Kernelspace**.

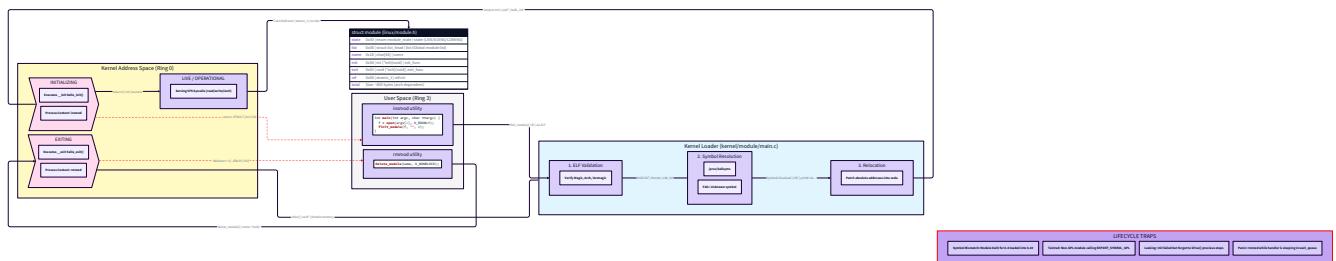
- **Userspace** is the restricted sandbox where your applications (like `ls`, `web browsers`, or `python`) run. They cannot touch hardware directly and have a limited view of memory.
- **Kernelspace** is the privileged area where the core operating system resides. It has full access to the CPU instructions and all physical hardware.

2. WHY you need it right now If you are writing a device driver or a kernel module, you are operating in Kernelspace. You cannot simply "pass a pointer" from a user application to the kernel and dereference it. Because each user process has its own isolated address space, the address `0x400500` in a user program might point to a string, but in the kernel, that same address could point to something entirely different—or nothing at all. You must use specific functions (like `copy_from_user`) to safely bridge this gap.

3. Key Insight: The "One-Way Mirror" Mental Model Think of the Kernel as the landlord of an apartment building and Userspace processes as the tenants. Each tenant lives in their own apartment (Address Space) and can't see inside others. The landlord lives in a private office (Kernelspace). The landlord can look into any apartment to move furniture or fix pipes, but the tenants can only interact with the landlord through a specific mail slot (System Calls). They can never walk into the landlord's office uninvited.

This isn't meant to frighten you — it's meant to orient you. Kernel development has a different *physics* than userspace development. The same discipline that makes you check every `malloc()` return value in userspace needs to be applied with ten times the rigor here, because the cost of failure is ten times higher. Once you internalize this, the entire design of the kernel API — the strange-looking macros, the mandatory error checks, the careful memory barriers — starts making sense. Every API is designed the way it is because someone suffered the alternative.

What a Module Actually Is



A kernel module is an ELF object file (`.ko` extension, for "kernel object") that the kernel's module loader reads, relocates, and links into the live kernel image at runtime. [[EXPLAIN:how-the-kernel-loads-and-links-modules-(.ko-files,-symbol-resolution)|How the kernel loads and links modules (.ko files, symbol resolution)]] The `.ko` file is structurally similar to a `.o` file from your compiler — it has code sections, data sections, relocation entries, and a symbol table. The critical difference is what it links *against*: not libc, not any userspace library, but the kernel's own symbol table — the list of all functions and variables the running kernel has explicitly exported for module use. Functions like `printk`, `kmalloc`, `class_create` — these are kernel symbols that your module will reference, and they get resolved at load time against the kernel that's actually running. This is why the kernel headers you compile against must exactly match the running kernel. If your headers say `struct file_operations` has 25 fields but the running kernel's struct has 26 fields, your module will write to wrong memory offsets when it fills that struct. The result is not a compile error — it's silent data corruption that may not manifest until some seemingly unrelated kernel path trips over your corrupted data.

The Kbuild System: Compiling Into the Kernel

You won't use a plain `Makefile` with `gcc` directly. The Linux kernel uses its own build system, **Kbuild**, which knows how to compile code that will be linked against the kernel. Kbuild handles the complex flags required to build kernel code correctly — position-independent code settings, stack protector flags, no standard library includes, kernel-specific warning flags, and more. {{DIAGRAM:diag-m1-kbuild-flow}} For an out-of-tree module (one that lives outside the kernel source tree), you write a two-line `Makefile` that delegates to Kbuild:

```
# Makefile – out-of-tree kernel module build  
obj-m += hello.o  
all:  
    make -C /lib/modules/$(shell uname -r)/build M=$(PWD) modules  
clean:  
    make -C /lib/modules/$(shell uname -r)/build M=$(PWD) clean
```

MAKEFILE

Let's read this carefully:

- `obj-m += hello.o` — tells Kbuild to compile `hello.c` and link it as a module (the `-m` suffix means "module", as opposed to `-y` which means "built into the kernel image"). Kbuild automatically maps `hello.o` to `hello.c` as the source.
- `-C /lib/modules/$(shell uname -r)/build` — changes to the kernel build directory for the **currently running kernel** (`uname -r` returns something like `6.8.0-47-generic`). The kernel build directory contains the Kbuild infrastructure and the kernel configuration used to build the running kernel.
- `M=$(PWD) modules` — tells Kbuild that the module source lives in the current directory (`$(PWD)`), and to build the `modules` target (which compiles out-of-tree modules). The practical implication: you must have the **kernel headers** installed for your running kernel. On Ubuntu/Debian:

```
sudo apt install linux-headers-$(uname -r)
```

BASH

On Fedora/RHEL:

```
sudo dnf install kernel-devel-$(uname -r)
```

BASH

If the headers aren't installed, the build directory won't exist and `make` will fail immediately.

Hardware Soul: Compiling kernel code requires disabling features like stack canaries in certain contexts, enabling specific code model flags (`-mcmodel=kernel` on `x86_64`) that constrain how code addresses memory, and ensuring no floating-point instructions are emitted (the kernel doesn't save/restore FPU state except when explicitly crossing to/from userspace). Kbuild applies all of these automatically. If you tried to compile with `gcc hello.c -o hello.ko` directly, you'd produce an object that *looks* correct but would corrupt the kernel when loaded.

Writing the Minimal Module

Here is the smallest valid kernel module:

```
// hello.c - minimal kernel module

#include <linux/module.h>      /* MODULE_LICENSE, MODULE_AUTHOR, etc. */

#include <linux/kernel.h>      /* printk, KERN_INFO */

#include <linux/init.h>        /* __init, __exit */

MODULE_LICENSE("GPL");

MODULE_AUTHOR("Your Name");

MODULE_DESCRIPTION("A minimal kernel module demonstrating load/unload lifecycle");

MODULE_VERSION("1.0");

static int __init hello_init(void)

{

    printk(KERN_INFO "hello: module loaded\n");

    return 0;

}

static void __exit hello_exit(void)

{

    printk(KERN_INFO "hello: module unloaded\n");

}

module_init(hello_init);

module_exit(hello_exit);
```

Every line here does something specific. Let's walk through each piece.

Include Headers — The Kernel's Standard Library

```
#include <linux/module.h>

#include <linux/kernel.h>

#include <linux/init.h>
```

These are not the headers from `/usr/include` that your userspace programs use. They come from the kernel source tree under `include/linux/`. Crucially, you **never** include `<stdio.h>`, `<stdlib.h>`,

`<string.h>`, or any other C standard library header. The C standard library is a userspace library — it doesn't exist in the kernel. The kernel has its own implementations of common utilities (`strlen`, `memcpy`, `sprintf`) in `<linux/string.h>`, `<linux/kernel.h>`, etc. If you accidentally include `<stdio.h>` in kernel code, the compiler will fail in spectacular ways — the standard library headers include syscall wrappers and type definitions that conflict with kernel type definitions.

Module Metadata Macros

```
MODULE_LICENSE("GPL");  
  
MODULE_AUTHOR("Your Name");  
  
MODULE_DESCRIPTION("A minimal kernel module demonstrating load/unload lifecycle");  
  
MODULE_VERSION("1.0");
```

C

These macros don't generate function calls or runtime initialization code. They inject **strings into dedicated ELF sections** in the `.ko` file. The `modinfo` tool reads these sections directly from the file, and the kernel's module loader reads `MODULE_LICENSE` to decide whether to grant access to GPL-only symbols.

`MODULE_LICENSE("GPL")` is not a legal formality — it is **functionally required**. The Linux kernel uses the GNU Public License

🔑 Foundation: GPL symbol export mechanism: how EXPORT_SYMBOL_GPL works and what "tainting" means

1. What it IS The Linux kernel is modular. To allow one module to use functions (symbols) from another, the kernel uses the `EXPORT_SYMBOL` macro. There are two tiers:

- **EXPORT_SYMBOL**: Makes a function available to any module, regardless of license.
- **EXPORT_SYMBOL_GPL**: Makes a function available *only* to modules that declare a GPL-compatible license.

Tainting is a state flag in the kernel. If you load a module that is proprietary (non-GPL) or if the kernel encounters a catastrophic hardware error, the kernel marks itself as "Tainted."

2. WHY you need it right now When developing a kernel module, your choice of license

(`MODULE_LICENSE("GPL")`) determines which internal APIs you can access. If you try to use a function exported via `EXPORT_SYMBOL_GPL` in a non-GPL module, the kernel will refuse to load your module. Furthermore, if you are debugging a crash and your kernel is "Tainted," kernel developers will often ignore your bug report because proprietary code makes it impossible to verify the kernel's internal state.

3. Key Insight: The "Void if Seal Broken" Sticker Think of "Tainting" like the warranty sticker on a piece of electronics. Loading a non-GPL module or forcing a module into the kernel is like peeling that sticker off. The device (the kernel) might still work perfectly fine, but the "manufacturer" (the Linux community) will no longer

support it because they can't be sure what happened inside the box once the "seal" was broken. to gate access to certain symbols. Functions like `__alloc_pages`, most of the crypto API, and many driver APIs are exported only to GPL-licensed modules via `EXPORT_SYMBOL_GPL()`. If your module declares any other license (or omits the macro), loading it will "taint" the kernel — marking it as running non-GPL code — and the symbol resolution for any `EXPORT_SYMBOL_GPL` symbol your module needs will fail with an `Unknown symbol` error. The accepted values for `MODULE_LICENSE` are: `"GPL"`, `"GPL v2"`, `"GPL and additional rights"`, `"Dual BSD/GPL"`, `"Dual MIT/GPL"`, `"Dual MPL/GPL"`, and `"Proprietary"`. The last one taints the kernel and blocks GPL-only symbols.

`__init` and `__exit` — Section Annotations

```
static int __init hello_init(void) { ... }

static void __exit hello_exit(void) { ... }
```

The `__init` and `__exit` macros are **section annotations** that place the decorated functions into special ELF sections (`.init.text` and `.exit.text` respectively). For functions in `.init.text`: after the module successfully initializes, the kernel **frees the memory** occupied by those functions. The init code runs exactly once — never again after initialization — so keeping it resident wastes RAM. On resource-constrained systems (embedded, server with thousands of modules), this matters. For functions in `.exit.text`: if the module is compiled as a built-in (not loadable, but baked into the kernel image at compile time), the exit functions are discarded entirely — a built-in module can never be unloaded, so the exit function is dead code. The `static` keyword matters too: it prevents these functions from being exported as global symbols, which keeps the module's symbol footprint minimal and prevents accidental name collisions with other kernel symbols.

`module_init()` and `module_exit()` — Registration

```
module_init(hello_init);

module_exit(hello_exit);
```

These macros register your init and exit functions with the kernel's module framework. They don't call the functions — they tell the kernel loader "when this module is loaded, call `hello_init`; when it's unloaded, call `hello_exit`." The return value of the init function is critical: returning `0` means success, returning a negative error code (like `-ENOMEM` for out-of-memory, or `-EBUSY` for a resource conflict) causes `insmod` to fail and

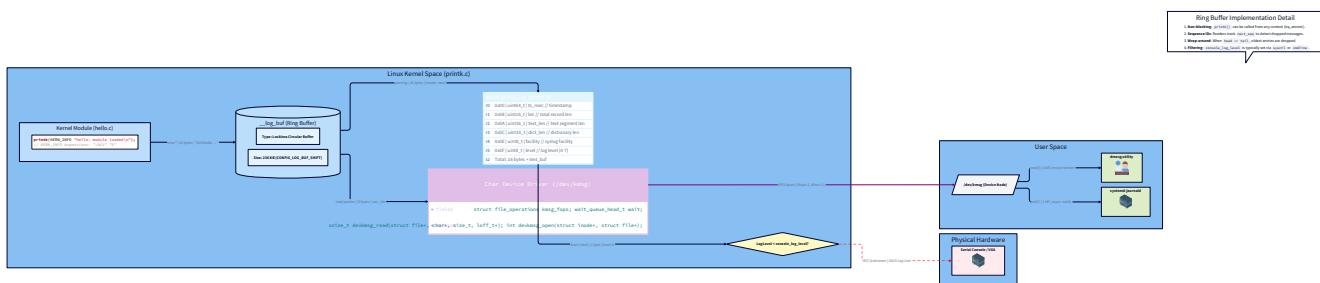
the module to not be loaded. The error code propagates back to the `insmod` process as the exit status. Never return a positive value from `init` — positive values have undefined behavior in this context. If `init` returns an error, the kernel guarantees that `exit` will not be called. This means your `init` function must clean up any resources it already allocated before returning the error — a partial initialization that returns an error but leaves resources allocated will leak them permanently.

`printk` — Logging Without a Terminal

Your first instinct when debugging might be to reach for `printf`. In the kernel, `printf` doesn't exist. Instead you use `printk`:

```
printk(KERN_INFO "hello: module loaded\n");
```

Note the unusual syntax: `KERN_INFO` is not a separate argument separated by a comma — it's a string **concatenated** with the format string at compile time. `KERN_INFO` expands to the string `"\001\006"` (SOH followed by '6', the log level number), and the C compiler concatenates adjacent string literals. The result is a single string: `"\001\006hello: module loaded\n"`.



The log levels in order from most to least severe:

Macro	Level	Numeric	Meaning
KERN_EMERG	0	Emergency	System is unusable
KERN_ALERT	1	Alert	Immediate action required
KERN_CRIT	2	Critical	Critical condition
KERN_ERR	3	Error	Error condition
KERN_WARNING	4	Warning	Warning condition
KERN_NOTICE	5	Notice	Normal but significant
KERN_INFO	6	Info	Informational
KERN_DEBUG	7	Debug	Debug-level message
<code>printf</code> writes its output to the kernel ring buffer — a circular buffer in kernel memory (typically 512KB to 4MB, configurable). The <code>dmesg</code> command reads this buffer by reading <code>/dev/kmsg</code> (which is itself a character device — you'll be building something structurally similar in Milestone 2). Because the buffer is circular, old messages get overwritten when the buffer fills up.			

The ring buffer is not `stdout`. `printf` doesn't write to a terminal. If the console log level is high enough (controlled by `/proc/sys/kernel/printk`), messages may also appear on the system console, but the primary output path is the ring buffer accessible via `dmesg`.

```
# After insmod hello.ko, check dmesg for your message
dmesg | tail -5
# Or follow in real-time with -w (like tail -f)
dmesg -w
```

BASH

The timestamp format in dmesg output ([12345.678901]) is seconds since boot, not wall clock time. Use `dmesg -T` to convert to human-readable timestamps.

Module Parameters: Runtime Configurability

A static module is limited. Real modules need to be tunable: buffer sizes, timeout values, debug levels, device addresses. The `module_param()` macro system exposes module-level variables as parameters settable at load time and (optionally) at runtime via `/sys`.

```
#include <linux/moduleparam.h>

static int buffer_size = 4096; /* default: 4KB */

module_param(buffer_size, int, 0644);

MODULE_PARM_DESC(buffer_size, "Size of the device buffer in bytes (default: 4096)");
```

The `module_param(name, type, permissions)` macro takes three arguments:

1. `name` : The variable name. Must be a static module-level variable — this is **not** a copy; the macro makes the variable itself accessible through sysfs.
2. `type` : The kernel's type tag for the parameter. Common types:
 - `int` , `uint` , `long` , `ulong` — integer types
 - `bool` — boolean (accepts `y` / `n` / `1` / `0`)
 - `charp` — char pointer (string, kernel allocates memory)
 - `byte` — unsigned char
3. `permissions` : The filesystem permission bits for the sysfs attribute. `0644` means owner can read/write, everyone else can read. Use `0` to make the parameter not appear in sysfs at all (load-time only). **Never use `0666`** — that lets any unprivileged user modify kernel state, which is a security vulnerability. After loading, the parameter appears at:

```
/sys/module/<module_name>/parameters/<parameter_name>
```

So for our module:

```
# Load with non-default parameter
sudo insmod hello.ko buffer_size=8192

# Check current value
cat /sys/module/hello/parameters/buffer_size

# Output: 8192

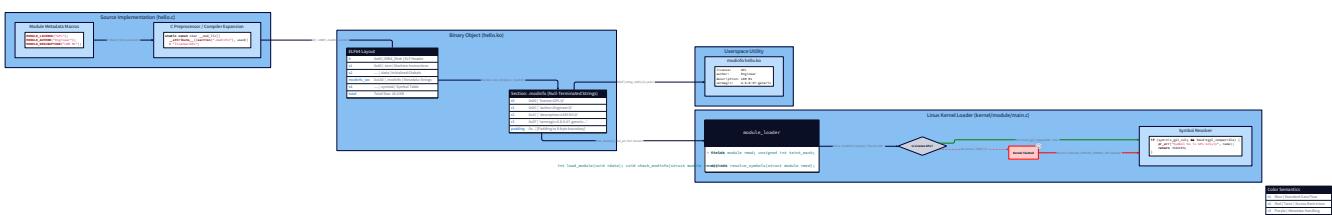
# If permissions allow writing, modify at runtime
echo 16384 | sudo tee /sys/module/hello/parameters/buffer_size
```

The sysfs attribute is directly backed by the `buffer_size` variable — reading it returns the current value, and (if writable) writing it updates the variable in the running kernel. This is the same mechanism that exposes hardware topology under `/sys/class` and `/sys/bus` — module parameters are first-class sysfs citizens, not a special case.

Design note: Module parameters are a double-edged sword. They make modules flexible, but they add code paths that depend on runtime values rather than compile-time constants. A module that worked correctly with `buffer_size=4096` might behave pathologically with `buffer_size=0` or `buffer_size=INT_MAX` if you don't validate. Always validate parameters in your init function.

The Module Metadata Flow

Let's trace exactly how the `MODULE_LICENSE`, `MODULE_AUTHOR`, and `MODULE_DESCRIPTION` macros become visible in `modinfo` output.



The macros expand to `__attribute__((section(".modinfo")))` annotated string arrays. For example, `MODULE_LICENSE("GPL")` roughly expands to:

```
static const char __UNIQUE_ID_license[] __attribute__((section(".modinfo"),used))
= "license=GPL";
```

The Kbuild system also injects additional metadata into `.modinfo`: the kernel version the module was compiled against (`vermagic`), the module dependencies, and any `module_param` descriptions. When you run `modinfo hello.ko`, it reads the `.modinfo` ELF section and prints the key-value pairs:

```
filename:      /path/to/hello.ko
version:       1.0
description:   A minimal kernel module demonstrating load/unload lifecycle
author:        Your Name
license:       GPL
srcversion:    A3B2C1D4E5F6A7B8C9D0E1F
depends:
retpoline:     Y
name:          hello
vermagic:      6.8.0-47-generic SMP preempt mod_unload modversions
parm:          buffer_size:Size of the device buffer in bytes (default: 4096) (int)
```

The `vermagic` field is the kernel's version check mechanism: when `insmod` loads the `.ko`, the kernel compares the embedded `vermagic` string against the running kernel's own version string. If they don't match, loading fails with `disagrees about version of symbol module_layout`. This is the safety guard that prevents you from loading a module compiled against kernel 6.7 into a kernel running 6.8.

The Complete Module: Putting It Together

Here is the full module source with parameter validation and proper cleanup structure:

```
// hello.c - Hello World kernel module with parameter

#include <linux/module.h>

#include <linux/kernel.h>

#include <linux/init.h>

#include <linux/moduleparam.h>

MODULE_LICENSE("GPL");

MODULE_AUTHOR("Your Name <you@example.com>");

MODULE_DESCRIPTION("Hello World kernel module with buffer_size parameter");

MODULE_VERSION("1.0");

/* Module parameter: buffer size with default of 4096 bytes */

static int buffer_size = 4096;

module_param(buffer_size, int, 0644);

MODULE_PARM_DESC(buffer_size, "Size of the device buffer in bytes (default: 4096)");

static int __init hello_init(void)

{

    /* Validate the parameter before using it */

    if (buffer_size <= 0 || buffer_size > (1024 * 1024)) {

        printk(KERN_ERR "hello: invalid buffer_size %d (must be 1 to 1048576)\n",

               buffer_size);

        return -EINVAL;

    }

    printk(KERN_INFO "hello: module loaded, buffer_size=%d\n", buffer_size);

    printk(KERN_INFO "hello: parameter visible at "

           "/sys/module/hello/parameters/buffer_size\n");

    return 0; /* 0 = success; negative = failure */

}

static void __exit hello_exit(void)
```

C

```
{  
  
    printk(KERN_INFO "hello: module unloaded\n");  
  
}  
  
module_init(hello_init);  
  
module_exit(hello_exit);
```

Building It

```
# In the directory containing hello.c and Makefile:  
  
make  
  
# Expected output (abbreviated):  
  
# make -C /lib/modules/6.8.0-47-generic/build M=/home/user/hello modules  
  
# CC [M] /home/user/hello/hello.o  
  
# MODPOST /home/user/hello/Module.symvers  
  
# CC [M] /home/user/hello/hello.mod.o  
  
# LD [M] /home/user/hello/hello.ko
```

The build produces several files, but `hello.ko` is the one you load. The `.mod.c` file is auto-generated by Kbuild and contains metadata; `.symvers` records the symbols your module references and their CRC values (used for module versioning checks).

Loading and Verifying

```
# Load the module                                         BASH
sudo insmod hello.ko

# Verify it's loaded

lsmod | grep hello

# Output: hello      16384  0

# Check dmesg for init message

dmesg | tail -3

# [12345.678901] hello: module loaded, buffer_size=4096

# [12345.678923] hello: parameter visible at /sys/module/hello/parameters/buffer_size

# Check module metadata

modinfo hello.ko

# Check parameter in sysfs

cat /sys/module/hello/parameters/buffer_size

# 4096

# Load with custom parameter

sudo rmmod hello

sudo insmod hello.ko buffer_size=8192

cat /sys/module/hello/parameters/buffer_size

# 8192

# Unload

sudo rmmod hello

dmesg | tail -1

# [12367.891234] hello: module unloaded
```

The Verification Script

A repeatable verification script is better than manual commands:

BASH

```
#!/bin/bash

# verify.sh – automated load/unload verification for hello module

set -e

MODULE_NAME="hello"

KO_FILE="./hello.ko"

echo "[1] Building module..."

make -s

echo "[2] Checking modinfo metadata..."

modinfo "$KO_FILE" | grep -E "^(license|author|description|parm):"

# Verify GPL license, author, description, and parameter are present

echo "[3] Loading module with default parameter..."

sudo insmod "$KO_FILE"

echo "[4] Verifying module is listed in lsmod..."

lsmod | grep -q "$MODULE_NAME" && echo "  PASS: module visible in lsmod"

echo "[5] Verifying init message in dmesg..."

dmesg | tail -5 | grep -q "${MODULE_NAME}: module loaded" && \
      echo "  PASS: init message in dmesg"

echo "[6] Verifying parameter in /sys..."

PARAM_PATH="/sys/module/${MODULE_NAME}/parameters/buffer_size"

[ -f "$PARAM_PATH" ] && echo "  PASS: parameter file exists"

cat "$PARAM_PATH" | grep -q "4096" && echo "  PASS: default value is 4096"

echo "[7] Unloading module..."

sudo rmmod "$MODULE_NAME"

echo "[8] Verifying exit message in dmesg..."

dmesg | tail -3 | grep -q "${MODULE_NAME}: module unloaded" && \
      echo "  PASS: exit message in dmesg"

echo "[9] Loading with custom parameter..."
```

```
sudo insmod "$KO_FILE" buffer_size=8192  
cat "$PARAM_PATH" | grep -q "8192" && echo " PASS: custom parameter accepted"  
sudo rmmod "$MODULE_NAME"  
echo ""  
echo "All checks passed."
```

Common Pitfalls and Their Consequences

This section covers the mistakes that will cost you time — and in some cases, machine reboots.

Missing `MODULE_LICENSE("GPL")` — Tainted Kernel, Missing Symbols

If you omit `MODULE_LICENSE`, the kernel marks itself as "tainted" upon loading your module and prints a warning:

```
hello: module license 'unspecified' taints kernel.  
Disabling lock debugging due to kernel taint
```

Beyond the warning, any `EXPORT_SYMBOL_GPL` symbol your module tries to use will fail to resolve at load time:

```
hello: Unknown symbol __alloc_pages (err -13)  
insmod: ERROR: could not insert module hello.ko: Unknown symbol in module
```

Mismatched Kernel Headers — Silent ABI Corruption

Building against `linux-headers-6.7.x` and loading into a `6.8.x` kernel (or vice versa) will fail with the `vermagic` check in most cases. But if someone bypasses this check (with `insmod --force`), they invite silent struct layout mismatches. Always use `uname -r` in your Makefile to ensure you build against the right headers.

`printk` Without Log Level — Invisible Messages

```
 printk("hello: this might be invisible\n"); // BAD: no log level
```

C

Without a log level prefix, the message defaults to `KERN_DEFAULT`, which maps to the kernel's default message level. If the current console log level is lower than the message level, the message won't appear on the console (though it's still in the ring buffer). Always specify a log level.

Module Parameter Permissions `0666` — Security Vulnerability

```
module_param(buffer_size, int, 0666); // BAD: world-writable
```

C

This allows any unprivileged user to `echo 0 > /sys/module/hello/parameters/buffer_size` and potentially put the module into an invalid state. Use `0644` (owner write, world read) or `0444` (world read, no write) for most parameters.

Returning Positive Value From `module_init` — Undefined Behavior

```
static int __init hello_init(void)  
{  
    return 1; // BAD: positive return from init  
}
```

C

The convention is strict: `0` for success, negative `errno` code for failure. Positive values are not a valid success indicator and will confuse the kernel's module loader.

Failing to Compile With `-Werror`

The kernel's Kbuild system can be configured to treat warnings as errors. If you're developing on a system where this isn't the default, enable it explicitly to catch issues early. In your Makefile:

```
ccflags-y := -Werror
```

MAKEFILE

A warning-free build is a requirement for any kernel code that might be submitted upstream, and it's a good discipline even for local modules. Warnings in kernel code are frequently indicators of real bugs, not style issues.

Three-Level View: What Happens During `insmod`

When you run `sudo insmod hello.ko`, more happens than just "the kernel loads your module." Here's the full picture: **Level 1 — User Command:** `insmod` reads `hello.ko` into memory, calls the `init_module()` system call (or the newer `finit_module()` with a file descriptor), passing the raw ELF data to the kernel.

Level 2 — Kernel Module Loader: The kernel's `load_module()` function (in `kernel/module/main.c`) performs these steps:

1. Validates the ELF header and sections
 2. Allocates executable memory for the module using the module allocator (backed by `vmalloc`-range memory on most architectures)
 3. Copies code and data sections to the allocated memory
 4. Applies relocations — fills in the addresses of kernel symbols your module references (e.g., resolves `printk` to its actual kernel address)
 5. Verifies `vermagic` and (if `CONFIG_MODVERSIONS`) CRC checksums for each symbol
 6. Checks that `MODULE_LICENSE` doesn't taint for GPL-only symbols needed
 7. Calls your `module_init` function in the context of the `insmod` process
 8. If init returns 0: adds module to the kernel's module list, makes it visible in `/proc/modules` and `lsmod`
 9. If init returns non-zero: frees the module memory, propagates error to `insmod`
- Level 3 — Hardware:** The module code and data land in the **kernel's virtual address space** — on x86_64, this is the upper half of the 64-bit address space (above `0xfffff80000000000` approximately, depending on KASLR). The CPU's page tables are updated to map these new pages as **kernel-mode executable** — setting the `U/S` bit to 0 (supervisor-only) and `NX` bit to 0 (executable). Any attempt by userspace code to access these addresses triggers a page fault and kills the process. Any bug in your module code that causes a fault in these pages triggers a kernel oops.

Knowledge Cascade: What You've Just Unlocked

Completing this milestone doesn't just teach you how to write "Hello World" in the kernel. It unlocks a constellation of related concepts: → **ELF Section Mechanics (same domain):** The `__init`, `__exit`, and `.modinfo` annotations you used are instances of a general ELF mechanism — placing data in named sections. The linker script that produces `vmlinux` (the kernel binary) uses the same mechanism to organize init callbacks, exception tables, and system call tables. When you later encounter

`__attribute__((section(".data.once"))))` or similar in kernel source, you'll recognize the pattern. → **Symbol Export and GPL Enforcement (cross-domain: licensing meets ABI):** The fact that `MODULE_LICENSE("GPL")`

technically controls API access — not just as a legal statement but as a runtime enforcement mechanism — is a fascinating intersection of intellectual property policy and systems design. The kernel uses

`EXPORT_SYMBOL` vs `EXPORT_SYMBOL_GPL` as a technical implementation of the GPL's copyleft requirement. This is the same design philosophy as "capabilities" in security: enforcement at the technical layer, not just the policy layer. → sysfs as Universal Control Plane (cross-domain: same pattern as Prometheus): The `/sys/module/<name>/parameters/` interface you just created is an instance of the sysfs virtual filesystem — a tree of kernel objects exposed as files. The same infrastructure exposes CPU frequency scaling under `/sys/devices/system/cpu/`, USB device properties under `/sys/bus/usb/`, and block device queue parameters under `/sys/block/<dev>/queue/`. When you look at a Kubernetes node's `/sys` tree or a Prometheus node exporter scraping hardware metrics, they're reading the same virtual files you just learned to write. The mental model — "make kernel state accessible as filesystem entries" — is one of Linux's most powerful design patterns. → printk Ring Buffer → You're Building a Character Device (forward): The `dmesg` command reads `/dev/kmsg`, which is a character device. The ring buffer architecture — a fixed-size circular buffer in kernel memory, with a reader that advances through it — is structurally identical to what you'll implement in Milestone 2. When you implement your `read` file operation and manage a `buf_pos` offset, you're implementing the same mechanism that `dmesg` uses to read the printk ring buffer. → Module Parameters → Kernel Configuration Philosophy: The sysfs-backed parameter you created is a micro-example of the kernel's general configuration philosophy: separate the policy (what values mean) from the mechanism (how they're stored and accessed). The same philosophy

appears in the kernel's `sysctl` interface (`/proc/sys/`), in device tree overlays, and in runtime kernel configuration via `kconfig`. Knowing one, you understand the pattern behind all of them.

Acceptance Criteria Checkpoint

Before moving to Milestone 2, verify:

- `make` completes without warnings using Kbuild against the running kernel headers
- `modinfo hello.ko` shows `license: GPL`, `author:` , `description:` , and `parm: buffer_size`
- `sudo insmod hello.ko` succeeds and `dmesg | tail` shows the init message with `KERN_INFO`
- `lsmod | grep hello` shows the module as loaded
- `/sys/module/hello/parameters/buffer_size` exists and reads `4096`
- `sudo insmod hello.ko buffer_size=8192` produces `8192` in the sysfs attribute
- `sudo insmod hello.ko buffer_size=-1` returns an error (init validates and returns `-EINVAL`)
- `sudo rmmod hello` succeeds and `dmesg | tail` shows the exit message
- `make` with `ccflags-y := -Werror` produces zero warnings If any of these fail, use `dmesg -w` in a separate terminal while running `insmod` / `rmmod` to see kernel messages in real time. If `insmod` fails with "Operation not permitted", you need `sudo`. If it fails with "File exists", the module is already loaded — run `rmmod` first.

Milestone 2: Character Device Driver

The Revelation: User Pointers Are Traps

You've just loaded a kernel module. Now a process is going to call `write(fd, buf, 4096)` and your driver needs to receive those 4096 bytes. The user hands you a pointer. Why can't you just do this?

```

/* THIS IS WRONG. DO NOT DO THIS. */

static ssize_t mydev_write(struct file *filp, const char __user *buf,
                         size_t count, loff_t *f_pos)

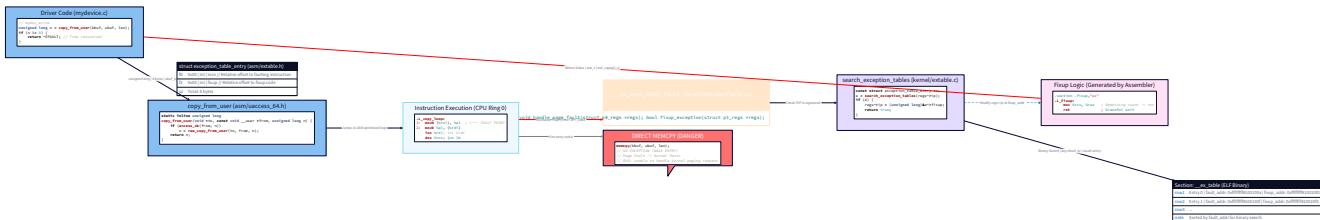
{
    memcpy(kernel_buffer, buf, count); /* ← WRONG: direct dereference of user pointer */

    return count;
}

```

Your instinct says: "The kernel can access *all* memory — it's in ring 0, it has the highest privilege. If the user gives me a pointer to their buffer at address `0x7fff1234`, I can just read from `0x7fff1234`." This instinct is half right and fully dangerous. The kernel *can* access that address in a mechanical sense — the CPU won't fault on privilege. But consider what happens when:

1. The user passes a **completely bogus pointer**: `write(fd, (void*)0xdeadbeef, 100)`. There's no page mapped at `0xdeadbeef` in the user's page tables. A direct `memcpy` from `0xdeadbeef` generates a page fault. In userspace, a page fault in process context just delivers SIGSEGV. But this page fault happens *in kernel context* — during your driver's write handler — and it kills the kernel.
2. The user passes a **pointer to a memory-mapped device register** and your driver blindly reads from it. Now you've let a userspace program trigger arbitrary hardware I/O through your driver.
3. The user passes a **valid pointer at the moment of the call**, but the page gets swapped out between your check and your copy. Classic TOCTOU — time-of-check/time-of-use race condition.
4. A 32-bit process running on a 64-bit kernel passes a pointer that looks valid but falls in the kernel's own address space. Direct dereference now reads kernel memory on behalf of a userspace program. The solution is `copy_from_user()` and `copy_to_user()`. But these functions don't solve the problem by *checking* the pointer first and then copying. They solve it by **installing a temporary exception handler before the copy**. Here's the real mechanism:



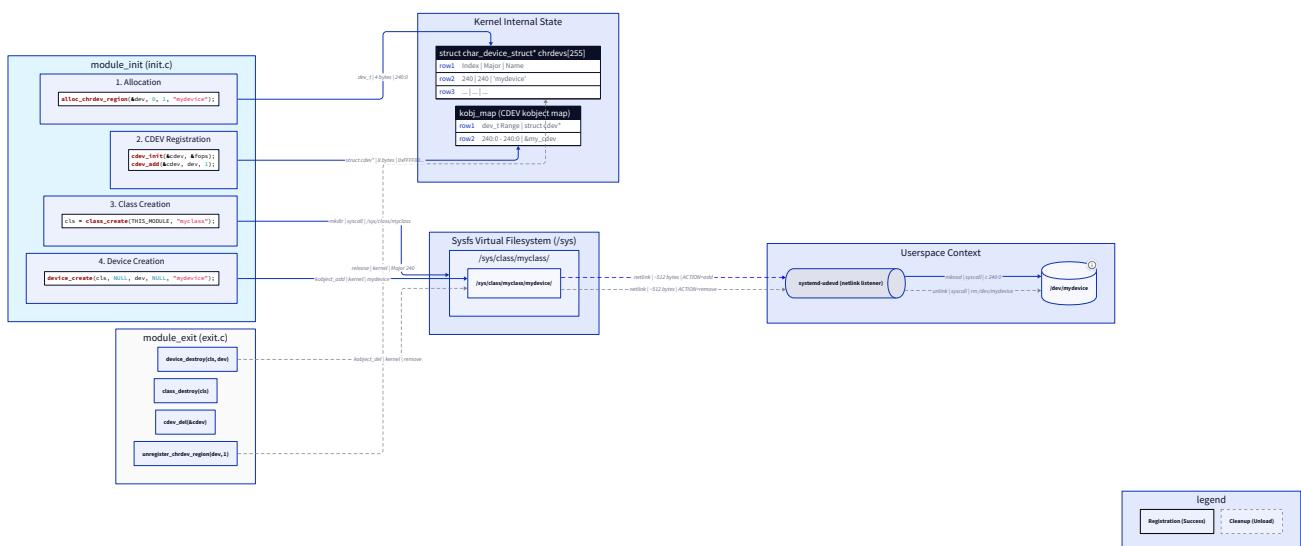
The kernel maintains an **exception table** — a sorted array of `(faulting_address, fixup_address)` pairs in a `.ex_table` ELF section. The assembly code inside `copy_from_user()` is annotated so that every load instruction has a corresponding entry in this table. When a page fault occurs anywhere in the kernel:

1. The page fault handler checks: *is this fault address in the exception table?*
2. If yes: instead of panicking, it jumps to the `fixup_address` — a small stub that sets the return value to indicate failure, zeroes the destination buffer, and returns gracefully.
3. If no: it's an unexpected kernel bug — oops/panic. This means `copy_from_user()` is a **trap-and-recover** mechanism, not a check-then-copy mechanism. The copy proceeds optimistically. If it faults, the kernel recovers cleanly and returns `-EFAULT` to your driver, which returns `-EFAULT` to the user process. The kernel never crashes. The user process gets an error code and continues (or not — their choice).

This is one of the most important pieces of kernel API design you'll encounter. Every single data transfer path in your driver — now and forever — goes through `copy_from_user` or `copy_to_user`. There are no exceptions. Not for "trusted" programs. Not for root processes. Not for "small" copies. The exception table mechanism exists to guarantee that even hostile or buggy userspace cannot crash the kernel through your driver.

The Architecture Before the Code

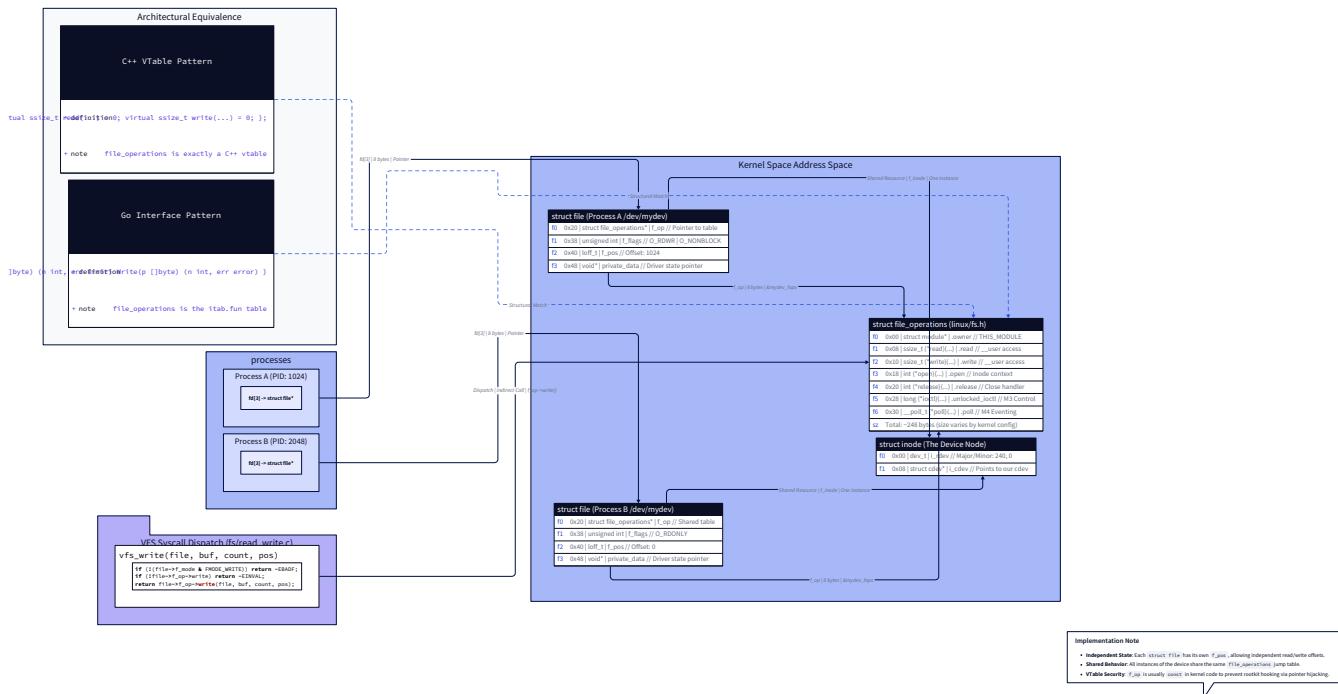
Before writing a single line, you need to see the full picture of what you're building.



When a userspace program opens `/dev/mydevice`, a chain of kernel machinery fires:

1. The VFS (Virtual File System — the kernel's unified interface over all filesystems and device files) looks at the device file's inode.

2. The inode contains a **device number** (`dev_t`) — a 32-bit integer encoding two values: the **major number** (identifies the driver) and the **minor number** (identifies which instance of that driver).
3. The kernel looks up the major number in a character device table, finds your `cdev` structure, and calls your `.open` function pointer.
4. From there, every `read()`, `write()`, `ioctl()`, and `close()` on that file descriptor routes through the function pointers in your `file_operations` struct. [[EXPLAIN:major/minor-numbers-and-the-vfs-layer|Major/minor numbers and the VFS layer]] This dispatch chain — from system call number to your function — is the VFS vtable pattern:



🔑 Foundation: The `file_operations` vtable pattern

What it IS The `struct file_operations` is a dispatch table (or Virtual Method Table) used by the Linux kernel to implement polymorphism in C. It is a structure filled with function pointers—such as `.read`, `.write`, `.open`, and `.release`—that link generic system calls to your specific driver's implementation. When a userspace program calls `read()` on a file descriptor, the kernel looks up the `file_operations` struct associated with that file and executes the function pointed to by the `.read` member.

WHY you need it right now In a character device driver, this is your primary interface. The kernel treats "everything as a file," but the kernel doesn't know what "reading" means for your specific hardware or virtual device. You must fill this struct to map standard POSIX API calls to your custom logic. Without it, your device exists in the system but remains unreachable and non-functional to userspace.

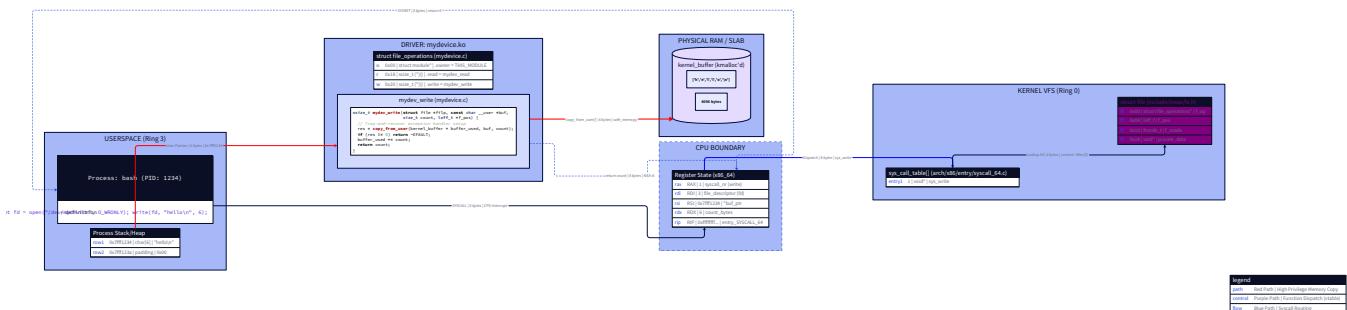
ONE key insight or mental model Think of the `file_operations` struct as a **Contract**: The kernel provides the "Interface" (the names of the functions), and you provide the "Implementation" (the actual logic).

It decouples the VFS (Virtual File System) from the hardware-specific details.

The `struct file_operations` is the kernel's equivalent of a C++ vtable or a Go interface: a struct of function pointers where each pointer implements one operation. Your driver fills in the ones it cares about; the ones left `NULL` get default behavior (or return `-ENOSYS` for "not implemented").

How `write(2)` Becomes Your Function

Let's trace the full path from a userspace `write()` call to your driver's write handler. This path demystifies why the function signatures look the way they do.





The function signature is dictated by this call chain. Your write handler receives:

- `struct file *filp` — the kernel's representation of the open file; carries state like `f_pos`, `f_flags`, and a pointer to your device-private data
- `const char __user *buf` — the userspace buffer address (annotated with `__user` to help Sparse, the kernel's static analysis tool, catch direct dereferences)
- `size_t count` — bytes requested
- `loff_t *f_pos` — pointer to the file position; you update this in read to advance the position. The `__user` annotation is parsed by **Sparse** (the kernel's static checker, invoked with `make C=1`). It doesn't affect compiled code, but Sparse will warn if you pass a `__user`-annotated pointer to a function expecting a kernel pointer (like `memcpy`). This turns potential runtime crashes into compile-time warnings.

Device Numbers: Dynamic Allocation

In the early days of Unix, device major numbers were statically assigned. The kernel source file `Documentation/admin-guide/devices.txt` lists these historical allocations — major 1 is `mem` (for `/dev/null`, `/dev/zero`, etc.), major 4 is `tty`, major 8 is SCSI disk. There are only 4096 major numbers,

and they were running out. Modern practice: **always use dynamic allocation**. The kernel picks an unused major number at load time:

```
static dev_t dev_num;           /* will hold our major+minor after allocation */

static struct cdev my_cdev;    /* kernel's char device structure */

static struct class *my_class; /* for automatic /dev/ node creation */

/* In module_init: */

ret = alloc_chrdev_region(&dev_num, 0, 1, "mydevice");

/*
 *          ^          ^  ^      ^
 *          |          |  |      device name in /proc/devices
 *          |          |  |      number of minors to reserve
 *          |          |      first minor number
 *
 *          output: dev_t with major+minor
 */

if (ret < 0) {

    printk(KERN_ERR "mydev: failed to allocate device number: %d\n", ret);

    return ret;
}

printk(KERN_INFO "mydev: allocated major=%d, minor=%d\n",
       MAJOR(dev_num), MINOR(dev_num));
```

`MAJOR(dev_num)` and `MINOR(dev_num)` are macros that extract the respective fields from the 32-bit `dev_t`. The encoding on Linux: top 12 bits are major, bottom 20 bits are minor. `MKDEV(major, minor)` constructs a `dev_t` from components. After this call succeeds, you've reserved a device number. The device doesn't exist yet — there's no `cdev` connected, no `/dev/` node. You've just claimed the number. You can verify it:

```
grep mydevice /proc/devices

# Output: 240 mydevice

# (240 is an example; actual number varies)
```

Critical cleanup rule: every successful `alloc_chrdev_region()` must be paired with `unregister_chrdev_region()` in your exit path and error paths. The kernel does *not* automatically reclaim device numbers when a module exits without cleanup.

The `cdev` Structure: Connecting Numbers to Operations

Allocating a number isn't enough. You need to tell the kernel: "when something opens the device with this number, use *these* file operations." That's the job of `struct cdev`:

```
static const struct file_operations mydev_fops = {  
    .owner    = THIS_MODULE,  
    .open     = mydev_open,  
    .release  = mydev_release,  
    .read     = mydev_read,  
    .write    = mydev_write,  
};  
  
/* In module_init, after alloc_chrdev_region: */  
  
cdev_init(&my_cdev, &mydev_fops);  
  
my_cdev.owner = THIS_MODULE;  
  
ret = cdev_add(&my_cdev, dev_num, 1);  
  
/*  
 *          ^          ^          ^  
 *          cdev        dev_t      number of minors this cdev handles */  
  
if (ret < 0) {  
    printk(KERN_ERR "mydev: cdev_add failed: %d\n", ret);  
    goto err_cdev;  
}
```

`cdev_add()` makes the device live — from this point, if anything opens the device number, the kernel will dispatch to your `mydev_fops`. This is why error handling order matters: you don't want the device reachable before it's fully initialized. The `.owner = THIS_MODULE` field prevents the module from being unloaded while the device has open file descriptors. `THIS_MODULE` is a pointer to the `struct module` that represents your module. The VFS increments the module's reference count when a file is opened and decrements it on close. If you omit `.owner`, `rmmmod` can unload your module while a process still has an open file descriptor — the next `read()` call would jump to the now-unmapped function and crash the kernel.

Automatic `/dev/` Nodes: `class_create` and `device_create`

You've allocated a number and registered a cdev. But there's still no `/dev/mydevice` file. In ancient times you'd run `mknod /dev/mydevice c 240 0` manually. Modern Linux uses `udev` — a userspace daemon that listens for hotplug events from the kernel and automatically creates device nodes.

 **Foundation: udev and the kernel uevent mechanism: how `device_create` triggers `/dev/` node creation**

What it IS The kernel `uevent` mechanism is a notification system where the kernel broadcasts "user events" to userspace whenever a device state changes (e.g., a device is created, removed, or changed). `udev` is a userspace daemon that listens for these events via a netlink socket. When you call `device_create()` in your driver, the kernel sends a uevent; `udev` catches it, looks up its rules (usually in `/etc/udev/rules.d/`), and automatically creates a corresponding device node in the `/dev/` directory with the correct permissions.

WHY you need it right now Manually creating device files using the `mknod` command is error-prone and tedious. By using the `class_create()` and `device_create()` functions in your module, you trigger the uevent mechanism. This ensures that as soon as your module is loaded (`insmod`), the device file (e.g., `/dev/mydevice`) appears instantly and automatically for the user, and disappears when the module is unloaded.

ONE key insight or mental model The kernel is the **Broadcaster** and `udev` is the **Assistant**. The kernel shouts "I've added a new device!" into a megaphone, and `udev` hears it and quickly runs to the `/dev/`

folder to put a new "doorway" (device node) there for userspace to use.

The kernel sends udev events through the **device model** — a hierarchical tree of `struct device` objects that represents all hardware and virtual devices in the system. `class_create()` and `device_create()` are your entry points into this system:

```
/* Create a device class – appears as /sys/class/<name> */

my_class = class_create(THIS_MODULE, "mydevice");

if (IS_ERR(my_class)) {

    ret = PTR_ERR(my_class);

    printk(KERN_ERR "mydev: class_create failed: %d\n", ret);

    goto err_class;

}

/* Create the device – triggers udev to create /dev/mydevice */

my_device = device_create(my_class, NULL, dev_num, NULL, "mydevice");

if (IS_ERR(my_device)) {

    ret = PTR_ERR(my_device);

    printk(KERN_ERR "mydev: device_create failed: %d\n", ret);

    goto err_device;

}
```

C

IS_ERR() and **PTR_ERR()** — **kernel error pointers**: The kernel uses a trick where error codes are encoded as pointer values in the high portion of the address space. `IS_ERR(ptr)` checks if the pointer is actually an error code (a value between `-MAX_ERRNO` and 0, cast to a pointer). `PTR_ERR(ptr)` extracts the integer error code. This allows functions that return pointers to signal errors without a separate output parameter. You'll see this pattern everywhere in kernel code for functions that return `struct * . device_create()` registers a `struct device` in the kernel's device tree and emits a **uevent** with the device's class, major/minor numbers, and name. udev (or systemd-udevd) receives this event and creates `/dev/mydevice` with the correct major/minor and permissions. You can watch this happen:

```
# In one terminal, watch udev events:  
  
udevadm monitor --kernel --udev  
  
# In another terminal, load your module:  
  
sudo insmod mydevice.ko  
  
# You'll see: add /class/mydevice/mydevice (class)  
  
# And udev creates /dev/mydevice automatically
```

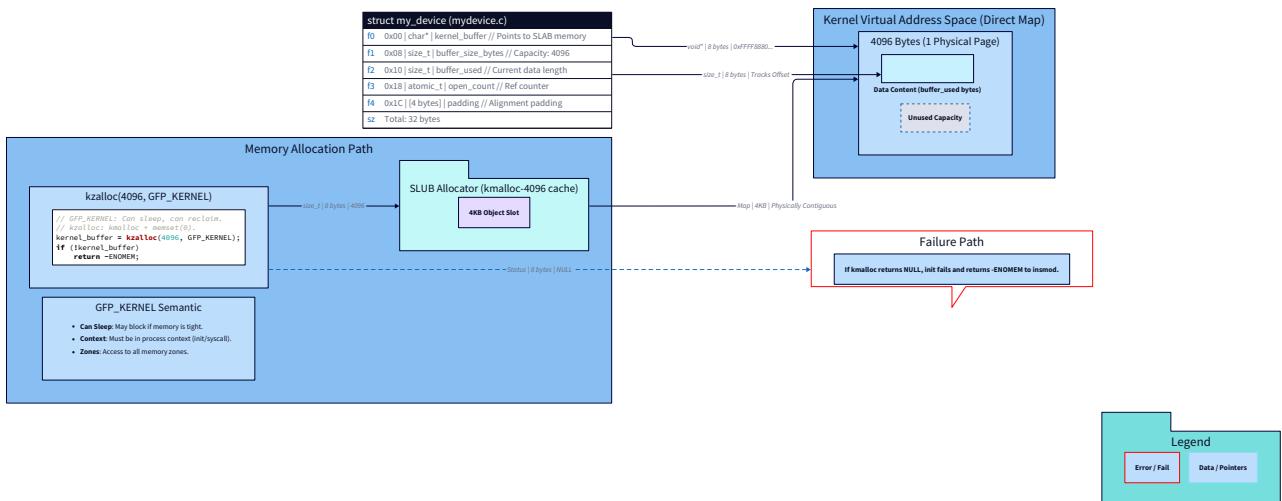
BASH

The Kernel Buffer: kmalloc and the SLAB Allocator

Your device needs a place to store data. In userspace you'd call `malloc()`. In the kernel you call `kmalloc()`:

```
static char *kernel_buffer;  
  
static size_t buffer_size_bytes = 4096;  
  
/* In module_init: */  
  
kernel_buffer = kzalloc(buffer_size_bytes, GFP_KERNEL);  
  
/*  
 * ^ zero-fills the allocation |  
 * (k + zalloc = malloc+memset(0))  
 * allocation size  
 * allocation flags */  
  
if (!kernel_buffer) {  
  
    printk(KERN_ERR "mydev: failed to allocate kernel buffer\n");  
  
    ret = -ENOMEM;  
  
    goto err_buffer;  
}
```

C



`kmalloc` draws from the **SLAB/SLUB allocator** — not a heap in the traditional sense.

🔑 Foundation: `kmalloc` and the SLAB/SLUB allocator: how kernel memory allocation differs from `malloc`

What it IS `kmalloc` is the kernel equivalent of userspace `malloc()`, used for allocating physically contiguous memory. Under the hood, it doesn't just grab raw pages; it uses the **SLAB** (or the more modern **SLUB**) allocator. This system manages "caches" of commonly used object sizes. Instead of constantly searching for free memory of arbitrary sizes, the SLUB allocator maintains lists of pre-allocated "slabs" (blocks of memory) divided into small, equal-sized "slots."

WHY you need it right now Kernel memory is a finite, precious resource that cannot be "swapped" to disk like userspace memory. When you need to store driver state or buffers, you use `kmalloc`. However, you must specify a **GFP (Get Free Page) flag**: usually `GFP_KERNEL` (which can sleep if memory is tight) or `GFP_ATOMIC` (which cannot sleep, used inside interrupt handlers). Understanding that `kmalloc` is backed by the SLUB allocator helps you realize why it is highly efficient for small, frequent allocations but can fail if you request massive, contiguous chunks.

ONE key insight or mental model The SLUB allocator is like a **Bento Box shop**. Instead of cooking every meal from scratch (searching raw memory), they have pre-portioned boxes for 32-byte, 64-byte, and 128-byte meals. When you call `kmalloc(40)`, the kernel just hands you a 64-byte "box" because it's the fastest way to serve you.

The `GFP_KERNEL` flag is the most important argument to understand right now. GFP stands for "Get Free Pages" — it describes the allocation context. `GFP_KERNEL` means:

- **Can sleep:** the allocator may block the current process if memory is scarce, waiting for pages to be freed or swapped out
- **Can reclaim:** the kernel may try to free other caches to satisfy this allocation

- **Can use all memory zones:** the allocator has access to all physical memory. This is correct for your `module_init` context — init runs in process context (the `insmod` process), so sleeping is fine. **The critical rule:** if you're ever in **interrupt context** (inside an interrupt handler) or holding a **spinlock** (a non-sleeping lock), you must use `GFP_ATOMIC` instead of `GFP_KERNEL`. `GFP_ATOMIC` cannot sleep, cannot reclaim, and may fail more often. Using `GFP_KERNEL` in atomic context causes the kernel to print a warning and potentially deadlock. For this milestone, you're only allocating in `module_init` where `GFP_KERNEL` is always correct.

Hardware Soul — kmalloc memory placement:

- `kmalloc` returns **physically contiguous** memory (pages from the buddy allocator, managed by SLAB/SLUB)
- This memory is in the **kernel's direct mapping** (on x86_64: the `0xffff888000000000` region, directly mapped to physical memory). This means kernel virtual address and physical address differ by a fixed offset — `virt_to_phys()` computes this.
- Your 4KB buffer fits in one page. For the L1 cache (typically 32KB, 8-way set-associative, 64-byte cache lines), a 4KB buffer occupies 64 cache lines and fits comfortably. Sequential access by your write/read handlers will be cache-friendly — hardware prefetcher will pull ahead as you access bytes in order.
- `kzalloc` (zero-filling variant) is marginally slower than `kmalloc` for the first use but ensures no data leakage from previously freed memory — a security property.

Implementing `file_operations`: Open and Release

Now for the actual file operations. Start with the simplest pair — open and release:

```

static atomic_t open_count = ATOMIC_INIT(0);

static int mydev_open(struct inode *inode, struct file *filp)
{
    atomic_inc(&open_count);

    printk(KERN_INFO "mydev: opened (count=%d)\n",
           atomic_read(&open_count));

    return 0;
}

static int mydev_release(struct inode *inode, struct file *filp)
{
    atomic_dec(&open_count);

    printk(KERN_INFO "mydev: released (count=%d)\n",
           atomic_read(&open_count));

    return 0;
}

```

Why `atomic_t` instead of a plain `int`? Because `open` and `release` can be called concurrently from multiple processes. Incrementing a plain integer is a read-modify-write operation — three machine instructions. Two CPUs doing it simultaneously can read the same value, both increment it to the same result, and write the same value back (both get `count+1` when the correct result is `count+2`). `atomic_inc` generates an atomic increment instruction (`lock xadd` on x86) that the CPU's bus lock makes indivisible.

🔑 Foundation: Linux atomic operations: `atomic_t`

What it IS Atomic operations are low-level instructions that are guaranteed to execute as a single, indivisible unit. The Linux kernel provides the `atomic_t` type (an opaque structure wrapping an integer) and functions like `atomic_inc()`, `atomic_dec()`, and `atomic_read()`. These use CPU-specific instructions (like `LOCK` prefixes on x86) to ensure that if two CPU cores try to increment the same variable simultaneously, one will strictly happen after the other, and no data will be lost.

WHY you need it right now In a multi-core kernel environment, a simple `i++` is dangerous. It involves three steps: read from memory, increment in register, write back to memory. If an interrupt occurs or another CPU accesses the variable mid-way, you get a "race condition" and corrupted data. In your driver, if you are

counting how many processes have opened your device or tracking a shared resource, you must use `atomic_t` to ensure thread safety without the heavy performance overhead of a Mutex or Spinlock.

ONE key insight or mental model An atomic operation is **Indivisible**. To the rest of the system, the variable jumps from `A` to `B` instantly; there is no "middle" state where the variable is being processed, making it immune to being interrupted by other tasks.

struct inode vs struct file : The inode represents the device file itself — it's shared between all processes and persists as long as the file exists. The `struct file` represents a single *open instance* — it's created fresh for each `open()` call and destroyed on the last `close()`. If three processes open `/dev/mydevice` simultaneously, they share one inode but have three distinct `struct file` instances. This is where `filp->f_pos` lives — per-file-descriptor, not per-device. **Private data pattern**: For more complex drivers, you'd store a pointer to your device-private struct in `filp->private_data` during open:

```
static int mydev_open(struct inode *inode, struct file *filp) C

{
    filp->private_data = &my_device_state; /* accessible in read/write/ioctl */

    atomic_inc(&open_count);

    return 0;
}
```

This is crucial for drivers that manage multiple device instances — each minor number maps to a different device state, and you need to

find the right one in your handlers.

Implementing Write: `copy_from_user` in Practice

```
static size_t buffer_used = 0; /* bytes currently in the buffer */

static ssize_t mydev_write(struct file *filp, const char __user *buf,
                           size_t count, loff_t *f_pos)

{
    size_t space_available;

    size_t bytes_to_copy;

    unsigned long not_copied;

    /* How much space remains in the buffer? */

    space_available = buffer_size_bytes - buffer_used;

    if (space_available == 0)

        return -ENOSPC; /* buffer full */

    /* Don't copy more than we have space for */

    bytes_to_copy = min(count, space_available);

    /*
     * copy_from_user(to, from, n):
     *
     *   to    = kernel destination address
     *
     *   from = userspace source address (__user pointer)
     *
     *   n    = bytes to copy
     *
     * Returns: number of bytes NOT copied (0 on success).
     *
     * Non-zero means a fault occurred; return -EFAULT.
     */
    not_copied = copy_from_user(kernel_buffer + buffer_used, buf, bytes_to_copy);
```

```

if (not_copied != 0) {

    /*
     * Partial copy: some bytes got through before the fault.
     *
     * We must account for the bytes that DID make it.
     *
     * bytes_copied = bytes_to_copy - not_copied
     *
     * But for simplicity in this milestone, treat any copy failure
     * as -EFAULT. More complex drivers handle partial copies.
     */

    return -EFAULT;
}

buffer_used += bytes_to_copy;

/* Note: we do NOT update *f_pos here.

 * Write position tracking in this device uses buffer_used, not f_pos.

 * f_pos is updated by the VFS for regular files; for character devices
 * the semantics are driver-defined. */

printk(KERN_DEBUG "mydev: write %zu bytes (buffer now %zu/%zu)\n",
       bytes_to_copy, buffer_used, buffer_size_bytes);

return bytes_to_copy; /* return number of bytes accepted */
}

```

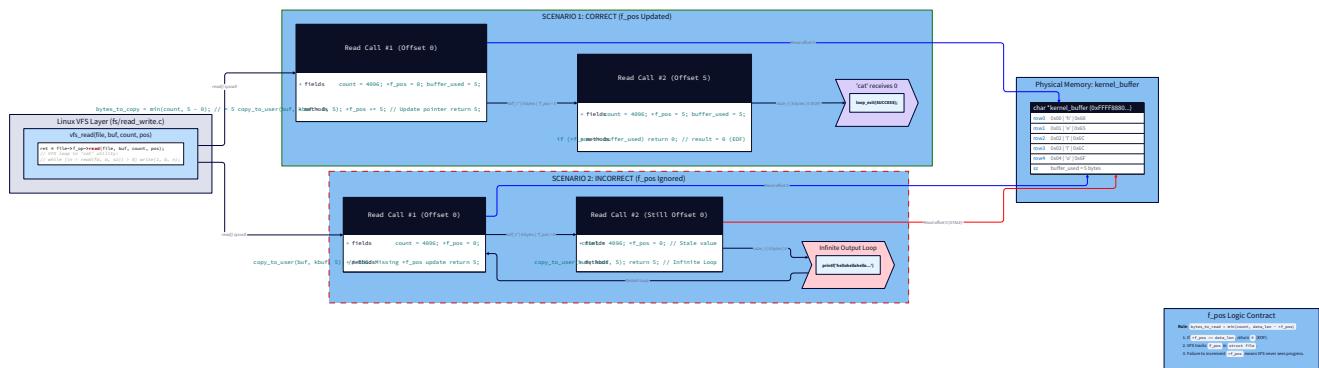
Return value semantics for write handlers:

- Positive: bytes successfully written. The VFS may call your handler again if the user requested more than you returned.
- `0`: nothing written but no error (unusual for write, but valid for non-blocking when buffer is full — we'll use `-ENOSPC` here for clarity)
- Negative: error code. `-EFAULT` for bad user pointer. `-ENOSPC` for no space. `-EINTR` for interrupted. The error propagates to the userspace `write()` return value. `min()` vs `min_t()`: The kernel defines `min(a, b)` as a type-safe macro that warns if `a` and `b` have different types. If you're comparing

values of different types (say `size_t` and `ssize_t`), use `min_t(type, a, b)` to specify the comparison type explicitly. Mismatched signed/unsigned comparisons are a source of subtle bugs (a negative `ssize_t` compared to a `size_t` compares as very large, not negative).

Implementing Read: f_pos Tracking and the EOF Contract

The read handler is where developers most often make the mistake of returning wrong values and creating infinite loops. Let's understand why this happens before writing the code.



When you run `cat /dev/mydevice`, `cat` calls `read(fd, buf, BUFSIZ)` in a loop until it gets a return value of `0` (EOF). This is the universal Unix contract: **read returns 0 to signal end of file, not an error**. If your read handler always returns a positive number (even 1 byte), `cat` will loop forever, generating an infinite stream of whatever data you return. This is the most common read handler bug in beginner kernel modules. The `f_pos` (file position) tracks how far into the "file" the reader has consumed. For your character device, it's tracking how far into `kernel_buffer` the current reader has gotten:

```
static ssize_t mydev_read(struct file *filp, char __user *buf,
                         size_t count, loff_t *f_pos)

{

    size_t bytes_available;

    size_t bytes_to_copy;

    unsigned long not_copied;

    /*

     * *f_pos is the current read position within the buffer.

     * If f_pos >= buffer_used, the reader has consumed all data: EOF.

     */

    if (*f_pos >= buffer_used)

        return 0; /* EOF: no more data - cat and read() loops terminate here */

    /* How many unread bytes remain? */

    bytes_available = buffer_used - *f_pos;

    /* Give the caller at most 'count' bytes */

    bytes_to_copy = min(count, bytes_available);

    /*

     * copy_to_user(to, from, n):

     *   to    = userspace destination address (__user pointer)

     *   from = kernel source address

     *   n    = bytes to copy

     *

     * Returns: number of bytes NOT copied (0 on full success).

     */

    not_copied = copy_to_user(buf, kernel_buffer + *f_pos, bytes_to_copy);

    if (not_copied != 0)

        return -EFAULT;
```

```

/* Advance the file position by the number of bytes we copied */

*f_pos += bytes_to_copy;

printf(KERN_DEBUG "mydev: read %zu bytes (f_pos now %lld/%zu)\n",
       bytes_to_copy, *f_pos, buffer_used);

return bytes_to_copy;

}

```

The f_pos dance:

- `*f_pos` starts at 0 when the file is first opened
- Each successful read advances `*f_pos` by the number of bytes returned
- When `*f_pos == buffer_used`, the next read returns 0 (EOF)
- After EOF, `cat` exits its read loop **Why `f_pos` is a pointer parameter**: The kernel stores `f_pos` in `filp->f_pos`. But for `pread()` (positional read), the VFS passes a *different* offset than `filp->f_pos` — it passes the caller's requested offset directly, without modifying the file's stored position. By taking a `loff_t *f_pos` parameter, your handler works correctly for both `read()` and `pread()` without modification. **The "two-process independence" property**: If process A and process B both open `/dev/mydevice`, they get separate `struct file` instances, each with its own `f_pos`. Process A reading doesn't advance process B's position. This is the correct semantics — it mirrors how two processes reading the same regular file maintain independent positions.

The Complete Device: Putting It All Together

Now let's write the full, complete driver. This is the production-quality reference you'll build and verify:

```
/* mydevice.c - Character device driver with read/write, f_pos tracking */
```

C

```
#include <linux/module.h>
```

```
#include <linux/kernel.h>
```

```
#include <linux/init.h>
```

```
#include <linux/fs.h>           /* file_operations, alloc_chrdev_region */
```

```
#include <linux/cdev.h>         /* cdev_init, cdev_add */
```

```
#include <linux/device.h>        /* class_create, device_create */
```

```
#include <linux/uaccess.h>       /* copy_to_user, copy_from_user */
```

```
#include <linux/slab.h>          /* kmalloc, kfree */
```

```
#include <linux/atomic.h>         /* atomic_t */
```

```
MODULE_LICENSE("GPL");
```

```
MODULE_AUTHOR("Your Name <you@example.com>");
```

```
MODULE_DESCRIPTION("Character device driver: read/write with kernel buffer");
```

```
MODULE_VERSION("1.0");
```

```
/* — Configuration —————— */
```

```
#define DEVICE_NAME      "mydevice"
```

```
#define CLASS_NAME       "mydevice_class"
```

```
#define BUFFER_SIZE      4096
```

```
/* — Module state —————— */
```

```
static dev_t          dev_num;           /* major + minor, set by alloc_chrdev_region */
```

```
static struct cdev    my_cdev;          /* kernel char device structure */
```

```
static struct class   *my_class;        /* device class for udev */
```

```
static struct device  *my_device;       /* device object */
```

```
static char            *kernel_buffer;  /* 4KB buffer for data storage */
```

```
static size_t           buffer_used = 0; /* bytes currently stored */
```

```
static atomic_t          open_count = ATOMIC_INIT(0);
```

```
/* — File Operations —————— */
```

```
static int mydev_open(struct inode *inode, struct file *filp)
{
    atomic_inc(&open_count);

    printk(KERN_INFO "mydev: open (count now %d)\n",
           atomic_read(&open_count));

    return 0;
}

static int mydev_release(struct inode *inode, struct file *filp)
{
    atomic_dec(&open_count);

    printk(KERN_INFO "mydev: release (count now %d)\n",
           atomic_read(&open_count));

    return 0;
}

static ssize_t mydev_read(struct file *filp, char __user *buf,
                         size_t count, loff_t *f_pos)
{
    size_t bytes_available;

    size_t bytes_to_copy;

    unsigned long not_copied;

    if (*f_pos >= buffer_used)

        return 0; /* EOF */

    bytes_available = buffer_used - *f_pos;

    bytes_to_copy = min(count, bytes_available);

    not_copied = copy_to_user(buf, kernel_buffer + *f_pos, bytes_to_copy);

    if (not_copied != 0)

        return -EFAULT;
```

```
*f_pos += bytes_to_copy;

return (ssize_t)bytes_to_copy;
}

static ssize_t mydev_write(struct file *filp, const char __user *buf,
                         size_t count, loff_t *f_pos)

{
    size_t space_available;

    size_t bytes_to_copy;

    unsigned long not_copied;

    space_available = BUFFER_SIZE - buffer_used;

    if (space_available == 0)

        return -ENOSPC;

    bytes_to_copy = min(count, space_available);

    not_copied = copy_from_user(kernel_buffer + buffer_used, buf, bytes_to_copy);

    if (not_copied != 0)

        return -EFAULT;

    buffer_used += bytes_to_copy;

    return (ssize_t)bytes_to_copy;
}

static const struct file_operations mydev_fops = {

    .owner     = THIS_MODULE,

    .open      = mydev_open,

    .release   = mydev_release,

    .read      = mydev_read,

    .write     = mydev_write,
};

/* —— Init / Exit ————— */
```

```
/*
 * Initialization uses goto-based error handling – the kernel idiom
 * for unwinding partial initialization in reverse order.
 *
 * Order of init: buffer → chrdev_region → cdev → class → device
 * Order of cleanup (exit and error paths): reverse
 */

static int __init mydev_init(void)
{
    int ret;

    /* 1. Allocate kernel buffer */

    kernel_buffer = kzalloc(BUFFER_SIZE, GFP_KERNEL);

    if (!kernel_buffer) {
        printk(KERN_ERR "mydev: failed to allocate buffer\n");

        return -ENOMEM;
    }

    /* 2. Reserve a character device number dynamically */

    ret = alloc_chrdev_region(&dev_num, 0, 1, DEVICE_NAME);

    if (ret < 0) {
        printk(KERN_ERR "mydev: alloc_chrdev_region failed: %d\n", ret);

        goto err_alloc_region;
    }

    printk(KERN_INFO "mydev: registered major=%d minor=%d\n",
           MAJOR(dev_num), MINOR(dev_num));

    /* 3. Initialize and add the cdev */

    cdev_init(&my_cdev, &mydev_fops);

    my_cdev.owner = THIS_MODULE;
```

```
ret = cdev_add(&my_cdev, dev_num, 1);

if (ret < 0) {

    printk(KERN_ERR "mydev: cdev_add failed: %d\n", ret);

    goto err_cdev_add;

}

/* 4. Create device class (appears in /sys/class/) */

my_class = class_create(THIS_MODULE, CLASS_NAME);

if (IS_ERR(my_class)) {

    ret = PTR_ERR(my_class);

    printk(KERN_ERR "mydev: class_create failed: %d\n", ret);

    goto err_class;

}

/* 5. Create device – triggers udev to create /dev/mydevice */

my_device = device_create(my_class, NULL, dev_num, NULL, DEVICE_NAME);

if (IS_ERR(my_device)) {

    ret = PTR_ERR(my_device);

    printk(KERN_ERR "mydev: device_create failed: %d\n", ret);

    goto err_device;

}

printk(KERN_INFO "mydev: initialized – /dev/%s created\n", DEVICE_NAME);

return 0;

/* Error unwind – reverse order of init */

err_device:

class_destroy(my_class);

err_class:

cdev_del(&my_cdev);

err_cdev_add:
```

```

    unregister_chrdev_region(dev_num, 1);

err_alloc_region:

    kfree(kernel_buffer);

    return ret;

}

static void __exit mydev_exit(void)

{

    device_destroy(my_class, dev_num);

    class_destroy(my_class);

    cdev_del(&my_cdev);

    unregister_chrdev_region(dev_num, 1);

    kfree(kernel_buffer);

    printk(KERN_INFO "mydev: exited cleanly\n");

}

module_init(mydev_init);

module_exit(mydev_exit);

```

The Goto Error Handling Pattern

Notice the `goto` labels in `mydev_init`. This is not sloppy code — it is the canonical kernel idiom for error handling with multiple resources. The structure ensures that if step N fails, steps 1 through N-1 are cleaned up in reverse order. Without this, a failure in `class_create` would leak the cdev and device number. The alternative — nested if-else chains — produces deeply indented, hard-to-read code. The goto pattern keeps the success path reading top-to-bottom, with the cleanup code separated at the bottom. Every experienced kernel developer uses this

pattern. If you see a kernel function without it, that function probably has a resource leak.

The Kbuild Makefile for This Module

```
# Makefile                                         MAKEFILE
obj-m += mydevice.o
ccflags-y := -Werror
all:
    make -C /lib/modules/$(shell uname -r)/build M=$(PWD) modules
clean:
    make -C /lib/modules/$(shell uname -r)/build M=$(PWD) clean
```

Verification: The echo/cat Round Trip

Build, load, and verify:

```

# Build

make

# Load

sudo insmod mydevice.ko

# Verify device number and class

grep mydevice /proc/devices

# Example: 240 mydevice

ls -la /dev/mydevice

# crw----- 1 root root 240, 0 Mar 1 12:00 /dev/mydevice

# Write data (echo appends a newline, so "hello\n" = 6 bytes)

echo "hello" | sudo tee /dev/mydevice

# Output: hello

# Read data back

sudo cat /dev/mydevice

# Output: hello

# Verify open_count via dmesg

dmesg | tail -10

# [ ...] mydev: open (count now 1)      ← from tee's open

# [ ...] mydev: release (count now 0)    ← from tee's close

# [ ...] mydev: open (count now 1)      ← from cat's open

# [ ...] mydev: release (count now 0)    ← from cat's close

# Unload

sudo rmmod mydevice

```

Why does `cat` terminate? When `cat` opens the device, `*f_pos` starts at 0. The first `read()` call returns "hello\n" (6 bytes) and advances `*f_pos` to 6. Since `buffer_used` is also 6, the second `read()` call hits `*f_pos >= buffer_used`

and returns 0 (EOF). `cat` sees EOF and exits. Why does a second `cat` read nothing? After the first `cat` finishes, `buffer_used` is still 6, but `*f_pos` for the new open is 0. So a second `cat` would read "hello\n" again. The buffer persists until overwritten or the module is unloaded. This is the intended behavior for a simple device.

Three-Level View: What Happens During `echo "hello" > /dev/mydevice`

Level 1 — Shell/Userspace: The shell forks, execs `echo`, which writes "`hello\n`" (6 bytes) to stdout. The shell has already redirected stdout to `/dev/mydevice` via `open("/dev/mydevice", O_WRONLY)`. The `write(fd, "hello\n", 6)` syscall fires. **Level 2 — Kernel/VFS:** `sys_write` → `ksys_write` → `vfs_write`. VFS looks up the `struct file` for `fd`, retrieves `f->f_op->write`, which is your `mydev_write`. Your handler runs in the context of the `echo` process. `copy_from_user` copies the 6 bytes from the `echo` process's stack (where "`hello\n`" lives) into `kernel_buffer`. The exception table ensures that if `echo`'s mapping disappears mid-copy, the kernel recovers with `-EFAULT`. **Level 3 — Hardware:** The copy involves two TLB lookups: one for the userspace source address (in `echo`'s page tables) and one for the kernel destination address (in the direct-mapped region). The CPU's hardware page walker traverses the page table hierarchy for the user address. The 6 bytes fit in a single 64-byte cache line. The copy itself is handled by optimized CPU instructions in the `copy_from_user` implementation — on x86_64, this uses `rep movsb` or SIMD instructions depending on size. The kernel buffer's cache line is brought into L1/L2 cache dirty; it will be

evicted to DRAM eventually but for read-back soon after, it's likely still hot.

Common Pitfalls and Their Consequences

Direct Dereference of User Pointer

```
/* WRONG – crashes kernel on bad pointer */

memcpy(kernel_buffer, buf, count);

/* RIGHT – recovers gracefully */

not_copied = copy_from_user(kernel_buffer, buf, count);

if (not_copied) return -EFAULT;
```

The compiler won't catch this. The static checker Sparse will (make C=1 during build). At runtime, a malicious or buggy program passing a bad pointer will kernel-panic without `copy_from_user`.

Returning Bytes Written Without Checking `copy_from_user` Return

```
/* WRONG – ignores partial copy */

copy_from_user(kernel_buffer, buf, count);

return count; /* wrong: may have only copied count - not_copied bytes */

/* RIGHT */

not_copied = copy_from_user(kernel_buffer, buf, count);

if (not_copied) return -EFAULT;

return count;
```

`copy_from_user` returns the number of bytes *not* copied — zero means full success. Treating a non-zero return as success tells the caller "I wrote N bytes" when you actually wrote fewer.

Not Returning 0 at EOF — Infinite read Loop

```
/* WRONG – cat loops forever */

static ssize_t mydev_read(struct file *filp, char __user *buf,
                         size_t count, loff_t *f_pos)

{
    copy_to_user(buf, kernel_buffer, buffer_used);

    return buffer_used; /* never returns 0, cat never terminates */

}

/* RIGHT – returns 0 when f_pos reaches end */

if (*f_pos >= buffer_used)

    return 0;
```

This is the single most common beginner read handler bug. If you ever find `cat /dev/yourdevice` hanging, this is why.

Not Updating `*f_pos` After Read

```
/* WRONG – f_pos never advances, infinite loop of same data */

*f_pos = 0; /* or just not updating it */

/* RIGHT */

*f_pos += bytes_to_copy;
```

Without advancing `f_pos`, every read call returns the same data from the same offset. `cat` sees a non-zero return value forever.

Forgetting Cleanup in Reverse Order

```
/* WRONG – unregister before destroy leaves dangling references */

unregister_chrdev_region(dev_num, 1); /* removes device number */

device_destroy(my_class, dev_num); /* tries to use the now-gone number */

/* RIGHT – destroy in reverse order of creation */

device_destroy(my_class, dev_num);

class_destroy(my_class);

cdev_del(&my_cdev);

unregister_chrdev_region(dev_num, 1);

kfree(kernel_buffer);
```

Out-of-order cleanup causes use-after-free in kernel data structures. The kernel may not crash immediately — the corruption may surface later as a mysterious oops in unrelated code.

Missing `.owner = THIS_MODULE` in `file_operations`

```
/* WRONG – module can be unloaded while device is open */

static const struct file_operations mydev_fops = {

    .open     = mydev_open,
    .read     = mydev_read,
    /* .owner missing */
};

/* RIGHT */

static const struct file_operations mydev_fops = {

    .owner     = THIS_MODULE, /* prevents rmmod while files are open */
    .open     = mydev_open,
    .read     = mydev_read,
};
```

Without `.owner`, `rmmode` succeeds while a process has the device open. The next `read()` jumps to a now-freed code address — instant kernel panic.

Knowledge Cascade: What You've Just Unlocked

→ VFS Dispatch Pattern = Interfaces Everywhere (cross-domain): The `file_operations` vtable you just filled in is the kernel's implementation of dynamic dispatch — the same pattern as C++ virtual methods, Go interfaces, and Rust trait objects. The call `f->f_op->read(filp, buf, count, f_pos)` is structurally identical to a C++ `virtual void read()` call that dispatches through a vtable pointer. When you next see a Go interface dispatch, a Rust `dyn Trait`, or a Java abstract class, you're looking at the same solution to the same problem: *dispatch to the right implementation at runtime without knowing the concrete type at compile time*. The kernel's vtable approach is the primordial form.

→ `f_pos` and `pread/pwrite` Atomicity (same domain, forward): The fact that `f_pos` is per-`struct file` (not per-inode) is why `pread(fd, buf, n, offset)` and `pwrite` exist. `pread` passes a separate offset without modifying `filp->f_pos`, allowing concurrent reads from different offsets without races — databases use `pread` exclusively for this reason. SQLite, PostgreSQL, RocksDB all call `pread` into their data files. The character device mechanism you just built is the same mechanism enabling atomic positioned I/O at scale.

→ udev Hotplug = Universal Device Discovery (cross-domain): `device_create()` sending a uevent is the same mechanism that fires when you plug in a USB drive and `/dev/sdb` appears, or when the kernel probes a PCIe device at boot and `/dev/nvme0n1` materializes. Android's `vold` (volume daemon), `systemd-udevd`, and the kernel's own

`kobject_uevent` are all operating the same infrastructure. You've just written code that participates in Linux's universal device discovery protocol. → Exception Table = Speculative Execution Safety Net (same domain): The `.ex_table` mechanism behind `copy_from_user` is conceptually related to how CPUs handle speculative execution faults: execute optimistically, recover if something goes wrong. The x86 `#GP` and `#PF` exception handlers consult the exception table the same way a CPU's microcode consults its own internal fault tables. More practically: this same mechanism powers `get_user()`, `put_user()`, and the `__copy_*` family — understanding it means you understand the safety foundation of the entire kernel-userspace boundary. → kmalloc to Slab Allocator to Memory Pressure (forward): The `kzalloc(4096, GFP_KERNEL)` call you made is trivial here, but it connects to one of the most important kernel subsystems: the memory reclaim path. Under memory pressure, `GFP_KERNEL` allocations can block while `kswapd` runs, pages are writeback to swap, slab caches are shrunk, and the OOM killer considers victims. In Milestone 3, when you add an ioctl for runtime buffer resize, you'll call `kfree` and `kmalloc` while the device is live — a perfect opportunity to think about what happens if the resize allocation fails at a moment when memory is under pressure. → This Is `/dev/null` (structural equivalence): `/dev/null`'s write handler is `return count` — accept all bytes, store nothing. Its read handler is `return 0` — always EOF. Your device is more capable. `/dev/zero`'s read handler fills the user buffer with zeros via `clear_user()` (a `copy_to_user` variant for zero-fill). `/dev/random` and `/dev/urandom` read from an entropy pool. The architecture is identical to what you just built — the only difference is what the read and write handlers actually do with the data.

Acceptance Criteria Checkpoint

Before moving to Milestone 3, verify every item:

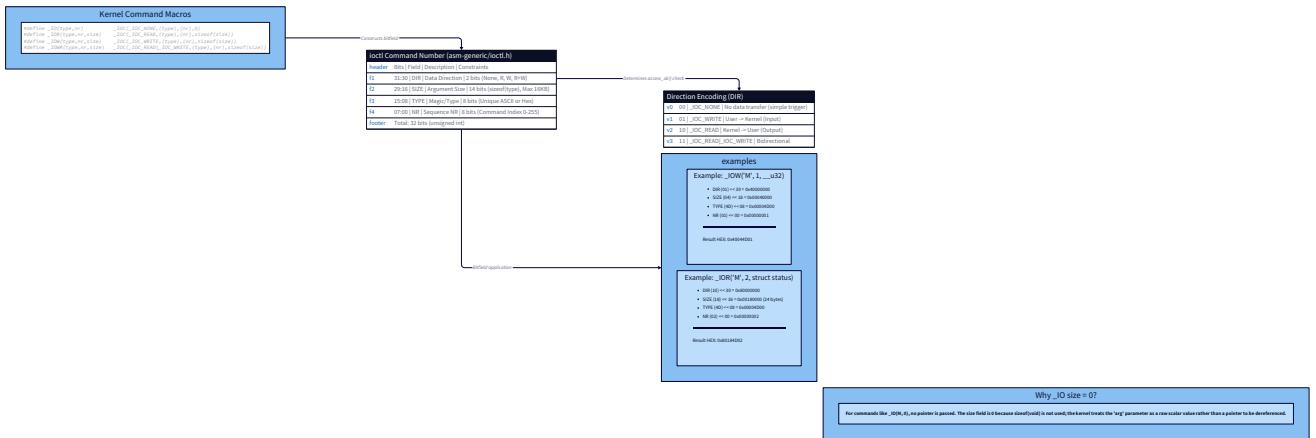
- `make` completes without warnings (`ccflags-y := -Werror active`)
- `sudo insmod mydevice.ko` succeeds; `dmesg | tail` shows `registered major=<N> minor=0`
- `grep mydevice /proc/devices` shows the major number
- `ls -la /dev/mydevice` shows the device node (created by udev, no manual `mknod`)
- `ls /sys/class/mydevice_class/` shows the device entry
- `echo "hello" | sudo tee /dev/mydevice` succeeds with output `hello`
- `sudo cat /dev/mydevice` outputs `hello` and terminates (no infinite loop)
- A second `sudo cat /dev/mydevice` also outputs `hello` (buffer persists between opens)
- `dmesg` shows open count increment and decrement for each `tee` and `cat` invocation
- Test with a malicious pointer: verify `-EFAULT` is returned and kernel does not panic
- `sudo rmmod mydevice` succeeds; `dmesg` shows clean exit message
- `/dev/mydevice` is removed after rmmod (udev responds to the remove uevent)
- `make C=1` (Sparse check) produces no warnings about `__user` annotation violations

Milestone 3: ioctl and /proc Interface

The Revelation: ioctl Is a Binary Protocol, Not a Switch Statement

You've built a character device that can read and write data. Now a process wants to *control* that device — resize its buffer, clear its contents, query its statistics. Your first instinct: add another write command with a special prefix. Write `"RESIZE:8192\n"` to the device and parse it in the write handler. This approach is used in production. Redis parses text commands over a socket. HTTP parses method strings. But it has a cost: you're multiplexing control and data through the same channel, which means your write handler now needs to answer "is this a data write or a control command?" before it can do anything. You've contaminated the data plane with control logic. Unix solved this problem in 1971 with `ioctl` — *input/output control*. The name is old; the insight is permanent. `ioctl` is a **separate system call** specifically for out-of-band device control, completely independent of the read/write data path. When a process calls `ioctl(fd, MYDEV_RESIZE, &new_size)`, it bypasses the VFS data path entirely and routes to your `unlocked_ioctl` handler with a structured command number and an argument. Here's where developers go wrong: they treat the command number as an arbitrary integer. They write `#define CMD_RESIZE 1`, `#define CMD_CLEAR 2`, and implement a `switch(cmd)` statement. This works — until it doesn't. **The problem with raw integers:** `ioctl(fd, 1, &new_size)` — is `1` a resize command for your driver? Or is it a terminal control command meant for the TTY layer? The kernel's ioctl dispatch has no way to know. If your device wraps another device

(a common pattern), the wrong driver might handle the wrong commands. Worse, `strace` will display your ioctl calls as `ioctl(3, 0x1, 0x7fff...)` — completely opaque. The solution is that **ioctl command numbers are encoded 32-bit values** — a binary protocol, not an enum. Every bit field has meaning:



31	16 15	8 7	2 1 0
	NR TYPE SIZE DIR		
	(cmd seq) (magic) (arg sz)		

bits 31:16 15:8 13:2 1:0

DIR (bits 1:0):

- 00 = _IO (no argument)
- 01 = _IOW (userspace → kernel, write to kernel)
- 10 = _IOR (kernel → userspace, read from kernel)
- 11 = _IOWR (bidirectional)

TYPE (bits 15:8): your unique magic number (e.g., 'M' = 0x4D)

NR (bits 31:16): sequential command number (0, 1, 2, ...)

SIZE (bits 13:2): sizeof(argument type), max 256 bytes

[[EXPLAIN:ioctl-command-number-encoding-(_iow/_ior-direction,-type,-nr,-size)]]

ioctl command number encoding (`_IOW/_IOR` direction, type, nr, size)]

The four macros that construct these values:

```

_I0(type, nr)          /* no argument */

_IOW(type, nr, argtype) /* userspace writes TO kernel: ioctl(fd, cmd, &data) passes data
in */

_IOR(type, nr, argtype) /* kernel writes TO userspace: ioctl(fd, cmd, &data) fills data
out */

_IOWR(type, nr, argtype) /* bidirectional */

```

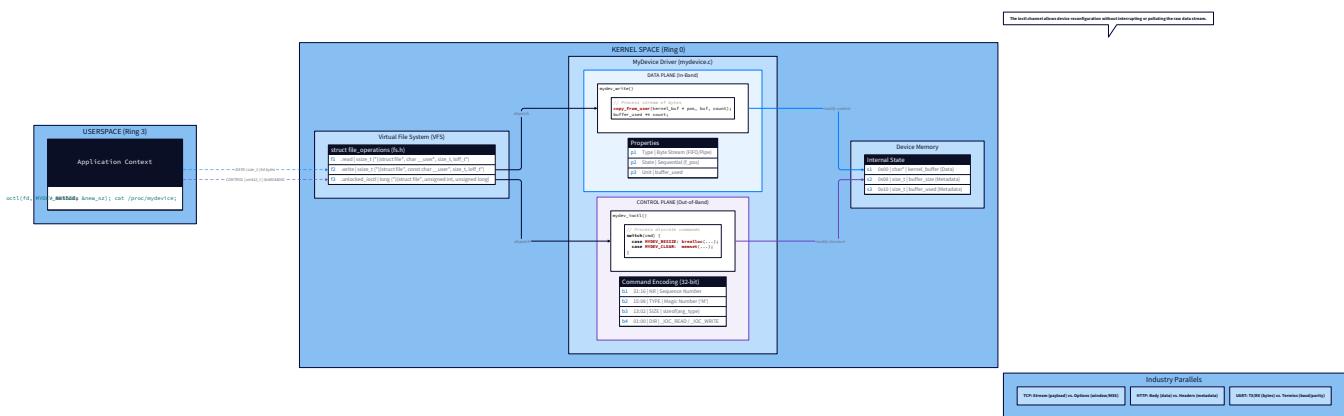
Each macro encodes the direction bits, your magic number, the command sequence number, and `sizeof(argtype)` into a single 32-bit constant. `_IOW('M', 0, int)` and `_IOW('N', 0, int)` produce

different 32-bit values — a command from a different driver with a different magic number will never accidentally match yours. This encoding exists because:

1. **strace decodes it symbolically**: with the right header, `strace` translates `ioctl(3, MYDEV_RESIZE, 0x7fff12340)` instead of `ioctl(3, 0xc0044d00, 0x7fff12340)`.
2. **The kernel can validate direction and size**: before your handler runs, the ioctl infrastructure knows whether the argument is incoming data, outgoing data, or both — and how large it is.
3. **Collision prevention**: your magic number namespaces your commands. `Documentation/userspace-api/ioctl/ioctl-number.rst` in the kernel tree is a registry of allocated magic numbers. Pick one that's not taken. Now let's build it.

The Architecture: Control Plane Meets Data Plane

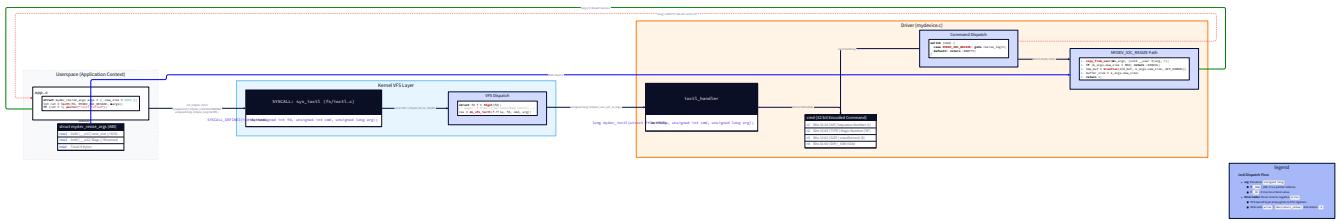
Before writing code, look at the full picture of what you're adding to your existing driver.



Your device now has two distinct interfaces:

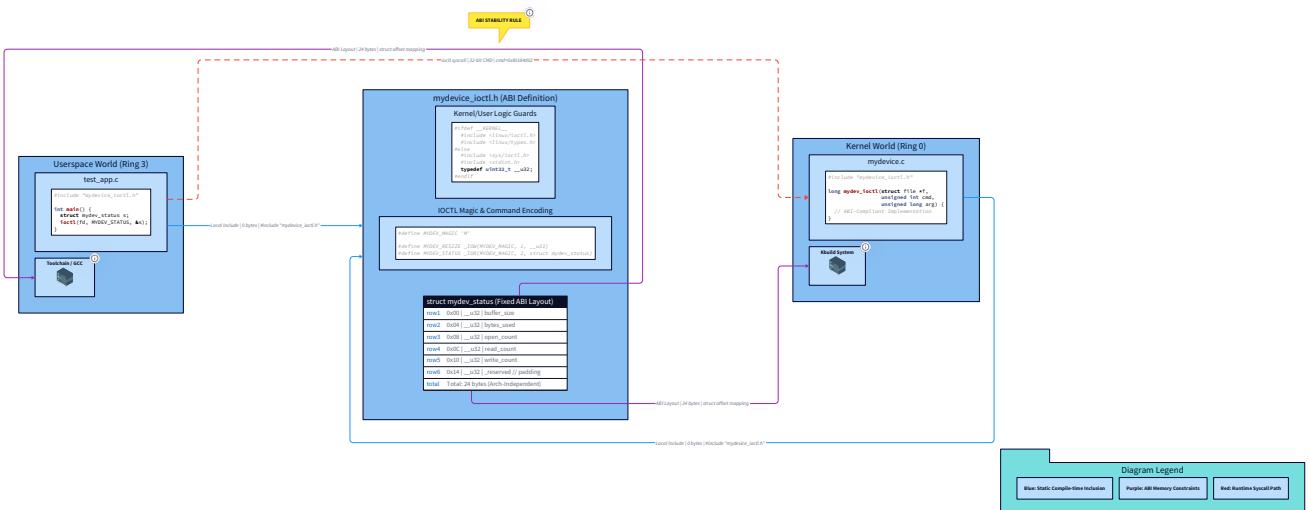
- **Data plane** (`read` / `write`): streams of bytes through the kernel buffer — the pipe
 - **Control plane** (`ioctl`): structured commands that change device behavior — the remote control
- This separation isn't arbitrary. It mirrors the same pattern across the entire systems stack:

TCP has a data stream (`send` / `recv`) and out-of-band control (socket options via `setsockopt`). HTTP has a body and headers. gRPC has streaming RPCs and unary control calls. Your character device has `read` / `write` and `ioctl`. The architectural pattern is identical: keep structured control commands out of the data stream so both can be reasoned about independently. The `/proc` entry you're adding is a third interface — **introspection**. It's read-only, always human-readable, and doesn't affect device state. It's the observability window into your driver's internals.



The Shared Header: One Definition, Two Worlds

The fundamental challenge of ioctl design is that the command definitions must be **identical** in kernel code and userspace code. If the kernel defines `MYDEV_RESIZE` as `0xc0044d00` and userspace computes `0xc0084d00` (because someone changed the argument size), the ioctl call silently fails — or worse, passes but transfers the wrong number of bytes. The solution: a single header file included by both.



```
/* mydevice_ioctl.h – shared between kernel module and userspace programs C

*
* This file defines the ioctl interface for mydevice. It is the ABI contract:
* once userspace programs link against these definitions, changing them
* breaks binary compatibility.
*
* Include guards prevent double-inclusion in both environments.
*/
#ifndef MYDEVICE_IOCTL_H

#define MYDEVICE_IOCTL_H

/* Guard against kernel-only headers when this file is included by userspace */

#ifndef __KERNEL__
# include <linux/ioctl.h>      /* __IOW, __IOR, __IOWR macros */
# include <linux/types.h>       /* __u32, __u64, etc. */
#else
# include <sys/ioctl.h>        /* userspace ioctl definitions */
# include <stdint.h>

    typedef uint32_t __u32;
    typedef uint64_t __u64;
#endif

/* Magic number: uniquely identifies this driver's ioctl namespace.
* 'M' (0x4D) is used here for illustration.
* In production: check Documentation/userspace-api/ioctl/ioctl-number.rst
* and pick an unregistered value.
*/
#define MYDEV_MAGIC 'M'
```

```

/*
 * Argument struct for status query.
 *
 * This struct is part of the ABI: fields CANNOT be reordered or removed
 *
 * without breaking userspace programs that compiled against this header.
 */

struct mydev_status {

    __u32 buffer_size;      /* total allocated buffer size in bytes */

    __u32 bytes_used;      /* bytes currently stored in buffer */

    __u32 open_count;       /* number of processes with device open */

    __u32 read_count;       /* total number of successful reads */

    __u32 write_count;      /* total number of successful writes */

    __u32 _reserved;        /* explicit padding for future use */

};

/*
 * ioctl command definitions.
 *
 * MYDEV_CLEAR: _IO - no argument; clears buffer contents
 *
 * MYDEV_RESIZE: _IOW - userspace sends new size (uint32); kernel resizes buffer
 *
 * MYDEV_STATUS: _IOR - kernel fills struct mydev_status; userspace reads it
 */

#define MYDEV_CLEAR    _IO(MYDEV_MAGIC,  0)

#define MYDEV_RESIZE   _IOW(MYDEV_MAGIC, 1, __u32)

#define MYDEV_STATUS   _IOR(MYDEV_MAGIC, 2, struct mydev_status)

#endif /* MYDEVICE_IOCTL_H */

```

Study this header carefully. Every design choice is deliberate: `__u32` instead of `uint32_t` or `unsigned int`: On different architectures and

compilation environments, `unsigned int` might be 16 or 32 bits. `_u32` is guaranteed 32-bit everywhere in Linux kernel headers. In userspace, we typedef it to `uint32_t` from `<stdint.h>` which has the same guarantee. The header conditionally provides whichever definition is appropriate.

`__KERNEL__` guard: The Kbuild system defines `__KERNEL__` when compiling kernel code. Userspace compilers don't define it. This lets you use one file in both contexts: the kernel side gets `<linux/ioctl.h>` (kernel macros), the userspace side gets `<sys/ioctl.h>` (POSIX definitions). The underlying math is the same — the macros produce identical bit patterns. `_reserved` in `mydev_status`: ABI stability practice. If you later add a field, you can replace `_reserved` instead of growing the struct (which would break the size encoding in `_IOR`). Explicit padding documents that you thought about future extensibility. The ABI contract: Once `mydev_status` is released and userspace programs link against it, the layout of this struct is frozen. The kernel's guarantee of stable userspace ABI means that a userspace binary compiled against this header in 2026 must still work with a kernel driver in 2036 — even if the driver's internal representation changes. This is why the kernel's internal structs change freely (they're not ABI) but the ioctl structs in public headers do not. This is the same concern that makes Protocol Buffer field numbers immutable and that governs JSON API versioning.

Implementing `unlocked_ioctl`

Add the ioctl handler to your driver. First, add the new fields to your device state:

```
/* Additional state in the kernel module – add to module-level globals */

static atomic_t read_count = ATOMIC_INIT(0);

static atomic_t write_count = ATOMIC_INIT(0);
```

Now the ioctl handler:

```
#include "mydevice_ioctl.h" /* shared definitions */

#include <linux/uaccess.h> /* copy_from_user, copy_to_user */

#include <linux/slab.h> /* kmalloc, kfree */

static long mydev_ioctl(struct file *filp, unsigned int cmd, unsigned long arg)

{

    int ret = 0;

    /*

     * Validate the magic number before processing any command.

     *

     * _IOC_TYPE(cmd) extracts bits [15:8] – the magic number.

     * If it doesn't match MYDEV_MAGIC, this command wasn't meant for us.

     * Return -ENOTTY: "inappropriate ioctl for device" – the POSIX standard

     * error for an unsupported ioctl command. NOT -EINVAL, NOT -ENOENT.

     * Tools like strace rely on -ENOTTY to identify unsupported commands.

     */

    if (_IOC_TYPE(cmd) != MYDEV_MAGIC)

        return -ENOTTY;

    /*

     * Validate the command number is within our range.

     * _IOC_NR(cmd) extracts bits [31:16] – the sequence number.

     * We've defined commands 0, 1, 2; reject anything above 2.

     */

    if (_IOC_NR(cmd) > 2)

        return -ENOTTY;

    switch (cmd) {

        /* — MYDEV_CLEAR —————— */

        case MYDEV_CLEAR:
```

```

/*
 * _IO: no argument. The 'arg' parameter is unused.
 *
 * Zero out the buffer contents and reset the usage counter.
 *
 * memset is safe here: kernel_buffer is a kernel pointer,
 * memset operates only in kernel address space.
 */

memset(kernel_buffer, 0, buffer_size_bytes);

buffer_used = 0;

printk(KERN_INFO "mydev: ioctl CLEAR – buffer cleared\n");

break;

/* — MYDEV_RESIZE —————— */
case MYDEV_RESIZE: {

/*
 * _IOW: userspace writes a __u32 to the kernel.
 * 'arg' is the userspace address of the __u32 new size value.
 * We must use copy_from_user to read it safely.
 */

__u32 new_size;

char *new_buffer;

/*
 * copy_from_user(kernel_dst, user_src, n)
 * Copies n bytes from userspace address 'arg' into kernel variable 'new_size'.
 * Returns bytes NOT copied; non-zero means fault → -EFAULT.
 */

if (copy_from_user(&new_size, (__u32 __user *)arg, sizeof(__u32)))

    return -EFAULT;
}

```

```
/* Validate the requested size */

if (new_size == 0 || new_size > (1U << 20)) { /* 0 to 1MB limit */

    printk(KERN_WARNING "mydev: ioctl RESIZE: invalid size %u\n", new_size);

    return -EINVAL;

}

/*
 * Allocate the new buffer before freeing the old one.

 * This is the "allocate-then-swap" pattern: if the new allocation
 * fails, we still have the old buffer intact – the device keeps working.

 */

new_buffer = kzalloc(new_size, GFP_KERNEL);

if (!new_buffer)

    return -ENOMEM;

/*
 * If new_size is smaller than current content, we must truncate.

 * Copy only min(buffer_used, new_size) bytes to preserve existing data.

*/

if (buffer_used > new_size) {

    printk(KERN_WARNING "mydev: RESIZE truncating %zu bytes to %u\n",

           buffer_used, new_size);

    buffer_used = new_size;

}

memcpy(new_buffer, kernel_buffer, buffer_used);

/* Swap: free old, install new */

kfree(kernel_buffer);

kernel_buffer = new_buffer;

buffer_size_bytes = new_size;
```

```

        printk(KERN_INFO "mydev: ioctl RESIZE - buffer resized to %u bytes\n", new_size);

        break;

    }

/* — MYDEV_STATUS ----- */

case MYDEV_STATUS: {

    /*

     * _IOR: kernel writes a struct mydev_status TO userspace.

     * 'arg' is the userspace address where we should write the struct.

    */

    struct mydev_status status;

    /* Fill the status struct from current device state */

    status.buffer_size = (u32)buffer_size_bytes;

    status.bytes_used = (u32)buffer_used;

    status.open_count = (u32)atomic_read(&open_count);

    status.read_count = (u32)atomic_read(&read_count);

    status.write_count = (u32)atomic_read(&write_count);

    status._reserved = 0;

    /*

     * copy_to_user(user_dst, kernel_src, n)

     * Copies n bytes from kernel struct 'status' to userspace address 'arg'.

     * Returns bytes NOT copied; non-zero → -EFAULT.

    */

    if (copy_to_user((struct mydev_status __user *)arg,
                    &status, sizeof(struct mydev_status)))

        return -EFAULT;

    break;
}

```

```

default:

/*
 * This should not be reached if the magic number and NR checks
 * above are correct. But defensive programming demands a default.
 */

return -ENOTTY;

}

return ret;
}

```

Register this handler in your `file_operations` struct:

```

static const struct file_operations mydev_fops = {

    .owner          = THIS_MODULE,
    .open           = mydev_open,
    .release        = mydev_release,
    .read           = mydev_read,
    .write          = mydev_write,
    .unlocked_ioctl = mydev_ioctl, /* ← add this */
};

C

```

Why `unlocked_ioctl`, not `ioctl`?

If you look at older kernel code (pre-2.6.36), you'll see a `.ioctl` field in `file_operations`. This has been removed. The historical `.ioctl` handler was called with the **Big Kernel Lock** (BKL) held — a global mutex that serialized all ioctl calls across the entire kernel. The BKL was a performance disaster removed in Linux 2.6.39. `unlocked_ioctl` means "your ioctl handler, without any BKL." You are responsible for your own locking (which we'll add in Milestone 4). For a single-buffer device, the operations in our three commands are simple enough that races are benign for now — but we'll tighten this in the concurrency milestone. There's also `.compat_ioctl` — called when a 32-bit process runs on a 64-bit kernel and makes an ioctl call. Pointer sizes differ between 32-bit and 64-bit, so struct layouts may differ. For our `mydev_status` (all `__u32` fields),

layout is identical across architectures. If you had 64-bit fields or pointers in your struct, you'd need `.compat_ioctl` to handle the translation.

The `-ENOTTY` Contract

The choice of `-ENOTTY` for unknown commands is not arbitrary — it's the POSIX-mandated error for "this ioctl is not applicable to this file descriptor." The name comes from "not a TTY" — the first `ioctl` commands were TTY-specific, so the error for "wrong device type" was named for the original use case. Today it means "this device doesn't support that command." Using `-EINVAL` instead (as many drivers do incorrectly) is technically wrong per POSIX and confuses tools that inspect ioctl return codes.

Validating Direction and Size Before Your Handler Runs

The `_IOC_DIR` and `_IOC_SIZE` fields in the encoded command number enable a useful pre-validation pattern. Before dispatching to your `switch`, you can verify that the userspace pointer is valid for the expected access:

```

/*
 * Optional but recommended: verify userspace pointer accessibility
 * based on the encoded direction and size.
 *
 * access_ok(addr, size) returns true if the address range [addr, addr+size)
 * is a valid userspace address. It does NOT guarantee the pages are present
 * (copy_from_user handles faults), but it catches obvious errors like
 * kernel addresses passed as userspace pointers.
 */

if (_IOC_DIR(cmd) & _IOC_READ) {

    /* Kernel will READ from user – user buffer must be writable by us */

    if (!access_ok((void __user *)arg, _IOC_SIZE(cmd)))

        return -EFAULT;

}

if (_IOC_DIR(cmd) & _IOC_WRITE) {

    /* Kernel will WRITE from user – user buffer must be readable by us */

    if (!access_ok((void __user *)arg, _IOC_SIZE(cmd)))

        return -EFAULT;

}

```

Note on direction naming confusion: The `_IOC_READ` and `_IOC_WRITE` bits in the command number describe the **userspace perspective**. `_IOC_READ` means "userspace is reading from the kernel" — which means the kernel writes to the user buffer. This is the `_IOR` direction. `_IOC_WRITE` means "userspace is writing to the kernel" — the `_IOW` direction. This reversal from the kernel's perspective trips up many developers. When in doubt, think from the userspace side. `access_ok()` is a fast, architecture-specific check that verifies the address doesn't fall in kernel space. On x86_64, it checks that `arg + size <= TASK_SIZE_MAX` (the upper bound of user address space). Note that modern kernels (5.x+) have simplified `access_ok` — it no longer takes a type parameter (it was `access_ok(VERIFY_READ, addr, size)` in older kernels). Update your code for the kernel version you're targeting.

The `/proc` Entry: Runtime Introspection

While ioctl gives processes *control* over the device, `/proc` gives *humans* (and monitoring tools) a window into the device's runtime state. The key insight:

`/proc` entries don't exist on disk. When you `cat /proc/mydevice`, the kernel calls your function to generate the content on-the-fly. The file system is a user interface metaphor, not actual storage. This is the virtual filesystem trick: the same VFS that dispatches `open` / `read` on real files dispatches them on `/proc` entries — but instead of reading from disk, it calls a kernel function. `cat /proc/cpuinfo` doesn't read a file; it calls a function that queries each CPU's registers and formats the result as text. You're now building the same mechanism.

The Problem with Raw proc Callbacks

Before `seq_file` existed, drivers implemented a `.read_proc` callback that was handed a `char *page` buffer and had to write into it. This was simple — until the data exceeded one page (4KB). A `/proc` entry with more than 4KB of content had to implement complex offset tracking manually. Most drivers got it wrong: they'd return the same content on every read (infinite loop) or miss content on reads that started mid-page.

🔑 Foundation: seq_file abstraction and iterator pattern

What it IS

The `seq_file` API is a standard Linux kernel abstraction designed to simplify the creation of virtual files (typically found in `/proc` or `/sys`) that output sequences of data. When a user-space program reads a file, the kernel must provide data in chunks. If the data is a long list—such as a list of all active network connections or loaded modules—managing the buffer offsets, multi-page reads, and partial transfers manually is complex and error-prone.

The `seq_file` interface uses the **Iterator Pattern**. It hides the "plumbing" of the filesystem (offsets, buffer sizes, and `read()` syscall mechanics) and asks the developer to provide only four functional callbacks:

1. `start` : Move to the beginning of the sequence or a specific position.
2. `next` : Move to the next item in the sequence.
3. `show` : Format the current item into the output buffer (usually using `seq_printf`).
4. `stop` : Clean up after the sequence is finished.

WHY the reader needs it right now

If you are writing a kernel module that needs to "leak" information to user space—like a list of custom device states or a debugging log—you shouldn't use the raw `file_operations.read` handler. Doing so requires you to track how many bytes the user has already read and handle cases where a single record is split across two separate `read()` calls.

By using `seq_file`, you focus entirely on your data structure (e.g., a linked list or an array). The kernel handles the "bookkeeping": it manages a temporary buffer, calls your `show` function until the buffer is full, stops, and resumes exactly where it left off when the user asks for more data. It ensures that even if a user-space tool like `cat` reads your 50KB list in 4KB chunks, the data remains consistent and the code remains simple.

ONE key insight or mental model

The "Stateful Cursor" Model. Think of `seq_file` as a cursor moving through a database. You don't tell the kernel "give me bytes 4096 to 8192 of this file." Instead, you tell the kernel, "I am currently at item #50 in my list; show it, then increment the cursor to #51."

The `seq_file` layer translates the user's **byte-based requests** (e.g., "I want the next 1KB") into **object-based iterations** (e.g., "Give me the next three items in your list"). This decoupling ensures that you never have to worry about a "partial print" cutting an integer or a string in half between two system calls.

`seq_file` solves this properly with an **iterator pattern**: your driver implements four functions (`start`, `next`, `stop`, `show`) that the `seq_file` infrastructure calls in a controlled loop, handling buffering, offset tracking, and partial reads correctly.



For simple single-record `/proc` entries (like device statistics), you use an even simpler helper that wraps `seq_file` into a single-function interface.

Implementation

Add the proc implementation to your module:

```
#include <linux/proc_fs.h>      /* proc_create, proc_remove */

#include <linux/seq_file.h>      /* seq_file, seq_printf, single_open */

#define PROC_NAME "mydevice"

/*
 * seq_show: called by the seq_file infrastructure to generate content.
 *
 * For a single-record /proc entry (all data fits in one logical "show"),
 * single_open() wraps this into a minimal seq_file iterator:
 *   start() returns a non-NULL pointer (signaling "data exists")
 *   show() calls this function to write content
 *   next() returns NULL (only one record)
 *   stop() does nothing
 *
 * seq_printf(m, fmt, ...) writes formatted text into the seq_file buffer.
 * It handles buffer management, overflow detection, and partial reads.
 * Never use sprintf or printk inside a seq_show function.
 */

static int mydev_proc_show(struct seq_file *m, void *v)
{
    seq_printf(m, "==== mydevice statistics ====\n");

    seq_printf(m, "buffer_size: %zu bytes\n", buffer_size_bytes);

    seq_printf(m, "bytes_used: %zu bytes\n", buffer_used);

    seq_printf(m, "open_count: %d\n", atomic_read(&open_count));

    seq_printf(m, "read_count: %d\n", atomic_read(&read_count));

    seq_printf(m, "write_count: %d\n", atomic_read(&write_count));

    return 0;
}
```

C

```
/*
 * mydev_proc_open: called when a process opens /proc/mydevice.
 *
 * single_open(file, show_fn, data) sets up the seq_file machinery
 * with a one-shot iterator that calls mydev_proc_show exactly once.
 * The 'data' pointer is available as m->private inside show_fn.
 */

static int mydev_proc_open(struct inode *inode, struct file *file)
{
    return single_open(file, mydev_proc_show, NULL);
}

/*
 * proc_fops: file operations for the /proc entry.
 *
 * We implement open with our wrapper, and delegate read, llseek,
 * and release to the seq_file infrastructure. These seq_* functions
 * know how to handle partial reads, seeks, and cleanup for seq_file.
 */

static const struct proc_ops mydev_proc_ops = {
    .proc_open    = mydev_proc_open,
    .proc_read    = seq_read,        /* seq_file's read handler */
    .proc_llseek   = seq_llseek,     /* seq_file's llseek handler */
    .proc_release  = single_release,/* cleans up single_open state */
};

static struct proc_dir_entry *proc_entry;
```

proc_ops vs file_operations : In kernels 5.6+, `/proc` entries use `struct proc_ops` instead of `struct file_operations`. The fields are similar (`proc_open`, `proc_read`, etc.) but `proc_ops` omits some fields that don't apply to proc entries and uses different defaults. If your target kernel is older than 5.6, use `file_operations` with a `.owner = THIS_MODULE` field. Register and remove the proc entry in your init/exit:

```
/* In mydev_init, after device_create succeeds: */

proc_entry = proc_create(PROC_NAME, 0444, NULL, &mydev_proc_ops);

/*
 *          ^
 *          ^      ^
 *          name    mode  parent  ops
 *
 * mode 0444: world-readable, nobody can write
 * parent NULL: create in /proc/ root
 * Returns NULL on failure (not IS_ERR – different error convention!) */

if (!proc_entry) {

    printk(KERN_ERR "mydev: proc_create failed\n");

    ret = -ENOMEM;

    goto err_proc;
}

printk(KERN_INFO "mydev: /proc/%s created\n", PROC_NAME);

/* Add to error unwind: */

err_proc:

    device_destroy(my_class, dev_num);

    /* ... rest of error chain ... */

/* In mydev_exit, before device_destroy: */

proc_remove(proc_entry);
```

Note the different error-checking convention for `proc_create` : it returns `NULL` on failure, not an `ERR_PTR()` value. You check `if (!proc_entry)`, not `if (IS_ERR(proc_entry))`. The kernel API is

not perfectly consistent — this is one of those inconsistencies you have to memorize. Now wire up the read/write counters in your existing handlers:

```
static ssize_t mydev_read(struct file *filp, char __user *buf,
                         size_t count, loff_t *f_pos)

{
    /* ... existing read logic ... */

    /* After a successful copy_to_user: */

    atomic_inc(&read_count);

    return (ssize_t)bytes_to_copy;
}

static ssize_t mydev_write(struct file *filp, const char __user *buf,
                         size_t count, loff_t *f_pos)

{
    /* ... existing write logic ... */

    /* After a successful copy_from_user: */

    atomic_inc(&write_count);

    return (ssize_t)bytes_to_copy;
}
```

How seq_file Actually Works Under the Hood

Understanding `seq_file` internally is worth a few minutes — it explains why it exists and why rolling your own offset tracking is a bad idea. When `cat /proc/mydevice` calls `read(fd, buf, 65536)`, the kernel calls your `proc_read` handler — which is `seq_read` from the `seq_file` library. Here's what `seq_read` does internally:

```

seq_read(filp, user_buf, count, f_pos):
1. If seq_file buffer is empty (first call or after exhaustion):
   call start(m, &m->index) → get first record (or NULL if done)
   if not NULL:
      call show(m, record) → write content into m->buf
      call next(m, record, &m->index) → advance to next record
      loop until done or m->buf is full
   call stop(m, record)
2. Copy min(m->buf used, count) bytes → copy_to_user(user_buf, m->buf, ...)
3. If m->buf has remaining content: update f_pos, return bytes copied
4. If m->buf exhausted: return 0 (EOF)

```

For `single_open`, the iterator is trivially simple: `start` returns a non-NULL dummy pointer once, `next` returns NULL immediately, `stop` does nothing. The `show` function (yours: `mydev_proc_show`) is called exactly once and writes all content. `seq_read` handles paging the output into userspace even if your output exceeds the kernel's internal buffer — it calls `start / show / next / stop` again from where it left off. This is why `seq_printf` is safe to call with arbitrarily large output: the `seq_file` infrastructure handles the case where your content doesn't fit in one `read()` call by maintaining position state internally. A hand-rolled proc callback that writes to a fixed `page` buffer and returns would truncate silently.

The Complete Updated Module

Here is the full module combining Milestone 2's character device with Milestone 3's ioctl and proc additions. This is the reference implementation:

```
/* mydevice.c - Character device with ioctl and /proc interface */

#include <linux/module.h>

#include <linux/kernel.h>

#include <linux/init.h>

#include <linux/fs.h>

#include <linux/cdev.h>

#include <linux/device.h>

#include <linux/uaccess.h>

#include <linux/slab.h>

#include <linux/atomic.h>

#include <linux/proc_fs.h>

#include <linux/seq_file.h>

#include <linux/string.h>

#include "mydevice_ioctl.h"      /* shared ioctl definitions */

MODULE_LICENSE("GPL");

MODULE_AUTHOR("Your Name <you@example.com>");

MODULE_DESCRIPTION("Character device: read/write, ioctl control, /proc introspection");

MODULE_VERSION("2.0");

/* — Configuration — */

#define DEVICE_NAME    "mydevice"

#define CLASS_NAME     "mydevice_class"

#define PROC_NAME      "mydevice"

/* — Device State — */

static dev_t          dev_num;

static struct cdev    my_cdev;

static struct class   *my_class;

static struct device *my_device;
```

```
static struct proc_dir_entry *proc_entry;

static char    *kernel_buffer;

static size_t   buffer_size_bytes = 4096;

static size_t   buffer_used       = 0;

static atomic_t open_count     = ATOMIC_INIT(0);

static atomic_t read_count     = ATOMIC_INIT(0);

static atomic_t write_count    = ATOMIC_INIT(0);

/* —— File Operations ————— */

static int mydev_open(struct inode *inode, struct file *filp)

{

    atomic_inc(&open_count);

    printk(KERN_INFO "mydev: open (count=%d)\n", atomic_read(&open_count));

    return 0;

}

static int mydev_release(struct inode *inode, struct file *filp)

{

    atomic_dec(&open_count);

    printk(KERN_INFO "mydev: release (count=%d)\n", atomic_read(&open_count));

    return 0;

}

static ssize_t mydev_read(struct file *filp, char __user *buf,

                         size_t count, loff_t *f_pos)

{

    size_t bytes_available, bytes_to_copy;

    unsigned long not_copied;

    if (*f_pos >= buffer_used)

        return 0; /* EOF */
```

```
bytes_available = buffer_used - *f_pos;

bytes_to_copy = min(count, bytes_available);

not_copied = copy_to_user(buf, kernel_buffer + *f_pos, bytes_to_copy);

if (not_copied)

    return -EFAULT;

*f_pos += bytes_to_copy;

atomic_inc(&read_count);

return (ssize_t)bytes_to_copy;

}

static ssize_t mydev_write(struct file *filp, const char __user *buf,
                           size_t count, loff_t *f_pos)

{

size_t space_available, bytes_to_copy;

unsigned long not_copied;

space_available = buffer_size_bytes - buffer_used;

if (space_available == 0)

    return -ENOSPC;

bytes_to_copy = min(count, space_available);

not_copied = copy_from_user(kernel_buffer + buffer_used, buf, bytes_to_copy);

if (not_copied)

    return -EFAULT;

buffer_used += bytes_to_copy;

atomic_inc(&write_count);

return (ssize_t)bytes_to_copy;

}

static long mydev_ioctl(struct file *filp, unsigned int cmd, unsigned long arg)

{
```

```
if (_IOC_TYPE(cmd) != MYDEV_MAGIC)

    return -ENOTTY;

if (_IOC_NR(cmd) > 2)

    return -ENOTTY;

switch (cmd) {

case MYDEV_CLEAR:

    memset(kernel_buffer, 0, buffer_size_bytes);

    buffer_used = 0;

    printk(KERN_INFO "mydev: ioctl CLEAR\n");

    break;

case MYDEV_RESIZE: {

    __u32 new_size;

    char *new_buffer;

    if (copy_from_user(&new_size, (__u32 __user *)arg, sizeof(__u32)))

        return -EFAULT;

    if (new_size == 0 || new_size > (1U << 20))

        return -EINVAL;

    new_buffer = kzalloc(new_size, GFP_KERNEL);

    if (!new_buffer)

        return -ENOMEM;

    if (buffer_used > new_size)

        buffer_used = new_size;

    memcpy(new_buffer, kernel_buffer, buffer_used);

    kfree(kernel_buffer);

    kernel_buffer = new_buffer;

    buffer_size_bytes = new_size;

    printk(KERN_INFO "mydev: ioctl RESIZE to %u bytes\n", new_size);
}
```

```
        break;

    }

    case MYDEV_STATUS: {

        struct mydev_status status = {

            .buffer_size = ( __u32 ) buffer_size_bytes,

            .bytes_used = ( __u32 ) buffer_used,

            .open_count = ( __u32 ) atomic_read( &open_count ),

            .read_count = ( __u32 ) atomic_read( &read_count ),

            .write_count = ( __u32 ) atomic_read( &write_count ),

            ._reserved = 0,

        };

        if ( copy_to_user( ( struct mydev_status __user * ) arg,
                           &status, sizeof( status ) ) )

            return -EFAULT;

        break;

    }

    default:

        return -ENOTTY;

    }

    return 0;
}

/* --- /proc Interface --- */

static int mydev_proc_show( struct seq_file *m, void *v)

{
    seq_printf(m, "==== mydevice statistics ====\n");

    seq_printf(m, "buffer_size: %zu bytes\n", buffer_size_bytes);

    seq_printf(m, "bytes_used: %zu bytes\n", buffer_used);
```

```

    seq_printf(m, "open_count: %d\n", atomic_read(&open_count));
    seq_printf(m, "read_count: %d\n", atomic_read(&read_count));
    seq_printf(m, "write_count: %d\n", atomic_read(&write_count));

    return 0;
}

static int mydev_proc_open(struct inode *inode, struct file *file)
{
    return single_open(file, mydev_proc_show, NULL);
}

static const struct proc_ops mydev_proc_ops = {
    .proc_open    = mydev_proc_open,
    .proc_read    = seq_read,
    .proc_lseek   = seq_lseek,
    .proc_release = single_release,
};

/* —— file_operations table ————— */

static const struct file_operations mydev_fops = {
    .owner        = THIS_MODULE,
    .open         = mydev_open,
    .release     = mydev_release,
    .read         = mydev_read,
    .write        = mydev_write,
    .unlocked_ioctl = mydev_ioctl,
};

/* —— Init / Exit ————— */

static int __init mydev_init(void)
{

```

```
int ret;

kernel_buffer = kzalloc(buffer_size_bytes, GFP_KERNEL);

if (!kernel_buffer) {

    return -ENOMEM;

}

ret = alloc_chrdev_region(&dev_num, 0, 1, DEVICE_NAME);

if (ret < 0) {

    printk(KERN_ERR "mydev: alloc_chrdev_region failed: %d\n", ret);

    goto err_region;

}

cdev_init(&my_cdev, &mydev_fops);

my_cdev.owner = THIS_MODULE;

ret = cdev_add(&my_cdev, dev_num, 1);

if (ret < 0) {

    printk(KERN_ERR "mydev: cdev_add failed: %d\n", ret);

    goto err_cdev;

}

my_class = class_create(THIS_MODULE, CLASS_NAME);

if (IS_ERR(my_class)) {

    ret = PTR_ERR(my_class);

    goto err_class;

}

my_device = device_create(my_class, NULL, dev_num, NULL, DEVICE_NAME);

if (IS_ERR(my_device)) {

    ret = PTR_ERR(my_device);

    goto err_device;

}
```

```
proc_entry = proc_create(PROC_NAME, 0444, NULL, &mydev_proc_ops);

if (!proc_entry) {
    ret = -ENOMEM;
    printk(KERN_ERR "mydev: proc_create failed\n");
    goto err_proc;
}

printk(KERN_INFO "mydev: initialized - /dev/%s and /proc/%s created\n",
       DEVICE_NAME, PROC_NAME);

return 0;

err_proc:
device_destroy(my_class, dev_num);

err_device:
class_destroy(my_class);

err_class:
cdev_del(&my_cdev);

err_cdev:
unregister_chrdev_region(dev_num, 1);

err_region:
kfree(kernel_buffer);

return ret;

}

static void __exit mydev_exit(void)
{
    proc_remove(proc_entry);
    device_destroy(my_class, dev_num);
    class_destroy(my_class);
    cdev_del(&my_cdev);
```

```
unregister_chrdev_region(dev_num, 1);

kfree(kernel_buffer);

printk(KERN_INFO "mydev: exited cleanly\n");

}

module_init(mydev_init);

module_exit(mydev_exit);
```

The Userspace Test Program

The test program lives in userspace but includes the shared header. This is the duality that makes the shared header design essential:

```
/* test_mydevice.c - Userspace test for ioctl and /proc interfaces */

#include <stdio.h>

#include <stdlib.h>

#include <fcntl.h>

#include <unistd.h>

#include <sys/ioctl.h>

#include <errno.h>

#include <string.h>

/*  
 * Include the shared header. The __KERNEL__ guard is NOT defined by  
 * the userspace compiler (gcc/clang), so this pulls in sys/ioctl.h  
 * and stdint.h on the userspace side.  
 */  
  
#include "mydevice_ioctl.h"  
  
#define DEVICE_PATH "/dev/mydevice"  
  
#define PROC_PATH    "/proc/mydevice"  
  
static void print_separator(const char *label)  
  
{  
    printf("\n---- %s ----\n", label);  
}  
  
static void read_proc(void)  
  
{  
    FILE *f = fopen(PROC_PATH, "r");  
  
    char line[256];  
  
    if (!f) {  
        perror("fopen /proc/mydevice");  
  
        return;  
    }
```

C

```
}

printf("Content of %s:\n", PROC_PATH);

while (fgets(line, sizeof(line), f))

    printf(" %s", line);

fclose(f);

}

int main(void)

{

    int fd, ret;

    /* Open the device */

    fd = open(DEVICE_PATH, O_RDWR);

    if (fd < 0) {

        perror("open " DEVICE_PATH);

        return EXIT_FAILURE;

    }

    printf("Opened %s (fd=%d)\n", DEVICE_PATH, fd);

    /* — Test 1: Write some data, check status ————— */

    print_separator("Test 1: Write and STATUS ioctl");

    const char *msg = "Hello from userspace!";

    ret = write(fd, msg, strlen(msg));

    if (ret < 0) {

        perror("write");

        goto out;

    }

    printf("Wrote %d bytes: \"%s\"\n", ret, msg);

    struct mydev_status status;

    /*
```

```
* MYDEV_STATUS is _IOR: kernel writes status TO us.

* We pass a pointer to our local struct; ioctl fills it.

*/
ret = ioctl(fd, MYDEV_STATUS, &status);

if (ret < 0) {

    perror("ioctl MYDEV_STATUS");

    goto out;

}

printf("Status: buffer_size=%u, bytes_used=%u, open_count=%u, "
       "reads=%u, writes=%u\n",
       status.buffer_size, status.bytes_used, status.open_count,
       status.read_count, status.write_count);

/* — Test 2: CLEAR ioctl */

print_separator("Test 2: CLEAR ioctl");

ret = ioctl(fd, MYDEV_CLEAR);

if (ret < 0) {

    perror("ioctl MYDEV_CLEAR");

    goto out;

}

printf("Buffer cleared.\n");

ret = ioctl(fd, MYDEV_STATUS, &status);

if (ret < 0) { perror("ioctl MYDEV_STATUS"); goto out; }

printf("After clear: bytes_used=%u (expected 0)\n", status.bytes_used);

/* — Test 3: RESIZE ioctl */

print_separator("Test 3: RESIZE ioctl");

__u32 new_size = 8192;

/*
```

```

    * MYDEV_RESIZE is _IOW: we write new_size TO the kernel.

    * Pass pointer to new_size; ioctl copies it into the kernel.

    */

ret = ioctl(fd, MYDEV_RESIZE, &new_size);

if (ret < 0) {

    perror("ioctl MYDEV_RESIZE");

    goto out;

}

printf("Resized buffer to %u bytes.\n", new_size);

ret = ioctl(fd, MYDEV_STATUS, &status);

if (ret < 0) { perror("ioctl MYDEV_STATUS"); goto out; }

printf("After resize: buffer_size=%u (expected 8192)\n", status.buffer_size);

/* — Test 4: Invalid ioctl returns -ENOTTY ————— */

print_separator("Test 4: Unknown ioctl → ENOTTY");

ret = ioctl(fd, _IO('Z', 99)); /* wrong magic number, unknown command */

if (ret < 0 && errno == ENOTTY)

    printf("PASS: unknown ioctl returned ENOTTY (errno=%d)\n", errno);

else

    printf("FAIL: expected ENOTTY, got ret=%d errno=%d\n", ret, errno);

/* — Test 5: Read /proc entry ————— */

print_separator("Test 5: /proc entry");

read_proc();

out:

close(fd);

return (ret < 0) ? EXIT_FAILURE : EXIT_SUCCESS;
}

```

Compile the test program with the shared header accessible:

```
# Compile: note we include the current directory for the shared header
#           BASH
gcc -Wall -Werror -I. -o test_mydevice test_mydevice.c

# Or with a Makefile target:

# test: test_mydevice.c mydevice_ioctl.h
#
#           BASH
gcc -Wall -Werror -I. -o test_mydevice test_mydevice.c
```

Run the full sequence:

```
# Build module
make

# Load
sudo insmod mydevice.ko

# Run tests (requires root for device access, or set device permissions)
sudo ./test_mydevice

# Check /proc directly
cat /proc/mydevice

# Verify strace decodes our ioctl (requires header in include path)
sudo strace ./test_mydevice 2>&1 | grep ioctl

# Unload
sudo rmmod mydevice
```

Expected `strace` output for the ioctl calls (with the header visible to strace — usually automatic if symbols match):

```
ioctl(3, MYDEV_STATUS, 0x7fff...) = 0
ioctl(3, MYDEV_CLEAR) = 0
ioctl(3, MYDEV_RESIZE, 0x7fff...) = 0
```

Three-Level View: What Happens During `ioctl(fd, MYDEV_STATUS, &status)`

Level 1 — Userspace: Your test program calls `ioctl(fd, MYDEV_STATUS, &status)`. The C library wraps this as the `ioctl(2)` syscall. The `cmd` argument is the 32-bit encoded value produced by `_IOR('M', 2, struct mydev_status)` — approximately `0x80184d02` (direction=read=2, size=24 bytes, type='M', nr=2). The `arg` is the stack address of `status`.

Level 2 — Kernel/VFS: `sys_ioctl` → `do_vfs_ioctl` → checks if `cmd` is a known VFS-level command (like `FIONREAD` or `FIOCLEX`) — it's not → calls `vfs_ioctl` → `filp->f_op->unlocked_ioctl(filp, cmd, arg)` → your `mydev_ioctl`. Your handler reads `open_count`, `read_count`, `write_count` from `atomic_t` variables, fills `struct mydev_status` on the kernel stack, and calls `copy_to_user`. The copy places the struct into the userspace stack frame.

Level 3 — Hardware: The `copy_to_user` for a 24-byte struct is a fast path. On `x86_64`, the struct fits in three cache lines (24 bytes < 64 bytes = one cache line, in fact). The CPU executes the copy with a few `mov` instructions and potential `rep movsb` for safety. The userspace page is likely hot in L1/L2 cache since the test program just passed its address. The `atomic_read` calls generate `mov` instructions with appropriate memory barriers — `atomic_read` on `x86` is just a regular load (`x86`'s memory model provides the necessary ordering guarantees without explicit barrier instructions in most cases). The total cost: maybe 200ns, dominated by the syscall overhead (~100ns) and cache accesses, not the copy itself.

Hardware Soul: The ioctl Path's Cache Footprint

When your ioctl handler runs, consider what memory it touches:

- `kernel_buffer` pointer: one word in the global data section — L1 hot after the first access

- `buffer_size_bytes`, `buffer_used`: adjacent `size_t` values in BSS — two cache lines apart at most, likely in the same 64-byte cache line if laid out consecutively by the compiler
- `open_count`, `read_count`, `write_count`: three `atomic_t` values (each a 4-byte `int`). If they're adjacent in memory, they share a cache line — 12 bytes total, fits in one 64-byte line. **Critical**: if multiple CPUs atomically increment these counters from different cores, each `atomic_inc` requires exclusive cache line ownership (the MESI cache coherency protocol's "M" (Modified) state). This is cache line bouncing — the line has to transfer between CPU caches on each write. For simple counters this is acceptable. For high-throughput code (thousands of operations per second), you'd use per-CPU counters and aggregate them on read. This is exactly how the kernel's own `percpu_counter` works. The `struct mydev_status` that you fill and copy:

```
struct mydev_status { /* total: 6 × 4 = 24 bytes */
    __u32 buffer_size;      /* offset 0 */
    __u32 bytes_used;      /* offset 4 */
    __u32 open_count;       /* offset 8 */
    __u32 read_count;       /* offset 12 */
    __u32 write_count;      /* offset 16 */
    __u32 _reserved;        /* offset 20 */
};

/* total: 24 bytes = 3/8 of one cache line */
```

24 bytes. One cache line (64 bytes) is more than enough. The `copy_to_user` for 24 bytes doesn't even need loop unrolling — the compiler generates a small sequence of register stores. The bottleneck in this ioctl is not the copy; it's the atomic reads and the syscall overhead.

Pitfalls: What Goes Wrong and Why

Using Wrong Magic Number — Silent Mismatch

```
/* WRONG - using a magic number someone else uses */

#define MYDEV_MAGIC 't' /* 't' is used by linux/tty.h TCGETS, TCSETS, etc. */

/* RIGHT - pick one from Documentation/userspace-api/ioctl/ioctl-number.rst */

#define MYDEV_MAGIC 'M' /* or whichever is unregistered */
```

If your magic number collides with a TTY magic number, a program calling your ioctl on a TTY might accidentally trigger your handler (or vice versa) when the numbers align. The magic number is the first line of

defense.

Returning Wrong Error for Unknown Command

```
/* WRONG - standard violation */

return -EINVAL;      /* "invalid argument" - ambiguous */

/* RIGHT - POSIX-mandated for unsupported ioctl */

return -ENOTTY;      /* "inappropriate ioctl for device" */
```

Programs testing for supported commands use the return code to distinguish "command not supported" (-ENOTTY) from "command supported but argument invalid" (-EINVAL). Returning -EINVAL for unknown commands breaks this distinction.

RESIZE Without Protecting Existing Data

```
/* WRONG - frees old buffer before copying to new */

kfree(kernel_buffer);

kernel_buffer = kmalloc(new_size, GFP_KERNEL);

if (!kernel_buffer) {

    /* original data is gone! kernel_buffer is now NULL! */

    return -ENOMEM;

}

/* RIGHT - allocate first, copy, then free */

new_buffer = kzalloc(new_size, GFP_KERNEL);

if (!new_buffer) return -ENOMEM; /* original buffer untouched */

memcpy(new_buffer, kernel_buffer, min(buffer_used, (size_t)new_size));

kfree(kernel_buffer);

kernel_buffer = new_buffer;
```

The "allocate-then-swap" pattern is a fundamental technique for making operations resilient to allocation failure. It applies equally to kernel code and userspace code (it's how `realloc` is supposed to be implemented for safety).

Not Handling Truncation on RESIZE-Down

```
/* WRONG – bytes_used may exceed new_size */

memcpy(new_buffer, kernel_buffer, buffer_used); /* writes past new_buffer end! */

/* RIGHT – truncate first */

if (buffer_used > new_size)

    buffer_used = new_size;

memcpy(new_buffer, kernel_buffer, buffer_used);
```

Writing `buffer_used` bytes into a `new_size`-byte buffer when `buffer_used > new_size` is a classic buffer overflow. In the kernel, this corrupts adjacent kernel memory — the consequences range from silent data corruption to immediate kernel panic, depending on what lives next to your buffer.

proc_create Error Check with IS_ERR

```
/* WRONG – proc_create returns NULL on error, not ERR_PTR */

if (IS_ERR(proc_entry)) {           /* always false! NULL != IS_ERR(NULL) */

    ret = PTR_ERR(proc_entry);

    goto err_proc;

}

/* RIGHT */

if (!proc_entry) {                  /* NULL check */

    ret = -ENOMEM;

    goto err_proc;

}
```

`IS_ERR(NULL)` returns false because `NULL` (0) is not in the error pointer range. The `proc_create` inconsistency (returning `NULL` vs `ERR_PTR`) is a historical accident in the kernel API. Always check the documentation for which convention a function uses.

Modifying buffer_size_bytes Without Updating /proc Show

If your proc show function reads `buffer_size_bytes` and your ioctl RESIZE modifies it, you need no extra work — both see the same global variable. But if you cache values (e.g., save `buffer_size_bytes` to a local at init time), your proc entry will show stale data after a resize. Always read live state in proc show handlers.

Knowledge Cascade: What You've Just Unlocked

→ ioctl = Control Plane / data-plane separation (cross-domain: network protocol design) The read/write vs. ioctl split you just implemented is the same architectural pattern as TCP's data stream vs. socket options (`setsockopt` / `getsockopt`). HTTP separates body (data plane) from headers (control plane). QUIC separates stream data from connection-level signals. gRPC offers streaming RPCs (high-throughput data) alongside unary calls (control). The insight — "don't mix control and data in the same channel" — appears everywhere because mixing them forces both sides to parse combined traffic, adding latency and complexity. When you next design a protocol or API and find yourself tempted to embed commands in the data stream, remember: TCP, HTTP, and your kernel driver all chose not to. → seq_file's iterator = Language-level iterators (cross-domain: language design) The `start` / `next` / `stop` / `show` interface `seq_file` gives you is structurally identical to Python's `__iter__` / `__next__` generator protocol, Rust's `Iterator` trait with its `next() -> Option<Item>` method, and database cursor iteration. The fundamental problem is the same: "produce a potentially large sequence lazily, without materializing it all at once." Python generators solve it with

`yield`. Rust iterators solve it with `next()`. `seq_file` solves it with `start / next / stop`. When you write a Python generator that yields database rows, you're implementing the same pattern as the `seq_file` infrastructure you just wired up. → `/proc` virtual filesystem = Everything-is-a-file philosophy (cross-domain: OS design) `/proc/mydevice` has no backing file on disk. Its "content" is generated by your `mydev_proc_show` function every time someone reads it. This is the same trick as `/proc/cpuinfo` (calls into the CPU topology code), `/proc/meminfo` (queries the page allocator), and `/proc/<pid>/status` (reads the task struct for that process). Plan 9 OS — Linux's philosophical ancestor in many ways — took this idea to its logical extreme: *everything* is a file, including network connections (open `/net/tcp/0/data` to make a TCP connection). FUSE (Filesystem in Userspace) implements the same idea in userspace: your code generates "file content" on-read. Understanding `/proc` means you understand FUSE, Plan 9, `/sys`, and the entire "virtual filesystem as control interface" philosophy. → Shared header ABI = Wire format stability (cross-domain: distributed systems) The `mydevice_ioctl.h` you just created is an ABI contract. `struct mydev_status`'s layout — field order, sizes, offsets — is now frozen. Any userspace binary compiled against this header must keep working against future kernel drivers. This is identical to Protocol Buffers: field numbers are immutable once assigned, because removing or reusing them breaks binary compatibility with old serialized data. It's the same as JSON API versioning (don't remove fields, only add). The Linux kernel's guarantee of stable userspace ABI is explicitly documented in `Documentation/process/stable-api-nonsense.rst` — the *internal* kernel API changes freely, but the userspace ABI (what you exposed in your header) is guaranteed stable. When you design any serialized format or

shared interface, you're facing the same problem your ioctl header just solved: separate the mutable internals from the immutable external contract. → Atomic counters → per-CPU counters → scalability (forward) Your `atomic_t read_count` increments on every read. With a single CPU, this is fine. With 16 CPUs all reading simultaneously, each `atomic_inc` causes the cache line holding `read_count` to bounce between CPU caches — a cache coherency round-trip that takes ~100 CPU cycles per operation. At 16 concurrent readers, your counter becomes a bottleneck. The kernel's solution: `percpu_counter` — each CPU maintains its own counter in CPU-local memory (no cache bouncing), and reads aggregate all CPU-local values. This is the same idea as Go's `sync.Pool` (per-goroutine allocation to avoid contention), or partitioned lock tables (one lock per shard, not one global lock). The pattern: "replace one shared resource with N per-CPU/per-thread resources, aggregate on read." You'll see it everywhere high-performance code needs shared counters.

Acceptance Criteria Checkpoint

Before moving to Milestone 4, verify every item:

- `mydevice_ioctl.h` compiles correctly when included in both kernel module (with `__KERNEL__` defined) and userspace C program (without `__KERNEL__`)
- `MYDEV_CLEAR`, `MYDEV_RESIZE`, `MYDEV_STATUS` produce distinct non-zero 32-bit values; verify by printing them in the test program
- `make` compiles the module with `ccflags-y := -Werror` and zero warnings
- `ioctl(fd, MYDEV_CLEAR)` zeroes the buffer; subsequent read returns 0 bytes (EOF immediately)
- `ioctl(fd, MYDEV_RESIZE, &new_size)` with `new_size=8192` changes buffer capacity; `MYDEV_STATUS` reports `buffer_size=8192`
- `ioctl(fd, MYDEV_RESIZE, &new_size)` with `new_size=0` returns `-EINVAL` (not a crash)
- `ioctl(fd, MYDEV_RESIZE, &new_size)` with content larger than `new_size` truncates safely (no overflow, no panic)

- `ioctl(fd, MYDEV_STATUS, &status)` fills all fields of `struct mydev_status` correctly; `read_count` and `write_count` increment after each read/write operation
- `ioctl(fd, _IO('Z', 99))` (wrong magic number) returns `-1` with `errno == ENOTTY`
- `cat /proc/mydevice` outputs all five statistics fields (`buffer_size`, `bytes_used`, `open_count`, `read_count`, `write_count`)
- `cat /proc/mydevice` called twice returns consistent data without hanging or double-printing
- `sudo rmmod mydevice` after `proc_create` removes `/proc/mydevice` cleanly (no stale proc entry)
- Test program `test_mydevice` compiles with `gcc -Wall -Werror -I. test_mydevice.c` and runs all four test cases without failure
- `strace ./test_mydevice 2>&1 | grep ioctl` shows ioctl calls with recognizable command names (not raw hex) when symbols are available

Milestone 4: Concurrent Access, Blocking I/O, and Poll Support

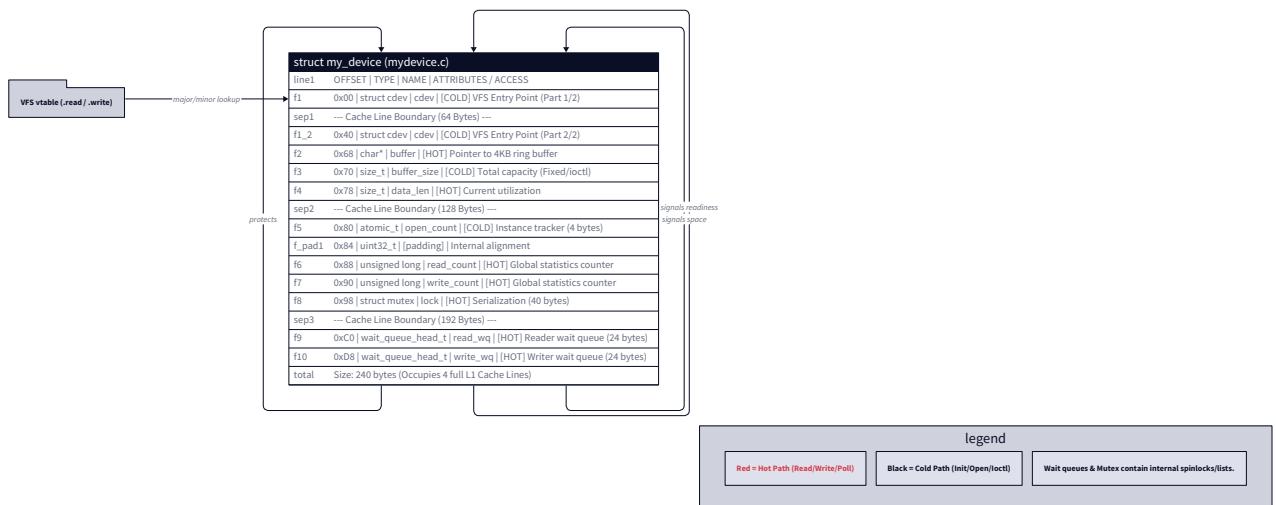
The Revelation: The Kernel Is Not a Thread — It's a State Machine

You've built a character device that works perfectly when one process uses it at a time. Now imagine four processes all calling `read()` simultaneously, three of them finding the buffer empty, and a fourth calling `write()` to fill it. In userspace, you'd reach for `pthread_mutex_lock()` and a condition variable. You'd think: "The reader locks the mutex, finds the buffer empty, calls `pthread_cond_wait()`, sleeps, wakes when the writer signals." Clean. Familiar. Here's the misconception that kernel newcomers carry from that model: they assume that "blocking" in the kernel works the same way, that `mutex_lock` in the kernel is just `pthread_mutex_lock` with a different name, and that `poll()` is just a kernel function that checks a boolean and returns "yes" or "no." Every part of that model needs surgery. **The first operation:** `mutex_lock()` in the kernel does sleep. But "sleeping" in the kernel is not a private affair between two threads. It is a state change on the **current process** — the actual Linux task (identified by its `struct task_struct`) that is executing your driver's code on behalf of the syscall. When your read handler calls `mutex_lock()` and the mutex is contended, the current process transitions from `TASK_RUNNING` (scheduled, consuming CPU cycles) to `TASK_INTERRUPTIBLE` (sleeping, removed from the run queue, not scheduled). The scheduler picks a different process to run. Your code is literally paused mid-function, in the kernel, and a completely different process starts running on that CPU. When the mutex becomes available, your process is placed back on the run queue and eventually resumes exactly where it paused — but potentially hundreds of milliseconds later, on a different CPU core. This is only safe in **process context** — code executing as part of a syscall from a user process. It is catastrophically illegal in **interrupt context** — code executing inside a hardware interrupt handler or a software interrupt (softirq, tasklet). In interrupt context, there is no "current process" to sleep. The CPU is responding to a hardware event, not running on behalf of

any process. If you call `mutex_lock()` and sleep in interrupt context, you corrupt the interrupt stack, the scheduler panics, and your machine halts. [[EXPLAIN:kernel-sleeping-primitives-vs-spinlocks-(process-context-vs-interrupt-context)|Kernel sleeping primitives vs spinlocks (process context vs interrupt context)]]

The second operation: `wait_event_interruptible()` is not a function that loops until something is true. It does something architecturally remarkable. It:

1. Adds the current process to a **wait queue** — a linked list of sleeping processes waiting for a condition
2. Sets the process state to `TASK_INTERRUPTIBLE`
3. Calls `schedule()` — voluntarily yields the CPU
4. Wakes up when another process calls `wake_up_interruptible()` on the wait queue
5. Re-checks the condition (spurious wakeup protection)
6. If condition is true: removes itself from the wait queue and returns 0
7. If a signal arrived: returns `-ERESTARTSYS` The writer calling `wake_up_interruptible()` doesn't resume the reader *immediately*. It transitions the reader from the wait queue back to the run queue and marks it `TASK_RUNNING`. The scheduler then decides when the reader actually gets CPU time — maybe microseconds later, maybe milliseconds. The writer continues running, finishes its write, releases the mutex, and returns. Only then does the scheduler get to pick the reader. **The third operation** — and this is the one that surprises everyone — `poll()` does not sleep inside your driver. Your `.poll` handler is called by the kernel's poll infrastructure, but it must return *immediately* with a readiness mask. It must never block. The sleeping happens *outside* your driver, in the VFS poll infrastructure. Your job is not to wait for readiness — your job is to **register your wait queue** with the poll table and **report current readiness**. The kernel's poll loop does the sleeping, wakes up when your wait queue fires, and calls your `.poll` handler again to re-check readiness. These three revelations reframe everything you're about to build. Let's build it.

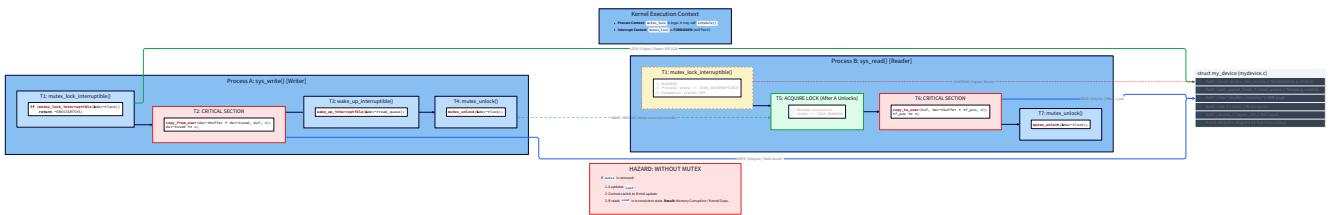


The Foundation: What You're Actually Building

Before any code, look at the structural transformation this milestone makes. In Milestone 3, your device had two unsynchronized global variables (`buffer_used`, `kernel_buffer`) accessed by handlers that could run simultaneously on different CPUs. Any concurrent access was undefined behavior — a data race that worked most of the time and produced subtle corruption the rest. This milestone adds:

- A `struct mutex dev_mutex` that serializes all access to `kernel_buffer` and `buffer_used`
- A `wait_queue_head_t read_queue` for blocking readers waiting for data
- A `wait_queue_head_t write_queue` for blocking writers waiting for space
- Logic to check `filp->f_flags & O_NONBLOCK` to implement non-blocking I/O
- A `.poll` file operation that uses `poll_wait()` to participate in `select / poll / epoll` The device's behavior changes fundamentally:

Scenario	Before (M3)	After (M4)
Two concurrent writes	Data race, corruption	Serialized by mutex
Read with empty buffer	Returns 0 (EOF) immediately	Blocks until data arrives
Read with O_NONBLOCK + empty	Returns 0 (EOF)	Returns -EAGAIN
Write with full buffer	Returns -ENOSPC	Blocks until space available
poll() on the fd	Not supported	Returns POLLIN when data exists
Ctrl+C during blocking read	Hangs forever (pre-M4)	Returns -EINTR/-ERESTARTSYS



Kernel Mutexes: The Right Lock for Process Context

What a Mutex Actually Is

The kernel has multiple locking primitives, each suited to a different context. Understanding *which* lock to use requires understanding what "context" means in the kernel. **Process context** is any code that executes as part of a system call on behalf of a user process. Your `mydev_read`, `mydev_write`, `mydev_ioctl`, and `mydev_open` handlers all run in process context — there is a `current` task, it has a process address

space, and it can sleep. **Interrupt context** is any code that executes in response to a hardware interrupt, softirq, or tasklet. There is no `current` task in the meaningful sense, no process address space, and **sleeping is forbidden** — the CPU must finish handling the interrupt and return to whatever it was doing before. The primitive you want for process context is `struct mutex`. It is a **sleeping lock**:

- If the mutex is available: `mutex_lock()` acquires it instantly (one atomic operation)
 - If the mutex is contended: `mutex_lock()` puts the current process to sleep until the mutex becomes available. The primitive for interrupt context (or any context where sleeping is forbidden) is `struct spinlock`. A spinlock **busy-waits** — it loops in a tight loop reading the lock variable until it becomes available. No sleeping, no `schedule()` call. But holding a spinlock while doing anything that sleeps (including `kmalloc(GFP_KERNEL)`, `copy_from_user`, `mutex_lock`) is catastrophic — it holds the CPU hostage. For your character device, **mutex is the correct choice**. Your file operation handlers run in process context. They may sleep. They copy data from/to userspace, which can fault and sleep.
- `GFP_KERNEL` allocations in ioctl handlers can sleep. Mutex handles all of this correctly.

```
#include <linux/mutex.h>

static DEFINE_MUTEX(dev_mutex);

/*
 * DEFINE_MUTEX(name) statically defines and initializes a mutex.
 *
 * Equivalent to:
 *
 *     struct mutex dev_mutex;
 *
 *     mutex_init(&dev_mutex);
 *
 *
 * Use for module-level mutexes. For dynamically allocated structs,
 *
 * use mutex_init(&obj->lock) in your allocation/init code.
 */

*/
```

The mutex API is small:

```
mutex_lock(&dev_mutex);          /* acquire; sleep if contended */

mutex_unlock(&dev_mutex);        /* release */

mutex_lock_interruptible(&dev_mutex); /* acquire; return -EINTR if signal arrives */

mutex_trylock(&dev_mutex);       /* acquire if available; return 0 if not */
```

For blocking I/O, use `mutex_lock_interruptible()` — it returns `-EINTR` if a signal arrives while the process is sleeping waiting for the mutex. This allows `Ctrl+C` to interrupt a write even if the mutex is contended.

Lock Ordering: if your driver ever needs two mutexes simultaneously (you won't here, but this principle is fundamental), you must acquire them in a consistent global order across all code paths. Acquiring A then B in one path and B then A in another path is a classic deadlock. Since this driver has one mutex, deadlock is only possible if your code calls `mutex_lock(&dev_mutex)` while already holding it — which deadlocks immediately because the kernel's mutex is not reentrant. You will see `WARNING: possible recursive locking detected` in dmesg if this happens.

Wait Queues: The Sleeping Heart of Blocking I/O

The Data Structure

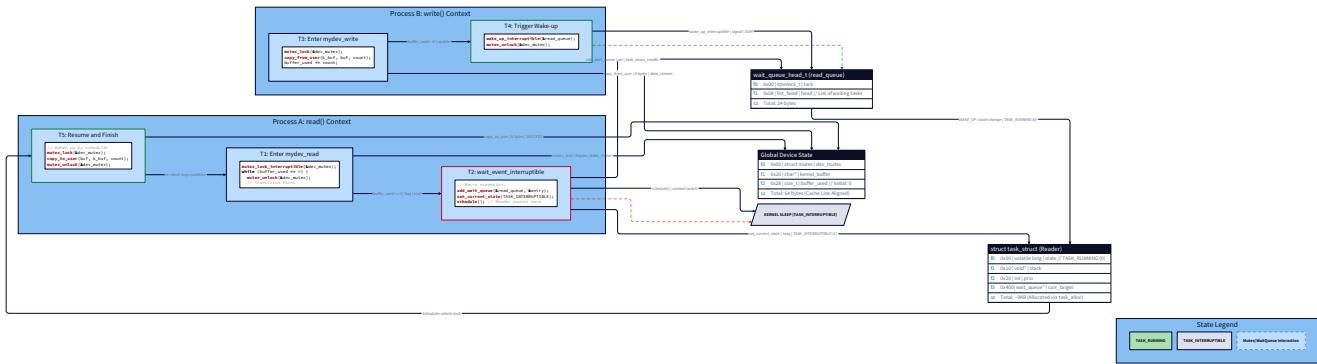
A **wait queue head** (`wait_queue_head_t`) is a linked list of sleeping processes, each waiting for a condition to become true. The list is protected by its own internal spinlock (which is why wait queue operations are safe to call from anywhere, including while holding your mutex — the internal spinlock protects the queue structure itself, and you release it before sleeping).

```
#include <linux/wait.h> C

static DECLARE_WAIT_QUEUE_HEAD(read_queue); /* for readers waiting on empty buffer */

static DECLARE_WAIT_QUEUE_HEAD(write_queue); /* for writers waiting on full buffer */

/*
 * DECLARE_WAIT_QUEUE_HEAD(name) statically defines and initializes.
 *
 * For dynamic allocation: init_waitqueue_head(&obj->wait_queue);
 */
```



`wait_event_interruptible` : The Full Mechanism

The macro `wait_event_interruptible(queue, condition)` is the standard pattern for blocking a process until a condition becomes true. Here is what it expands to, conceptually:

```
wait_event_interruptible(wq, condition):
    1. Check condition – if true, return 0 immediately (fast path, no sleeping)
    2. Prepare a wait_queue_entry: allocate a struct, set it to point to current task
    3. add_wait_queue(wq, &entry)           – insert into the wait queue
    4. set_current_state(TASK_INTERRUPTIBLE) – mark process as sleeping
    5. Check condition AGAIN (re-check after state change, before schedule)
        – this closes the TOCTOU window between adding to queue and sleeping
    6. If condition false AND no signal pending:
        schedule()                         – yield CPU; sleep here
        set_current_state(TASK_INTERRUPTIBLE) – re-arm for next iteration
        goto step 5
    7. set_current_state(TASK_RUNNING)   – re-arm as runnable
    8. remove_wait_queue(wq, &entry)     – remove from wait queue
    9. If signal_pending(current): return -ERESTARTSYS
    10. Return 0 (condition became true)
```

The double-check at step 5 is critical. Between step 3 (adding to wait queue) and step 6 (calling schedule), a writer might already have written data and called `wake_up_interruptible()`. Without the re-check, the reader would sleep on a now-non-empty buffer, potentially forever. This is the **lost wakeup problem**, and `wait_event_interruptible` handles it correctly. [[EXPLAIN:wait-queue-internals-and-the-thundering-herd-problem|Wait queue internals and the thundering herd problem]]

The Condition and the Mutex: A Critical Pattern

There's a subtle interaction between the mutex and the wait queue that confuses many developers. Your read handler:

1. **Acquires the mutex** (to safely check `buffer_used`)
 2. **Checks the condition** (`buffer_used == 0`)
 3. If empty: **must release the mutex before sleeping** You cannot hold the mutex while calling `wait_event_interruptible()` — that would block all writers from acquiring the mutex to add data,

causing a deadlock where the reader sleeps holding the lock that writers need. The pattern looks like this:

```
/* WRONG — deadlock: holds mutex while sleeping */

mutex_lock(&dev_mutex);

wait_event_interruptible(read_queue, buffer_used > 0); /* writers can't lock! */

/* ... read ... */

mutex_unlock(&dev_mutex);

/* RIGHT — release mutex before sleeping, re-acquire after waking */

mutex_lock(&dev_mutex);

while (buffer_used == 0) {

    mutex_unlock(&dev_mutex);

    if (wait_event_interruptible(read_queue, buffer_used > 0))

        return -ERESTARTSYS;

    mutex_lock(&dev_mutex);

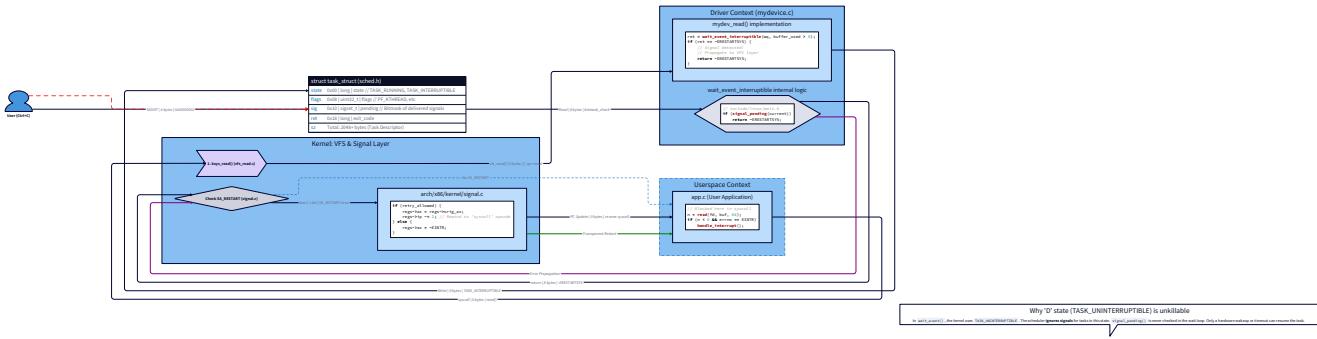
}

/* now holding mutex and buffer_used > 0 */
```

But wait — is there a race between releasing the mutex and adding ourselves to the wait queue? **No** — `wait_event_interruptible` checks the condition *before* sleeping. The sequence is:

1. Check `buffer_used == 0` with mutex held — yes, empty
2. Release mutex
3. `wait_event_interruptible` checks `buffer_used > 0` — still empty (writer hasn't run yet)
4. Adds to wait queue, sets `TASK_INTERRUPTIBLE`, calls `schedule()`
5. Writer acquires mutex, adds data, calls `wake_up_interruptible()`
6. Reader wakes, re-acquires mutex
7. Condition is now true — reader proceeds If the writer runs between steps 2 and 3 (after mutex release but before `wait_event_interruptible`), the writer adds data and calls `wake_up_interruptible()` on an empty wait queue (reader hasn't added itself yet). Then in step 3, `wait_event_interruptible` checks the condition — `buffer_used > 0` is now true — and returns immediately without sleeping. No lost wakeup.

-ERESTARTSYS: The Signal That Travels Through Time

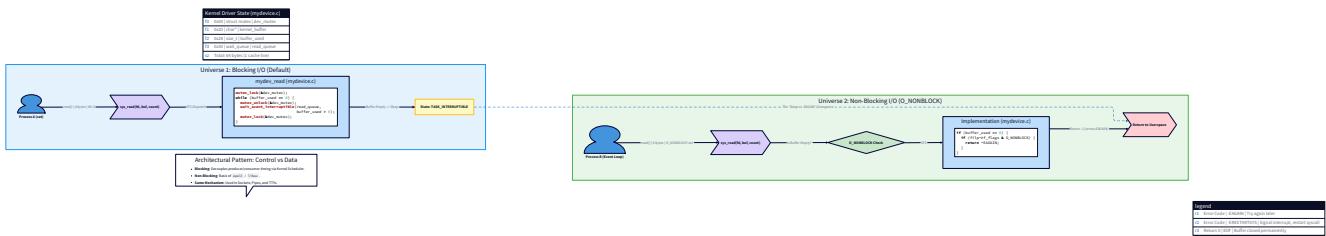


When a user presses `Ctrl+C` while a process is blocked in `read()` on your device, a `SIGINT` signal is delivered. Here is the chain of events:

1. The kernel delivers `SIGINT` to the sleeping process by setting a flag in its `task_struct` and calling `wake_up_process(task)`
2. The process transitions from `TASK_INTERRUPTIBLE` to `TASK_RUNNING`
3. `wait_event_interruptible` detects a pending signal via `signal_pending(current)` and returns `-ERESTARTSYS` instead of `0`
4. Your read handler propagates this: `if (ret == -ERESTARTSYS) return -ERESTARTSYS;`
5. The kernel's signal handling code intercepts the `-ERESTARTSYS` return from your syscall. Here's the subtle part: `-ERESTARTSYS` is **not** `-EINTR`. The kernel uses `-ERESTARTSYS` as an internal signal to its own syscall restart machinery. The kernel then checks the signal handler:
 - If the signal handler was installed with `SA_RESTART` flag: the syscall is **automatically restarted** — the process calls `read()` again transparently, the user process never sees the interruption
 - If the signal handler was installed without `SA_RESTART` (or the signal has default disposition): the kernel converts `-ERESTARTSYS` to `-EINTR` and returns that to the user process, which sees `errno == EINTR`. This is exactly why `read()` sometimes returns `-1` with `errno == EINTR` in userspace programs — the signal handler interrupted the blocking syscall, and the `SA_RESTART` flag was not set. Your `signal-handler` prerequisite project encountered this behavior; now you're seeing the kernel mechanism that produces it.

Why not return `-EINTR` directly from your driver? Because `-EINTR` prevents syscall restart entirely. `-ERESTARTSYS` gives the kernel's signal machinery the opportunity to restart the syscall if `SA_RESTART` allows it. Returning `-ERESTARTSYS` from your handler is the correct, `SA_RESTART`-aware approach. The kernel translates it to `-EINTR` for userspace if restart isn't appropriate.

O_NONBLOCK: Two Universes of I/O Behavior



When a process opens your device with `O_NONBLOCK` (e.g., `open("/dev/mydevice", O_RDWR | O_NONBLOCK)`), it is making a contract with your driver: **"Never block me. If the operation can't complete immediately, return an error."** The flag is stored in `filp->f_flags` and remains for the lifetime of the open file descriptor. You check it in your read and write handlers:

```
/* Check for O_NONBLOCK before sleeping */

if (filp->f_flags & O_NONBLOCK) {

    /* Can't complete immediately - return -EAGAIN */

    mutex_unlock(&dev_mutex);

    return -EAGAIN;
}
```

`-EAGAIN` means "try again" — the operation failed not because of an error but because it would have blocked. The POSIX error is often referred to as `EWOULDBLOCK` (which on Linux has the same value as `EAGAIN`). The caller is expected to either retry later, use `poll()` / `select()` to wait for readiness, or handle the condition gracefully. The complete non-blocking contract for your device:

Condition	Blocking (default)	Non-blocking (O_NONBLOCK)
Read, buffer empty	Sleep until data available	Return -EAGAIN
Read, buffer has data	Return data immediately	Return data immediately
Write, buffer full	Sleep until space available	Return -EAGAIN
Write, buffer has space	Write and return count	Write and return count
O_NONBLOCK is the foundation of event-driven I/O . The <code>select()</code> , <code>poll()</code> , and <code>epoll()</code> system calls use O_NONBLOCK internally — they mark fds as non-blocking, check readiness via your <code>.poll</code> handler, sleep if nothing is ready, and then call <code>read()</code> / <code>write()</code> only when the poll says the operation won't block. If your driver doesn't implement O_NONBLOCK correctly, the entire event-driven I/O machinery breaks for your device.		

The `.poll` File Operation: Not What You Think

The Architectural Insight

The biggest conceptual mistake people make with the kernel's `.poll` handler is thinking it works like this:

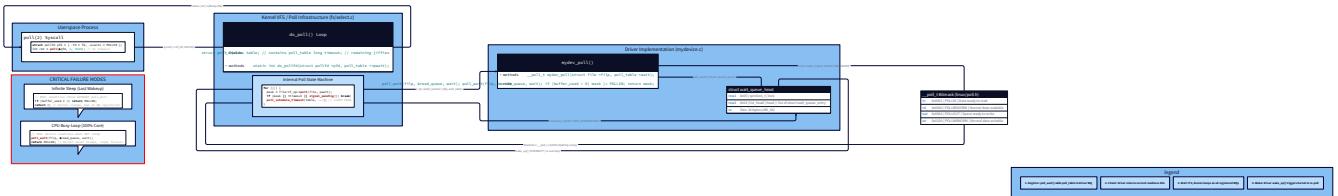
```
// WRONG mental model
poll_handler():
    while (buffer_empty):
        sleep(10ms)
    return POLLIN
```

It actually works like this:

```
// CORRECT mental model
poll_handler():
    register_interest(caller, my_wait_queue) // "wake me if this queue fires"
    return current_readiness_mask           // what's ready RIGHT NOW
```

Your `.poll` handler is not responsible for sleeping. It is responsible for two things:

1. Registering the caller's interest in your wait queue via `poll_wait()`
2. Returning a bitmask of currently-ready events The sleeping happens in the VFS poll infrastructure (in `fs/select.c` for `select / poll`, in `fs/eventpoll.c` for `epoll`). That infrastructure calls your `.poll` handler, inspects the returned mask, and if nothing is ready, sleeps on the wait queues you registered. When your `wake_up_interruptible()` fires (from a write), it wakes the infrastructure, which calls your `.poll` handler again to get the new mask. [[EXPLAIN:poll/select/epoll-kernel-side-protocol|poll/select/epoll kernel-side protocol]]



`poll_wait()` : The Registration Function

```
#include <linux/poll.h>

static __poll_t mydev_poll(struct file *filp, poll_table *wait)
{
    __poll_t mask = 0;

    /*
     * poll_wait(file, wait_queue_head, poll_table)
     *
     * Registers this file's wait queue with the poll_table.
     *
     * When the wait queue is woken (by wake_up_interruptible),
     * the poll infrastructure wakes the sleeping poll/select/epoll
     * and re-checks readiness by calling this .poll handler again.
     *
     * CRITICAL: poll_wait() does NOT sleep. It returns immediately.
     *
     * The sleeping is done by the poll infrastructure, not by your driver.
     */

    poll_wait(filp, &read_queue, wait);
    poll_wait(filp, &write_queue, wait);

    /*
     * After registering interest, report CURRENT readiness.
     *
     * This must be checked after poll_wait() to avoid a race:
     *
     * data might have arrived between the poll() call and poll_wait().
     *
     * If we checked before poll_wait and found empty, then data arrived
     * before the registration, we'd miss the wakeup and wait forever.
     */

    mutex_lock(&dev_mutex);

    if (buffer_used > 0)
```

C

```

        mask |= POLLIN | POLLRDNORM; /* data available for reading */

    if (buffer_used < buffer_size_bytes)

        mask |= POLLOUT | POLLWRNORM; /* space available for writing */

    mutex_unlock(&dev_mutex);

    return mask;

}

```

The Mask Bits: A Quick Reference

Bit	Meaning	When to set
POLLIN	Data available to read	buffer_used > 0
POLLRDNORM	Normal data readable (same as POLLIN for most devices)	Same as POLLIN
POLLOUT	Space available to write	buffer_used < buffer_size_bytes
POLLWRNORM	Normal data writable (same as POLLOUT for most devices)	Same as POLLOUT
POLLERR	Error condition	On error state
POLLHUP	Hangup (device disconnected)	On device removal
POLLPRI	Urgent data (OOB for TCP, etc.)	Rarely needed
Always return both `POLLIN`	POLLRDNORM together and POLLOUT	POLLWRNORM together. Some applications check for POLLRDNORM specifically (POSIX portable code), others check POLLIN. Returning both costs nothing and maximizes compatibility.

The race that poll_wait() prevents: Imagine a sequence without proper registration: (1) you check `buffer_used == 0` — empty. (2) A writer adds data and calls `wake_up_interruptible()`. (3) You call `poll_wait()` — but the wakeup already happened, nobody was on the queue. (4) The poll infrastructure sleeps — forever. By calling `poll_wait()` before checking the condition, you register interest first. If a wakeup fires between the registration and the condition check, the infrastructure sees the wakeup and re-checks. No events can be missed.

The Complete Implementation

Now let's assemble the full driver. This is the complete, production-quality reference combining everything from milestones 2, 3, and 4.

Updated Device State

The device state grows to include the mutex and wait queues:

```
/* mydevice.c - Complete character device with concurrency support */

#include <linux/module.h>

#include <linux/kernel.h>

#include <linux/init.h>

#include <linux/fs.h>

#include <linux/cdev.h>

#include <linux/device.h>

#include <linux/uaccess.h>

#include <linux/slab.h>

#include <linux/atomic.h>

#include <linux/proc_fs.h>

#include <linux/seq_file.h>

#include <linux/mutex.h>          /* mutex_lock, mutex_unlock */

#include <linux/wait.h>           /* wait_event_interruptible, wake_up_interruptible */

#include <linux/poll.h>            /* poll_wait, __poll_t, POLLIN, POLLOUT */

#include <linux/string.h>

#include "mydevice_ioctl.h"

MODULE_LICENSE("GPL");

MODULE_AUTHOR("Your Name <you@example.com>");

MODULE_DESCRIPTION("Character device: mutex protection, blocking I/O, poll support");

MODULE_VERSION("3.0");

/* — Configuration — */

#define DEVICE_NAME    "mydevice"

#define CLASS_NAME     "mydevice_class"

#define PROC_NAME      "mydevice"

#define BUFFER_SIZE    4096

/* — Device State — */
```

C

```
static dev_t           dev_num;

static struct cdev     my_cdev;

static struct class   *my_class;

static struct device  *my_device;

static struct proc_dir_entry *proc_entry;

static char   *kernel_buffer;

static size_t  buffer_size_bytes = BUFFER_SIZE;

static size_t  buffer_used       = 0;

/*
 * Synchronization primitives:
 *
 * dev_mutex:    protects kernel_buffer, buffer_used, buffer_size_bytes.
 *
 *             Must be held for any read or write of these fields.
 *
 * read_queue:   blocking readers sleep here when buffer_used == 0.
 *
 *             Woken by write handler after adding data.
 *
 * write_queue:  blocking writers sleep here when buffer is full.
 *
 *             Woken by read handler after consuming data.
 */

static DEFINE_MUTEX(dev_mutex);

static DECLARE_WAIT_QUEUE_HEAD(read_queue);

static DECLARE_WAIT_QUEUE_HEAD(write_queue);

static atomic_t open_count = ATOMIC_INIT(0);

static atomic_t read_count = ATOMIC_INIT(0);

static atomic_t write_count = ATOMIC_INIT(0);
```

Open and Release

Open and release remain simple — the open count is atomic and needs no mutex:

```
static int mydev_open(struct inode *inode, struct file *filp) C

{
    atomic_inc(&open_count);

    printk(KERN_INFO "mydev: open (count=%d)\n", atomic_read(&open_count));

    return 0;
}

static int mydev_release(struct inode *inode, struct file *filp)

{
    atomic_dec(&open_count);

    printk(KERN_INFO "mydev: release (count=%d)\n", atomic_read(&open_count));

    return 0;
}
```

The Write Handler: Producer with Wake-up

```
static ssize_t mydev_write(struct file *filp, const char __user *buf,
                           size_t count, loff_t *f_pos)

{

    ssize_t ret = 0;

    size_t space_available, bytes_to_copy;

    unsigned long not_copied;

    /*

     * Acquire the mutex before checking or modifying buffer state.

     * Use mutex_lock_interruptible so Ctrl+C can interrupt even the

     * mutex acquisition phase (if the mutex is contended and we're sleeping

     * waiting to acquire it).

    */

    if (mutex_lock_interruptible(&dev_mutex))

        return -ERESTARTSYS;

    /*

     * Wait until there is space in the buffer.

     * This loop handles:

     *   - O_NONBLOCK: return -EAGAIN immediately

     *   - Blocking: sleep until space available or signal arrives

     *   - Spurious wakeups: re-check condition after waking

    */

    while (buffer_used >= buffer_size_bytes) {

        /* Buffer is full */

        if (filp->f_flags & O_NONBLOCK) {

            ret = -EAGAIN;

            goto out_unlock;
        }
    }
}
```

```
}

/*
 * Must release the mutex before sleeping.
 *
 * If we slept holding it, no reader could ever consume data
 * to free space – classic deadlock.
 */

mutex_unlock(&dev_mutex);

/*
 * wait_event_interruptible(queue, condition):
 *
 * Sleep until (buffer_used < buffer_size_bytes) becomes true,
 * or until a signal arrives (returns -ERESTARTSYS).
 *
 * The condition is re-evaluated inside the macro before sleeping –
 * this prevents lost wakeups if a reader freed space between our
 * mutex_unlock and this call.
 */

if (wait_event_interruptible(write_queue,
                             buffer_used < buffer_size_bytes)) {

    return -ERESTARTSYS;
}

/* Re-acquire mutex to re-check condition safely */

if (mutex_lock_interruptible(&dev_mutex))

    return -ERESTARTSYS;

}

/* Now holding mutex, buffer has space */

space_available = buffer_size_bytes - buffer_used;

bytes_to_copy = min(count, space_available);
```

```
not_copied = copy_from_user(kernel_buffer + buffer_used, buf, bytes_to_copy);

if (not_copied) {

    ret = -EFAULT;

    goto out_unlock;

}

buffer_used += bytes_to_copy;

atomic_inc(&write_count);

ret = (ssize_t)bytes_to_copy;

/*
 * Wake up any sleeping readers.
 *
 * wake_up_interruptible() wakes all TASK_INTERRUPTIBLE waiters on
 * read_queue. They will re-check buffer_used > 0 and, if true,
 * proceed to read. If multiple readers wake, only the one that
 * acquires the mutex first will find data; others re-sleep.
 *
 * The wake happens while we still hold dev_mutex. This is intentional
 * and safe: the woken reader cannot acquire dev_mutex until we release
 * it below, so there is no race on buffer_used between the wake and
 * the unlock.
 */

wake_up_interruptible(&read_queue);

out_unlock:

mutex_unlock(&dev_mutex);

return ret;

}
```

The Read Handler: Consumer with Blocking

```
static ssize_t mydev_read(struct file *filp, char __user *buf,
                         size_t count, loff_t *f_pos)

{

    ssize_t ret = 0;

    size_t bytes_available, bytes_to_copy;

    unsigned long not_copied;

    if (mutex_lock_interruptible(&dev_mutex))

        return -ERESTARTSYS;

    /*

     * Wait until there is data to read.

     *

     * NOTE: This driver uses the buffer as a FIFO stream, not a

     * seekable file. The f_pos tracks position within current content.

     * When buffer_used reaches f_pos (all written content consumed),

     * the reader blocks waiting for more content.

     *

     * If the device is used as a pipe (write then read), *f_pos starts

     * at 0 and advances. When *f_pos == buffer_used, we block.

     */

    while (*f_pos >= buffer_used) {

        if (filp->f_flags & O_NONBLOCK) {

            /*

             * Non-blocking: return -EAGAIN (not 0/EOF) when no data is

             * available yet. 0 would mean "end of stream, no more data

             * ever" – which is wrong for a device that might get writes.

             * -EAGAIN means "try again, data might come later."
        }
    }
}
```

```

        */

    ret = -EAGAIN;

    goto out_unlock;

}

mutex_unlock(&dev_mutex);

if (wait_event_interruptible(read_queue, buffer_used > *f_pos)) {

    return -ERESTARTSYS;

}

if (mutex_lock_interruptible(&dev_mutex))

    return -ERESTARTSYS;

}

/* Now holding mutex, buffer has data past *f_pos */

bytes_available = buffer_used - *f_pos;

bytes_to_copy = min(count, bytes_available);

not_copied = copy_to_user(buf, kernel_buffer + *f_pos, bytes_to_copy);

if (not_copied) {

    ret = -EFAULT;

    goto out_unlock;

}

*f_pos += bytes_to_copy;

atomic_inc(&read_count);

ret = (ssize_t)bytes_to_copy;

/*
 * Wake up waiting writers now that we've consumed data and freed space.
 */

wake_up_interruptible(&write_queue);

out_unlock:

```

```
    mutex_unlock(&dev_mutex);

    return ret;

}
```

The Poll Handler

```
static __poll_t mydev_poll(struct file *filp, poll_table *wait) C

{
    __poll_t mask = 0;

    /*
     * Register this file's wait queues with the poll infrastructure.
     *
     * This MUST happen before checking the condition.
     *
     * Sequence with proper registration:
     *
     * 1. poll_wait registers read_queue and write_queue
     *
     * 2. We check conditions and find nothing ready
     *
     * 3. poll infrastructure sleeps on these queues
     *
     * 4. Write handler adds data, calls wake_up_interruptible(read_queue)
     *
     * 5. poll infrastructure wakes, calls .poll again – we return POLLIN
     *
     * Without step 1, step 4's wakeup would have no registered waiter
     *
     * and step 5 would never happen.
     */
    poll_wait(filp, &read_queue, wait);
    poll_wait(filp, &write_queue, wait);

    /*
     * Now check current readiness. Hold the mutex to get a consistent
     * view of buffer_used and buffer_size_bytes.
     *
     * NOTE: We use mutex_lock() here, NOT mutex_lock_interruptible().
     *
     * The .poll handler should not return -ERESTARTSYS or similar – the
     * poll infrastructure does not handle those return codes from .poll.
    }
```

```

* The .poll handler must complete and return a mask.

*
* This means .poll can block briefly on the mutex (until the read or
* write handler finishes its critical section), but this is acceptable
* because the critical section is short.

*/
mutex_lock(&dev_mutex);

if (buffer_used > 0)

    mask |= POLLIN | POLLRDNORM;

if (buffer_used < buffer_size_bytes)

    mask |= POLLOUT | POLLWRNORM;

mutex_unlock(&dev_mutex);

return mask;

}

```

Updated `file_operations` Table

```

static const struct file_operations mydev_fops = {

    .owner          = THIS_MODULE,
    .open           = mydev_open,
    .release        = mydev_release,
    .read           = mydev_read,
    .write          = mydev_write,
    .unlocked_ioctl = mydev_ioctl, /* unchanged from M3 */
    .poll           = mydev_poll,  /* new in M4 */
};

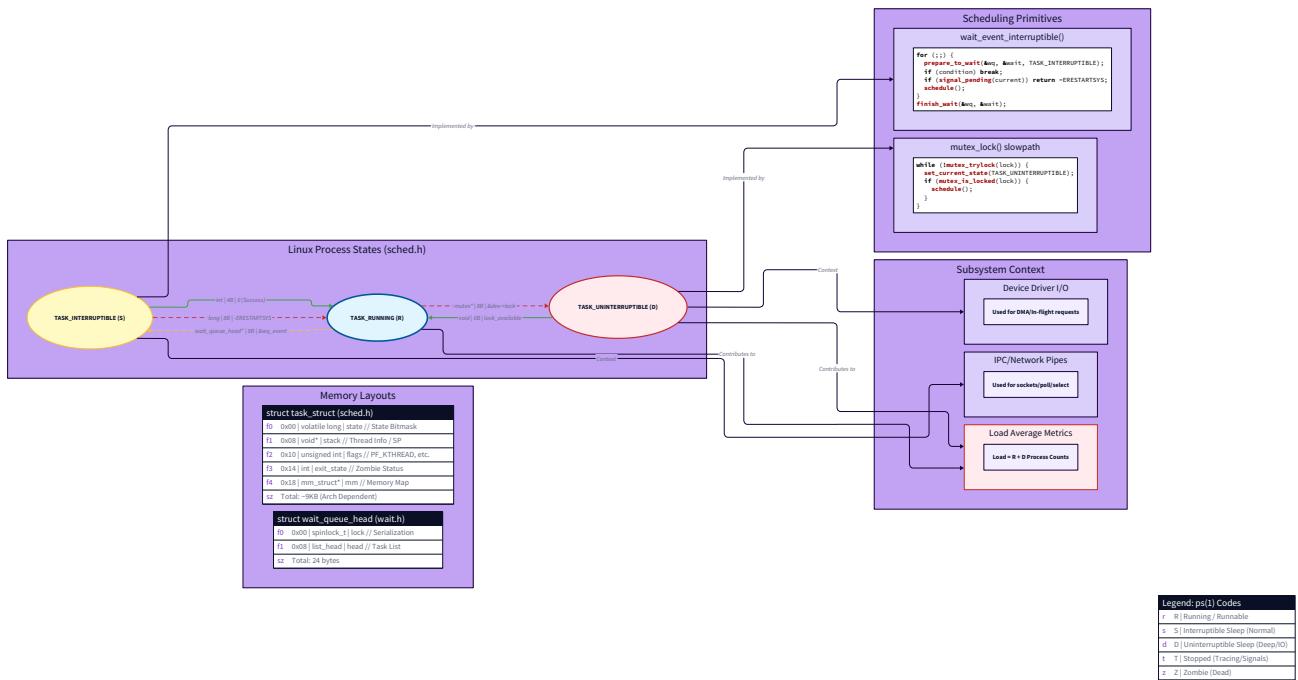
C

```

Init and Exit (Updated)

The init/exit structure from Milestone 3 doesn't change — the mutex and wait queues are statically initialized and don't need explicit init/cleanup. `DEFINE_MUTEX` and `DECLARE_WAIT_QUEUE_HEAD` do all the initialization at compile time.

Process States and the Scheduler: The Big Picture



Understanding what `wait_event_interruptible` does to the Linux scheduler demystifies a lot of "why does my server behave this way" questions. Linux processes have several states visible in `ps` and `top`:

State	ps Code	Meaning	Kernel State
Running	R	On a CPU or runnable, on the run queue	TASK_RUNNING
Interruptible sleep	S	Sleeping, can be woken by signals	TASK_INTERRUPTIBLE
Uninterruptible sleep	D	Sleeping, signals can NOT wake it	TASK_UNINTERRUPTIBLE
Stopped	T	SIGSTOP or ptrace	TASK_STOPPED
Zombie	Z	Exited, waiting for parent to reap	EXIT_ZOMBIE
When your read handler calls <code>wait_event_interruptible()</code> , the process transitions to S (interruptible sleep) . This is the normal, healthy blocking state — <code>Ctrl+C</code> can wake it.			
The D state (uninterruptible sleep) is used for I/O waits that cannot be interrupted mid-operation — for example, waiting for a disk read to complete. Processes in D state do not respond to <code>SIGKILL</code> . This is why a hung NFS mount can produce unkillable processes: they're in D state waiting for network I/O that never completes.			
This directly explains Linux load average behavior: Linux load average counts <i>both</i> R (running/runnable) and D (uninterruptible sleep) processes. A server with 8 CPUs can have a load average of 64 if 64 processes are simultaneously blocked in D state on disk I/O — all 64 are "waiting for the CPU to be able to do useful work for them," even though none are consuming CPU cycles right			

State	ps Code	Meaning	Kernel State
now. A database that shows 0% CPU but load average of 20 is I/O-bound, not CPU-bound. Your <code>wait_event_interruptible</code> produces <code>S</code> state processes — they contribute 0 to load average when sleeping. <code>wait_event</code> (uninterruptible variant) produces <code>D</code> state — it contributes to load average.			

The Thundering Herd Problem

Imagine 100 reader processes are sleeping on `read_queue`, all waiting for data. One writer calls `write()` and adds 100 bytes. The writer calls `wake_up_interruptible(&read_queue)`. `wake_up_interruptible` wakes all sleeping processes on the queue. All 100 readers transition to `TASK_RUNNING` and compete to acquire `dev_mutex`. The one that wins reads all 100 bytes and releases the mutex. The other 99 acquire the mutex in turn, find `buffer_used == 0`, and go back to sleep. This is the thundering herd problem: a single wakeup event causes O(N) processes to wake, contend, and re-sleep. For N=100 with a fast mutex, this is 99 unnecessary context switches and 99 unnecessary mutex acquisitions. The kernel provides `wake_up_interruptible_nr(queue, nr)` to wake at most `nr` waiters, and `WQ_FLAG_EXCLUSIVE` to designate certain waiters as "exclusive" — `wake_up_interruptible` wakes exclusive waiters one at a time. For your single-buffer device, this optimization isn't necessary — the buffer either has data for all readers (stream device) or for one reader (pipe-like device). But knowing this exists tells you why `nginx`'s `accept_mutex` exists (to prevent all workers from thundering on a new connection) and why `epoll` with `EPOLLET | EPOLLONESHOT` is useful for high-connection servers. For your driver: since the buffer can hold multiple bytes and

multiple readers can each read different portions, `wake_up_interruptible` waking all readers is actually correct behavior — each reader that wakes finds some data available. The "unnecessary wakeup" scenario only occurs when readers compete for the same bytes, which doesn't happen with your `f_pos-per-file` design.

Updating ioctl for Thread Safety

The ioctl handler from Milestone 3 accesses `kernel_buffer`, `buffer_used`, and `buffer_size_bytes` without holding the mutex. Now that concurrent access is possible, every access to shared state needs the lock:

```
static long mydev_ioctl(struct file *filp, unsigned int cmd, unsigned long arg) C

{
    int ret = 0;

    if (_IOC_TYPE(cmd) != MYDEV_MAGIC)

        return -ENOTTY;

    if (_IOC_NR(cmd) > 2)

        return -ENOTTY;

    switch (cmd) {

        case MYDEV_CLEAR:

            if (mutex_lock_interruptible(&dev_mutex))

                return -ERESTARTSYS;

            memset(kernel_buffer, 0, buffer_size_bytes);

            buffer_used = 0;

            /*

             * After clearing, wake writers: the buffer now has maximum space.

             * No need to wake readers: clear reduces bytes_used, not increases it.

            */

            wake_up_interruptible(&write_queue);

            mutex_unlock(&dev_mutex);

            break;

        case MYDEV_RESIZE: {

            __u32 new_size;

            char *new_buffer;

            if (copy_from_user(&new_size, (__u32 __user *)arg, sizeof(__u32)))

                return -EFAULT;

            if (new_size == 0 || new_size > (1U << 20))

                return -EINVAL;
    }
}
```

```

new_buffer = kzalloc(new_size, GFP_KERNEL);

if (!new_buffer)

    return -ENOMEM;

if (mutex_lock_interruptible(&dev_mutex)) {

    kfree(new_buffer);

    return -ERESTARTSYS;

}

if (buffer_used > new_size)

    buffer_used = new_size;

memcpy(new_buffer, kernel_buffer, buffer_used);

kfree(kernel_buffer);

kernel_buffer = new_buffer;

buffer_size_bytes = new_size;

/*
 * A resize might free space (if shrinking, content was truncated)
 * or add space (if growing). Wake both queues to let waiters
 * re-evaluate their conditions.
 */

wake_up_interruptible(&read_queue);

wake_up_interruptible(&write_queue);

mutex_unlock(&dev_mutex);

break;

}

case MYDEV_STATUS: {

    struct mydev_status status;

    if (mutex_lock_interruptible(&dev_mutex))

        return -ERESTARTSYS;

```

```

status.buffer_size = (__u32)buffer_size_bytes;

status.bytes_used = (__u32)buffer_used;

status.open_count = (__u32)atomic_read(&open_count);

status.read_count = (__u32)atomic_read(&read_count);

status.write_count = (__u32)atomic_read(&write_count);

status._reserved = 0;

mutex_unlock(&dev_mutex);

if (copy_to_user((struct mydev_status __user *)arg,
                 &status, sizeof(status)))

    return -EFAULT;

break;

}

default:

    return -ENOTTY;

}

return ret;
}

```

Note the pattern for `MYDEV_STATUS` : we copy the data under the mutex, then call `copy_to_user` *outside* the mutex. Why? Because `copy_to_user` can sleep (if the user page is not present and needs to be faulted in). Sleeping while holding a mutex is legal in process context, but it holds the mutex longer than necessary, increasing contention. Copy to a local

stack struct under the mutex, then copy to userspace without the mutex — this is the canonical pattern.

Userspace Poll Test

Let's write a userspace test that exercises the poll functionality:

```
/* poll_test.c - Test poll() on /dev/mydevice */

#include <stdio.h>

#include <stdlib.h>

#include <fcntl.h>

#include <unistd.h>

#include <poll.h>

#include <string.h>

#include <errno.h>

#define DEVICE "/dev/mydevice"

int main(void)

{

    int fd;

    struct pollfd pfd;

    int ret;

    char buf[64];

    fd = open(DEVICE, O_RDWR | O_NONBLOCK);

    if (fd < 0) { perror("open"); return 1; }

    pfd.fd      = fd;

    pfd.events  = POLLIN | POLLOUT;

    /* — Test 1: Initial state – should be writable, not readable — */

    ret = poll(&pfd, 1, 0);    /* timeout=0: immediate check, no sleep */

    printf("Initial poll: revents=0x%llx\n", pfd.revents);

    if (pfd.revents & POLLOUT) printf("  PASS: POLLOUT set (buffer has space)\n");

    if (!(pfd.revents & POLLIN)) printf("  PASS: POLLIN clear (buffer is empty)\n");

    /* — Test 2: Write data, then poll – should become readable — */

    write(fd, "hello world", 11);

    ret = poll(&pfd, 1, 0);
```

C

```

printf("After write poll: revents=0x%x\n", pfd.revents);

if (pfd.revents & POLLIN) printf(" PASS: POLLIN set (data available)\n");

if (pfd.revents & POLLOUT) printf(" PASS: POLLOUT set (still has space)\n");

/* — Test 3: Non-blocking read ————— */

ret = read(fd, buf, sizeof(buf));

if (ret > 0) {

    buf[ret] = '\0';

    printf("Read %d bytes: \"%s\"\n", ret, buf);

}

/* — Test 4: poll with timeout — wait for data from another process */

printf("Waiting for data (5s timeout)...");

printf("Run in another terminal: echo test | sudo tee /dev/mydevice\n");

pfd.events = POLLIN;

ret = poll(&pfd, 1, 5000); /* 5 second timeout */

if (ret == 0) {

    printf("Timed out - no data arrived\n");

} else if (ret > 0 && (pfd.revents & POLLIN)) {

    ret = read(fd, buf, sizeof(buf));

    buf[ret > 0 ? ret : 0] = '\0';

    printf("Data arrived: \"%s\"\n", buf);

}

close(fd);

return 0;

}

```

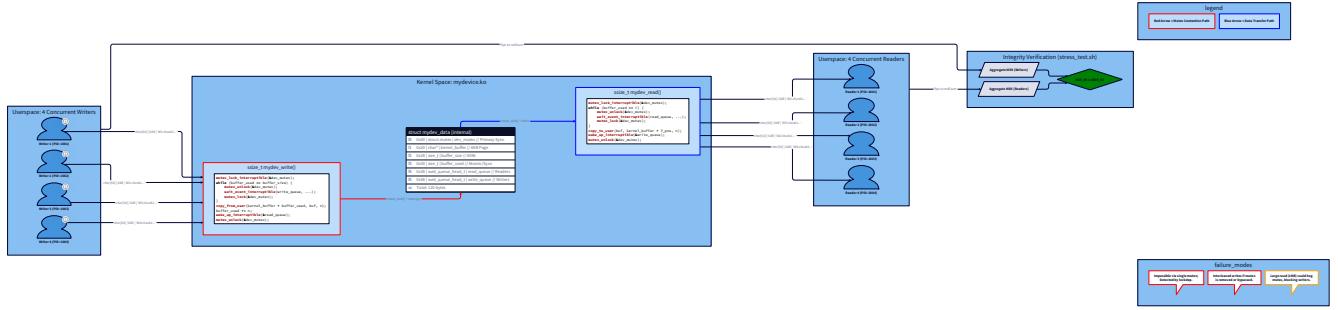
Compile and run:

```
gcc -Wall -Werror -o poll_test poll_test.c
```

BASH

```
sudo ./poll_test
```

The Stress Test: 4 Writers + 4 Readers



The acceptance criteria require a stress test with 4 concurrent writer processes and 4 concurrent reader processes, with data integrity verified by checksums. Here is a complete test harness:

BASH

```
#!/bin/bash

# stress_test.sh – Concurrent stress test for mydevice

#
# Architecture:

#   4 writer processes: each writes N chunks of fixed-size data

#   4 reader processes: each reads until it has collected N*chunk_size bytes total

#
# Data integrity: each writer generates deterministic data (byte pattern
# based on writer ID + sequence number). All written bytes are checksummed
# with CRC or sum. All read bytes are checksummed. Totals must match.

#
# Limitations of this simple test: because multiple writers interleave their
# writes, the ORDER of bytes in the buffer is non-deterministic. This test
# verifies that no bytes are lost or corrupted (checksum of all written
# data == checksum of all read data), not that they arrive in order.

set -e

DEVICE="/dev/mydevice"

NUM_WRITERS=4

NUM_READERS=4

WRITES_PER_WRITER=100

CHUNK_SIZE=64          # 64 bytes per write – fits in buffer without constant blocking

WRITE_DIR=$(mktemp -d)

READ_DIR=$(mktemp -d)

echo "[stress] Device: $DEVICE"

echo "[stress] Writers: $NUM_WRITERS, Readers: $NUM_READERS"

echo "[stress] Each writer: $WRITES_PER_WRITER writes × $CHUNK_SIZE bytes"

TOTAL_BYTES=$((NUM_WRITERS * WRITES_PER_WRITER * CHUNK_SIZE))
```

```

echo "[stress] Total bytes: $TOTAL_BYTES"

# — Writer function ——————  

writer() {  

    local id=$1  

    local outfile="$WRITE_DIR/writer_${id}.dat"  

    # Generate deterministic data: writer_id repeated as bytes  

    # Each writer uses a unique byte value (0x41='A', 0x42='B', etc.)  

    local byte_val=$((0x41 + id))    # 'A', 'B', 'C', 'D'  

    local byte_char  

    byte_char=$(printf "\\$(printf '%03o' $byte_val)")  

    for i in $(seq 1 ${WRITES_PER_WRITER}); do  

        # Create a chunk: CHUNK_SIZE bytes all equal to byte_val  

        printf "${byte_char}%.0s" $(seq 1 ${CHUNK_SIZE}) | \  

            sudo tee -a "$outfile" >> /dev/null  

        printf "${byte_char}%.0s" $(seq 1 ${CHUNK_SIZE}) | \  

            sudo tee -a "$DEVICE" > /dev/null  

        # Small delay to increase interleaving  

        sleep 0.$((RANDOM % 5))  

    done  

    echo "[writer $id] done (${WRITES_PER_WRITER} writes × ${CHUNK_SIZE} bytes)"  

}  

# — Reader function ——————  

reader() {  

    local id=$1  

    local outfile="$READ_DIR/reader_${id}.dat"  

    local bytes_per_reader=$((TOTAL_BYTES / NUM_READERS))  

    local collected=0

```

```
while [ $collected -lt $bytes_per_reader ]; do

    # Read whatever is available (blocks until data present)

    local chunk

    chunk=$(sudo dd if="$DEVICE" bs=64 count=1 2>/dev/null)

    if [ -n "$chunk" ]; then

        printf "%s" "$chunk" >> "$outfile"

        collected=$((collected + ${#chunk}))

    fi

done

echo "[reader $id] done (${collected} bytes collected)"

}

# — Run stress test ——————  
echo "[stress] Starting writers and readers..."  
  
# Start readers first (they'll block waiting for data)  
  
for i in $(seq 0 $((NUM_READERS-1))); do  
  
    reader $i &  
  
done  
  
# Start writers  
  
for i in $(seq 0 $((NUM_WRITERS-1))); do  
  
    writer $i &  
  
done  
  
# Wait for all background jobs  
  
wait  
  
echo "[stress] All processes completed"  
  
# — Checksum verification ——————  
echo "[stress] Computing checksums..."  
  
# Checksum of all written data (concatenate all writer output files)
```

```

WRITE_CHECKSUM=$(cat "$WRITE_DIR"/writer_*.dat | md5sum | cut -d' ' -f1)

READ_CHECKSUM=$(cat "$READ_DIR"/reader_*.dat | md5sum | cut -d' ' -f1)

echo "[stress] Written data MD5: $WRITE_CHECKSUM"

echo "[stress] Read data MD5:      $READ_CHECKSUM"

WRITE_TOTAL=$(cat "$WRITE_DIR"/writer_*.dat | wc -c)

READ_TOTAL=$(cat "$READ_DIR"/reader_*.dat | wc -c)

echo "[stress] Written bytes: $WRITE_TOTAL (expected $TOTAL_BYTES)"

echo "[stress] Read bytes:      $READ_TOTAL (expected $TOTAL_BYTES)"

if [ "$WRITE_TOTAL" -eq "$READ_TOTAL" ]; then

    echo "[stress] PASS: byte counts match"

else

    echo "[stress] FAIL: byte count mismatch!"

    rm -rf "$WRITE_DIR" "$READ_DIR"

    exit 1

fi

# Note: due to interleaving, the exact byte ordering may differ between

# written and read streams. The meaningful verification is:

# 1. No bytes are lost: total written == total read

# 2. No kernel oops, warnings, or panics in dmesg during the test

# 3. dmesg shows no mutex warnings, deadlock warnings, or data corruption

echo "[stress] Checking dmesg for kernel warnings..."

if dmesg | tail -50 | grep -iE "(BUG|WARNING|OOPS|panic|deadlock|corruption)" ; then

    echo "[stress] FAIL: kernel warnings detected"

    exit 1

else

    echo "[stress] PASS: no kernel warnings"

fi

```

```
rm -rf "$WRITE_DIR" "$READ_DIR"  
echo "[stress] All tests PASSED"
```

For a more rigorous byte-level integrity check (verifying that every specific byte written by writer N is received by some reader, with no mutations), use this C test program:

```
/* stress_c.c - C-based stress test with per-byte verification */

#include <stdio.h>

#include <stdlib.h>

#include <string.h>

#include <fcntl.h>

#include <unistd.h>

#include <pthread.h>

#include <stdint.h>

#include <errno.h>

#include <sys/types.h>

#define DEVICE          "/dev/mydevice"

#define NUM_WRITERS    4

#define NUM_READERS    4

#define WRITES_PER_W   200

#define CHUNK_SIZE     32

/*
 * Verification strategy:
 *
 * Each writer writes a recognizable 4-byte header (0xWW + seq number)
 * followed by (CHUNK_SIZE - 4) bytes of the writer's fill byte.
 *
 * Readers collect all bytes and verify structure.
 *
 * Simpler approach: each writer writes N*CHUNK_SIZE bytes of a unique
 * fill value. Total written per writer-value is tracked. Total read
 * per value must match. This proves no bytes are added or lost.
 */

static uint64_t bytes_written[NUM_WRITERS] = {0};

static uint64_t bytes_read_by_value[256] = {0};
```

C

```

static pthread_mutex_t stats_mutex          = PTHREAD_MUTEX_INITIALIZER;

static void *writer_thread(void *arg)
{

    int id = *(int*)arg;

    int fd = open(DEVICE, O_WRONLY);

    if (fd < 0) { perror("open writer"); return NULL; }

    uint8_t fill = (uint8_t)('A' + id); /* each writer uses unique byte */

    uint8_t chunk[CHUNK_SIZE];

    memset(chunk, fill, CHUNK_SIZE);

    uint64_t written = 0;

    for (int i = 0; i < WRITES_PER_W; i++) {

        ssize_t n;

        do {

            n = write(fd, chunk + (CHUNK_SIZE - (written % CHUNK_SIZE)),

                      CHUNK_SIZE);

            if (n < 0 && errno != EAGAIN) { perror("write"); goto done; }

        } while (n < 0);

        written += n;
    }
}

done:

close(fd);

pthread_mutex_lock(&stats_mutex);

bytes_written[id] = written;

pthread_mutex_unlock(&stats_mutex);

printf("Writer %d done: %lu bytes of 0x%02X\n", id, (unsigned long)written, fill);

return NULL;
}

```

```
static void *reader_thread(void *arg)
{
    (void)arg;

    int fd = open(DEVICE, O_RDONLY);

    if (fd < 0) { perror("open reader"); return NULL; }

    uint64_t target = (uint64_t)NUM_WRITERS * WRITES_PER_W * CHUNK_SIZE / NUM_READERS;

    uint64_t collected = 0;

    uint8_t buf[256];

    uint64_t local_count[256] = {0};

    while (collected < target) {

        ssize_t n = read(fd, buf, sizeof(buf));

        if (n < 0) {

            if (errno == EINTR) continue;

            perror("read"); break;

        }

        for (ssize_t i = 0; i < n; i++)
            local_count[buf[i]]++;

        collected += n;
    }

    close(fd);

    pthread_mutex_lock(&stats_mutex);

    for (int i = 0; i < 256; i++)
        bytes_read_by_value[i] += local_count[i];

    pthread_mutex_unlock(&stats_mutex);
}

return NULL;
}

int main(void)
```

```
{\n\n    pthread_t writers[NUM_WRITERS], readers[NUM_READERS];\n\n    int ids[NUM_WRITERS];\n\n    printf("Stress test: %d writers, %d readers, %d writes each of %d bytes\\n",\n\n        NUM_WRITERS, NUM_READERS, WRITES_PER_W, CHUNK_SIZE);\n\n    for (int i = 0; i < NUM_READERS; i++)\n\n        pthread_create(&readers[i], NULL, reader_thread, NULL);\n\n    for (int i = 0; i < NUM_WRITERS; i++) {\n\n        ids[i] = i;\n\n        pthread_create(&writers[i], NULL, writer_thread, &ids[i]);\n\n    }\n\n    for (int i = 0; i < NUM_WRITERS; i++) pthread_join(writers[i], NULL);\n\n    for (int i = 0; i < NUM_READERS; i++) pthread_join(readers[i], NULL);\n\n    printf("\\n==== Verification ===\\n");\n\n    int pass = 1;\n\n    for (int i = 0; i < NUM_WRITERS; i++) {\n\n        uint8_t fill = (uint8_t)('A' + i);\n\n        printf("Writer %d (0x%02X '%c'): wrote %lu, read %lu - %s\\n",\n\n            i, fill, (char)fill,\n\n            (unsigned long)bytes_written[i],\n\n            (unsigned long)bytes_read_by_value[fill],\n\n            bytes_written[i] == bytes_read_by_value[fill] ? "PASS" : "FAIL");\n\n        if (bytes_written[i] != bytes_read_by_value[fill]) pass = 0;\n\n    }\n\n    printf("\\n%\\s\\n", pass ? "ALL CHECKS PASSED" : "SOME CHECKS FAILED");\n\n    return pass ? 0 : 1;\n}
```

```
}
```

Compile and run:

```
gcc -Wall -Werror -pthread -o stress_c stress_c.c  
sudo insmod mydevice.ko  
sudo chmod 666 /dev/mydevice    # or run as root  
sudo ./stress_c  
  
# Watch for kernel panics/oops during the test:  
dmesg -w &
```

BASH

Three-Level View: Blocking Read in Depth

When process A calls `read()` on your empty device and blocks, while process B writes data: Level 1 — Applications: Process A calls `read(fd, buf, 100)` — blocks. Process B calls `write(fd, data, 50)` — completes immediately. Process A's `read()` returns 50 bytes. Level 2 — Kernel / Scheduler: Process A calls `sys_read` → `mydev_read` → finds `buffer_used == 0` → releases `dev_mutex` → calls `wait_event_interruptible(read_queue, ...)` → sets state `TASK_INTERRUPTIBLE` → calls `schedule()`. The scheduler removes A from the run queue and runs another task. Process B calls `sys_write` → `mydev_write` → acquires `dev_mutex` → copies 50 bytes → sets `buffer_used = 50` → calls `wake_up_interruptible(&read_queue)`. This sets A's state back to `TASK_RUNNING` and places A on the run queue. B releases the mutex and returns. The scheduler eventually selects A to run. A resumes inside `wait_event_interruptible`, re-acquires `dev_mutex`, re-checks `buffer_used > 0` (now true), exits the wait loop, copies 50 bytes to userspace via `copy_to_user`, returns 50. Level 3 — Hardware: When `schedule()` is called, the CPU performs a context switch: saves the current register state

(instruction pointer, stack pointer, general-purpose registers, FPU state if dirty) to process A's `thread_struct` in its `task_struct`. Loads process B's saved register state. Switches the page table base register (`CR3` on x86_64) to B's page directory. The TLB is largely flushed (with PCID optimization, partially invalidated). Process B's code now runs from where it last left off. The `wake_up_interruptible()` call generates a store-release memory barrier before adding A back to the run queue, ensuring the `buffer_used = 50` write is visible to all CPUs before A might observe it. On x86, this is a `LOCK` prefixed instruction; on ARM64, it's an explicit `stlr` (store-release) or `dmb` (data memory barrier). Without this, A might wake up and read a stale `buffer_used` value from its cache.

Hardware Soul: The Cache Line Dance

Every time your write handler executes `buffer_used += bytes_to_copy` and then calls `wake_up_interruptible`, a specific sequence of cache events fires: `buffer_used access` (a `size_t` at some address in BSS):

- Write handler running on CPU 0 has the cache line containing `buffer_used` in **M (Modified) state** in CPU 0's L1 cache
- If CPU 1 is simultaneously in the read handler checking `buffer_used`, the MESI protocol requires CPU 0's L1 to supply the value to CPU 1 via the L3 (shared last-level cache) or inter-processor communication. This cache line invalidation and transfer costs ~40-100 ns on modern x86 — cheap for one access, expensive if done millions of times per second `kernel_buffer writes` (the actual data bytes):
- Sequential writes to `kernel_buffer + buffer_used` at offsets 0, 32, 64, etc. are perfectly cache-friendly — the hardware prefetcher detects the sequential access pattern and brings cache lines ahead of the write pointer. With a 4KB buffer, the entire buffer fits in L1 cache (typically 32KB, so 1/8 full)
- `copy_from_user` on x86_64 uses `rep movsb` for this size, which the hardware optimizes as a fast string copy with near-peak memory bandwidth `wake_up_interruptible` accesses the wait queue head (`read_queue`):
- The spinlock inside the wait queue head must be acquired. On a multicore system where many processes might be calling `read()` simultaneously, this spinlock is a potential contention point. The lock acquisition uses `lock xchg` or `lock cmpxchg`, requiring exclusive cache line ownership — this bounces the cache line between cores if multiple cores try to acquire it simultaneously

- For your driver at moderate concurrency (4-8 processes), this is negligible. At thousands of operations per second on 64 cores, this would be the bottleneck **Context switches** (when a reader goes to sleep and a writer wakes it):
- A context switch saves/restores ~200 bytes of register state to the `task_struct`. The `task_struct` itself is typically in L2 or L3 cache for recently-active tasks. A full context switch costs ~2-4 µs on modern hardware, dominated by cache pollution from loading the new process's working set, TLB refill (even with PCID, some entries are invalidated), and pipeline flushing

Common Pitfalls in Concurrent Kernel Code

Holding a Mutex While Sleeping in `wait_event_interruptible`

```
/* DEADLY DEADLOCK – do not do this */

mutex_lock(&dev_mutex);

wait_event_interruptible(read_queue, buffer_used > 0);

/* The writer needs dev_mutex to update buffer_used.

   The writer calls mutex_lock(&dev_mutex) – blocks.

   Reader waits for buffer_used to change – which requires writer to run.

   Writer waits for mutex – which reader holds.

   Classic circular wait = deadlock. Machine hangs. */
```

The fix is always: release the mutex before waiting, reacquire after.

Not Re-checking the Condition After Waking (Manual Wait Loop)

If you implement a manual wait loop instead of `wait_event_interruptible`:

```

/* BUGGY - spurious wakeup not handled */

mutex_unlock(&dev_mutex);

wait_event_interruptible(read_queue, buffer_used > 0);

mutex_lock(&dev_mutex);

/* WARNING: another reader might have consumed all data between
   the wakeup and our mutex_lock. buffer_used might now be 0 again.

   We must re-check! */

/* CORRECT */

mutex_unlock(&dev_mutex);

if (wait_event_interruptible(read_queue, buffer_used > 0))

    return -ERESTARTSYS;

mutex_lock(&dev_mutex);

/* wait_event_interruptible guarantees the condition is true when it
   returns 0, because it re-checks under the wait queue's internal lock. */

```

Returning 0 Instead of -EAGAIN for Non-Blocking Empty Read

```

/* WRONG - 0 signals EOF; programs like cat will exit normally */

if (filp->f_flags & O_NONBLOCK)

    return 0; /* wrong! cat thinks the "file" is done */

/* CORRECT */

if (filp->f_flags & O_NONBLOCK)

    return -EAGAIN; /* "no data now, try later" */

```

0 from `read()` means **end of file** — there will never be more data. `-EAGAIN` (as `errno = EAGAIN`) means "no data now, but there might be later." These have completely different semantics to userspace programs. Returning 0 for an empty device means `cat` exits, poll-based event loops deregister the fd, and `read()` loops terminate. `-EAGAIN` means "keep the fd alive, check again."

Calling `poll_wait()` After the Condition Check

```
/* WRONG – race window between condition check and registration */

mutex_lock(&dev_mutex);

if (buffer_used > 0)

    mask |= POLLIN;

mutex_unlock(&dev_mutex);

poll_wait(filp, &read_queue, wait); /* ← too late! wakeup could have been missed */

return mask;

/* CORRECT – register first, then check */

poll_wait(filp, &read_queue, wait); /* register interest first */

mutex_lock(&dev_mutex);

if (buffer_used > 0)

    mask |= POLLIN;

mutex_unlock(&dev_mutex);

return mask;
```

Not Using `mutex_lock()` (Non-Interruptible) in `.poll`

The `.poll` handler must return a mask — it cannot return an error code like `-ERESTARTSYS`. If you use `mutex_lock_interruptible()` in `.poll` and a signal arrives, you can't correctly signal this to the caller. Use `mutex_lock()` (non-interruptible) in `.poll`. The critical section in `.poll` is intentionally short (just reading two variables), so the wait should be microseconds at most.

Forgetting to Wake `write_queue` After Read

```
/* INCOMPLETE – writers sleeping on a full buffer never wake after a read */

static ssize_t mydev_read(...)

{
    /* ... read bytes ... */

    atomic_inc(&read_count);

    mutex_unlock(&dev_mutex);

    return bytes_to_copy;

    /* MISSING: wake_up_interruptible(&write_queue) */
}

C
```

If the buffer was full before the read, writers are sleeping on `write_queue`. After the read frees space, you must call `wake_up_interruptible(&write_queue)` or writers remain asleep indefinitely. Always wake the opposite queue after modifying `buffer_used`.

Knowledge Cascade: What You've Just Unlocked

→ Process States Explain Load Average (cross-domain: system administration) Now you understand that `wait_event_interruptible` puts processes in `TASK_INTERRUPTIBLE` (`S` state), while `wait_event` (uninterruptible) uses `TASK_UNINTERRUPTIBLE` (`D` state). Linux's load average counts both `R` and `D` processes — which is why a disk-bound server shows high load average at 0% CPU. When `top` shows many `D` processes, they're waiting in uninterruptible sleep (typically for disk or network I/O). When `top` shows many `S` processes (like your readers sleeping on an empty device), they contribute 0 to load average because they're interruptibly

sleeping. This is a subtle but important distinction: a server running your blocking-read device will show `S`-state processes during blocking reads, not `D`-state, so it won't inflate load average. A disk driver using `wait_event` (uninterruptible) *will* inflate load average. → `-ERESTARTSYS`

Connects to the signal-handler Prerequisite (same domain: kernel signal path) In your `signal-handler` prerequisite project, you encountered `errno == EINTR` when a signal interrupted a syscall. Now you know the mechanism: your driver returned `-ERESTARTSYS`, the kernel's syscall exit path detected a pending signal, and if `SA_RESTART` wasn't set on the signal handler, it converted `-ERESTARTSYS` to `-EINTR` before returning to userspace. This is why `EINTR`-safe loops in C (wrapping `read()` in a `do { n = read(...); } while (n < 0 && errno == EINTR)`) are necessary: your driver is correctly propagating signal interruption, and userspace needs to decide whether to restart or handle it. → Your `.poll` Handler Is the Foundation of Every Event Loop (cross-domain: Node.js, nginx, Go)

The `poll_wait()` registration mechanism you just implemented is exactly what `epoll` uses internally. When an epoll fd monitors your device with `epoll_ctl(epfd, EPOLL_CTL_ADD, devfd, ...)`, the kernel calls your `.poll` handler via `ep_item_poll()`. Your `poll_wait()` registers epoll's internal wait queue callback with `read_queue` and `write_queue`. When your write handler calls `wake_up_interruptible(&read_queue)`, it fires the epoll callback, which wakes the `epoll_wait()` call, which returns the device fd as ready. This is the O(1) event notification that makes `epoll` scale to millions of file descriptors: no polling loop, just wait queue callbacks. Node.js's `libuv`, Go's `netpoller`, and nginx all sit on top of this exact mechanism. The `.poll` handler you just wrote is the leaf node in a chain that enables all modern high-performance I/O. → Your Driver Is Now a Go Channel (cross-domain: concurrent programming) The bounded buffer you just

built — with a mutex protecting shared state, blocking producers when full, blocking consumers when empty, and waking the other side on state change — is the exact implementation of a Go buffered channel (`(make(chan T, N))`). The Go runtime implements channels with a mutex, a wait queue for senders (your `write_queue`), and a wait queue for receivers (your `read_queue`). Java's `ArrayBlockingQueue`, Python's `queue.Queue`, and POSIX pipes all implement this same producer-consumer pattern.

When you next read Go's channel implementation in `runtime/chan.go`, you'll see `sudog` (waiting goroutine descriptors) in a send/receive queue — the Go-level equivalent of your `wait_queue_entry_t` instances in `read_queue` and `write_queue`. The kernel C implementation you just wrote is the conceptual ancestor of all of them. → Thundering Herd Explains nginx's `accept_mutex` and Linux 4.5's `EPOLLEXCLUSIVE` (cross-domain: web servers)

The thundering herd you encountered with `wake_up_interruptible` waking all readers is the exact problem that afflicted early multi-process web servers: when a new TCP connection arrives, all worker processes (sleeping in `accept()`) wake up, but only one gets the connection. The others re-sleep after doing useless work. nginx added `accept_mutex` to serialize connection acceptance. Linux 4.5 added `EPOLLEXCLUSIVE` to epoll — when multiple processes epoll-wait on the same fd with `EPOLLEXCLUSIVE`, only one is woken per event. For your device, `wake_up_interruptible_nr(&read_queue, 1)` would wake exactly one reader — appropriate if the device is pipe-like (one reader gets all the bytes). `wake_up_interruptible` (all waiters) is appropriate if the device is broadcast-like. The choice of which wake mechanism to use is an architectural decision about your device's concurrency semantics. → Kernel Synchronization Patterns → Database Locking (cross-domain: databases)

The pattern you just implemented — "acquire lock, check

condition, release lock, sleep, re-acquire lock, re-check condition" — is identical to how database engines implement row-level blocking waits. PostgreSQL's `LockAcquire()` function uses a heavyweight lock mechanism with a wait list; `LockRelease()` calls `ProcWakeup()` to wake a waiting process — exactly your `wake_up_interruptible` after modifying `buffer_used`. MySQL's InnoDB uses a similar structure with `lock_wait_suspend_thread()` corresponding to your `wait_event_interruptible`. The concurrency primitives you just implemented in 200 lines of C are the conceptual foundation of MVCC, row locking, and deadlock detection in every major database.

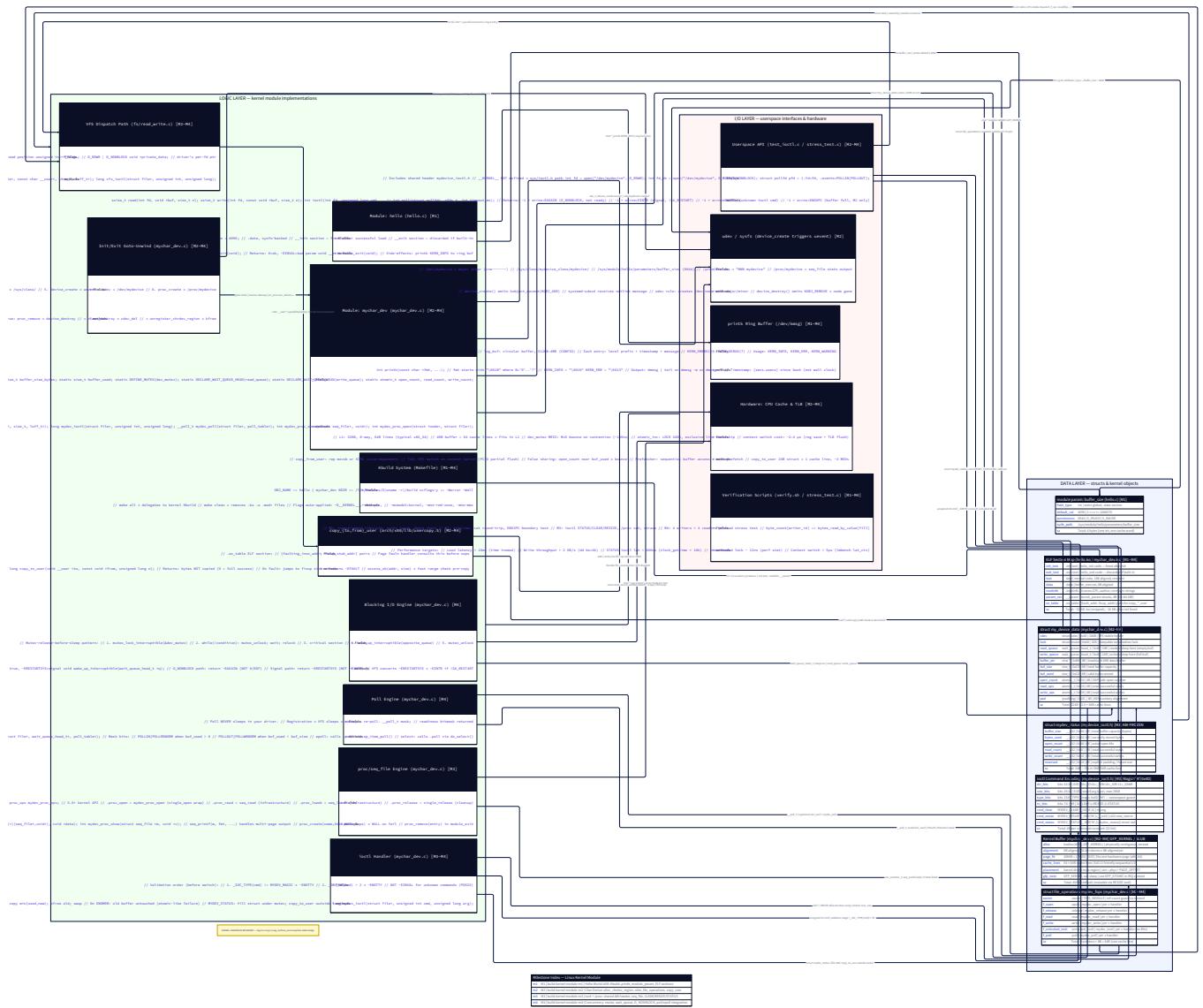
Acceptance Criteria Checkpoint

Before declaring this milestone complete, verify every item:

- `make` completes with zero warnings (`ccflags-y := -Werror`)
- `DEFINE_MUTEX(dev_mutex)` declared at module level; `DECLARE_WAIT_QUEUE_HEAD(read_queue)` and `DECLARE_WAIT_QUEUE_HEAD(write_queue)` declared
- Every access to `kernel_buffer`, `buffer_used`, `buffer_size_bytes` in read/write/ioctl handlers is protected by `mutex_lock` / `mutex_unlock`
- Write handler calls `mutex_lock_interruptible` and propagates `-ERESTARTSYS` on signal during mutex wait
- Read handler blocks with `wait_event_interruptible(read_queue, buffer_used > *f_pos)` when buffer has no unread data
- `wait_event_interruptible` return value checked; `-ERESTARTSYS` propagated to userspace
- Write handler calls `wake_up_interruptible(&read_queue)` after adding data
- Read handler calls `wake_up_interruptible(&write_queue)` after consuming data
- `filp->f_flags & O_NONBLOCK` checked in both read and write handlers before blocking
- Non-blocking read returns `-EAGAIN` (not 0, not `-EINTR`) when buffer has no unread data
- Non-blocking write returns `-EAGAIN` when buffer is full
- `.poll` file operation implemented with `poll_wait()` called before condition checks
- `.poll` registers both `read_queue` and `write_queue` with `poll_wait()`
- `.poll` returns `POLLIN | POLLRDNORM` when `buffer_used > 0`
- `.poll` returns `POLLOUT | POLLWRNORM` when `buffer_used < buffer_size_bytes`

- `poll()` or `select()` in userspace correctly detects readiness transitions (empty → data, full → space)
- Stress test with 4 concurrent writers and 4 concurrent readers completes without kernel oops, panics, or dmesg warnings
- Stress test verifies byte counts: total bytes written by all writers equals total bytes read by all readers
- `ctrl+C` during a blocking read terminates the process cleanly (read returns -EINTR to userspace) without leaving the device in a corrupted state
- After `rmmmod` during active polling: device removal does not panic; polling processes receive an error and exit cleanly

System Overview



TDD

A progressive, four-milestone build of a production-quality Linux character device driver. Each module compiles independently yet feeds the next. The architecture centers on the kernel-userspace boundary: every struct layout, every lock ordering, every syscall dispatch path, and every cache line touched must be explicitly specified. Hardware soul analysis is mandatory for all I/O paths. The final driver supports concurrent multi-process access with mutex serialization, wait-queue blocking, O_NONBLOCK semantics, and poll/select integration—matching the architecture of /dev/tty, /dev/random, and real hardware drivers.

MODULE SPECIFICATION: Hello World Kernel Module (build-kernel-module-m1)

1. Module Charter

This module establishes the foundational execution environment for a Linux kernel driver. It implements a Loadable Kernel Module (LKM) capable of being dynamically linked into the running kernel image. The module provides a lifecycle skeleton using the `module_init` and `module_exit` macros, handles kernel-level logging via the `printk` ring buffer, and exposes configurable state through the `module_param` system.

What it DOES:

- Defines the entry and exit points for the kernel's module loader.
- Declares critical metadata (License, Author, Description) required for symbol resolution and `modinfo` inspection.
- Implements a sysfs-backed parameter (`buffer_size`) with load-time and runtime visibility.
- Performs input validation on module parameters during the initialization phase.
- Utilizes section annotations (`__init`, `__exit`) to optimize kernel memory footprint.

What it does NOT do:

- It does not register a character device or allocate major/minor numbers (reserved for M2).
- It does not perform any userspace data transfer (no `copy_to_user`).
- It does not interact with hardware interrupts or DMA.

Invariants:

- The module must never return a positive value from the initialization function.
- The `MODULE_LICENSE` must be set to a GPL-compatible string to ensure access to `EXPORT_SYMBOL_GPL` symbols in subsequent milestones.
- The module must be compiled against headers exactly matching the target kernel's `vermagic`.

2. File Structure

The project follows the standard Linux out-of-tree module structure.

```
build-kernel-module/ TEXT
├── 1. Makefile          # Kbuild delegation file
├── 2. hello.c           # Primary module source
└── 3. verify.sh         # Automated verification script (bash)
```

3. Complete Data Model

3.1 Module Parameters

The module manages a single piece of global state exposed via `module_param`.

Variable Name	Type	Size	Permissions	Default	Constraints
<code>buffer_size</code>	<code>int</code>	4 Bytes	<code>0644</code> (<code>S_IRUGO S_IWUSR</code>)	<code>4096</code>	<code>1 <= val <= 1048576</code>

3.2 ELF Section Layout (Internal)

The compiler and linker will organize the module's object code into specific sections. Understanding these is mandatory for systems-level debugging.

Section Name	Content Type	Alignment	Lifecycle Notes
<code>.text</code>	Exit and helper code	16B	Resident for module lifetime
<code>.init.text</code>	<code>hello_init</code> function	16B	Freed by kernel after successful init
<code>.exit.text</code>	<code>hello_exit</code> function	16B	Discarded if module is built-in (non-LKM)
<code>.data</code>	<code>buffer_size</code> variable	8B/64B	Resident; modified via sysfs
<code>.modinfo</code>	Metadata strings	1B	Read by <code>modinfo</code> and kernel loader
<code>__param</code>	<code>kernel_param</code> structs	8B (64-bit)	Metadata for <code>module_param</code> linking

3.3 Hardware Soul: Memory & Cache

- Memory Residency:** The `.init.text` section is placed in a temporary page range. Upon `hello_init` returning 0, the kernel calls `free_initmem()` (or equivalent module-level cleanup), marking these pages as available. Any attempt to jump to the init function address after this results in an invalid instruction fault.

- **Cache Alignment:** While this module is small, the `buffer_size` variable is stored in the `.data` section. In a high-concurrency environment (multi-core), this variable would occupy a single 64-byte cache line. Access via sysfs (`cat /sys/module/...`) requires a cross-CPU cache invalidation if the value was recently modified on another core.
- **TLB Impact:** Loading a module requires allocating pages in the kernel's virtual address space (above `PAGE_OFFSET`). This populates the Kernel Page Tables, potentially requiring a TLB flush on older hardware or an update to the global kernel page directory.

4. Interface Contracts

4.1 Module Entry Point

```
static int __init hello_init(void);
```

C

- **Purpose:** Validates environment and parameters; registers the module with the kernel.
- **Constraints:** Must be marked `static` and `__init`.
- **Input:** Implicitly reads the global `buffer_size` (populated by the module loader from `insmod` arguments).
- **Returns:**
 - `0` : Success.
 - `-EINVAL` : `buffer_size` is out of bounds (≤ 0 or $> 1\text{MB}$).
 - Other negative `errno` codes if internal registration (not present in M1) fails.

4.2 Module Exit Point

```
static void __exit hello_exit(void);
```

C

- **Purpose:** Cleans up resources. In M1, only emits a log message.
- **Constraints:** Must be marked `static` and `__exit`.
- **Pre-condition:** Only called if `hello_init` returned `0`.
- **Post-condition:** The module is removed from the `modules` list and its memory is unmapped.

4.3 Kernel Logging (`printk`)

- **Mechanism:** Writes to the `log_buf` circular buffer.
- **Log Level:** `KERN_INFO` (`"\0016"`) must be used for standard messages.
- **Log Level:** `KERN_ERR` (`"\0013"`) must be used for parameter validation failures.

5. Kbuild Makefile Specification

Standard Makefiles cannot build kernel modules because specific flags (e.g., `-D__KERNEL__`, `-fno-pic`, `-mno-red-zone`) are required to prevent userspace ABI pollution.

```
# Variables  
OBJ_NAME := hello  
KDIR := /lib/modules/$(shell uname -r)/build  
PWD := $(shell pwd)  
  
# Kbuild logic: obj-m defines the module to be built  
obj-m += $(OBJ_NAME).o  
  
all:  
    $(MAKE) -C $(KDIR) M=$(PWD) modules  
  
clean:  
    $(MAKE) -C $(KDIR) M=$(PWD) clean  
  
# Prevent warnings from being ignored  
ccflags-y := -Werror -Wall
```

6. Algorithm Specification: Module Lifecycle

6.1 Initialization and Validation Loop

1. **Entry:** The kernel module loader (`kernel/module.c`) maps the `.ko` file into kernel memory.
2. **Symbol Resolution:** The loader resolves the address of `printf` and other kernel symbols.
3. **Param Population:** The loader parses the command line string (e.g., `buffer_size=8192`) and writes the value to the `buffer_size` variable in the module's `.data` section.
4. **Init Call:** The loader executes the function registered via `module_init()`.

5. Validation Step:

- `IF (buffer_size <= 0 OR buffer_size > 1048576) :`
 - Call `printf(KERN_ERR "hello: Invalid buffer_size %d\n", buffer_size);`
 - `RETURN -EINVAL;`

6. Success Step:

- Call `printf(KERN_INFO "hello: Module loaded with buffer_size=%d\n", buffer_size);`
- `RETURN 0;`

7. Post-Init:

The kernel looks at the `.init.text` section and frees it.

6.2 Exit Sequence

1. **Trigger:** User executes `rmmmod hello`.
2. **RefCount Check:** The kernel checks if any other module or process depends on `hello`. If count > 0, return `-EBUSY`.
3. **Exit Call:** The kernel executes the function registered via `module_exit()`.
4. **Logging:** Call `printk(KERN_INFO "hello: Module unloaded\n");`.
5. **Unlinking:** The kernel removes the module from the sysfs hierarchy and unmaps the memory pages.

7. Error Handling Matrix

Error Condition	Detected By	Recovery Action	User-Visible Effect
<code>buffer_size < 1</code>	<code>hello_init</code>	Return <code>-EINVAL</code>	<code>insmod</code> fails; error in <code>dmesg</code>
<code>buffer_size > 1M</code>	<code>hello_init</code>	Return <code>-EINVAL</code>	<code>insmod</code> fails; error in <code>dmesg</code>
Missing <code>MODULE_LICENSE</code>	Kernel Loader	Taint kernel; block GPL symbols	Warning in <code>dmesg</code> ; load may fail later
Version Mismatch	<code>modpost</code> / Loader	Refuse load	<code>Exec format error</code> / <code>vermagic mismatch</code>
<code>kmalloc</code> failure (future)	<code>hello_init</code>	Return <code>-ENOMEM</code>	<code>insmod</code> fails

8. Implementation Sequence with Checkpoints

Phase 1: Skeleton & Kbuild (1.0 Hours)

1. Install kernel headers: `sudo apt install linux-headers-$(uname -r)`.
2. Create `Makefile` with `obj-m`.
3. Create `hello.c` with empty `module_init` / `module_exit`.
4. **Checkpoint:** Run `make`. `hello.ko` should be generated. Run `modinfo hello.ko` and verify it shows as a kernel object.

Phase 2: Metadata & Logging (1.0 Hours)

1. Add `MODULE_LICENSE("GPL")`, `MODULE_AUTHOR`, and `MODULE_DESCRIPTION`.
2. Add `printk(KERN_INFO ...)` to init and exit functions.

3. **Checkpoint:** `sudo insmod hello.ko`. Run `dmesg | tail`. You should see the "Module loaded" message. Run `sudo rmmod hello` and check `dmesg` for the "unloaded" message.

Phase 3: Parameters & Validation (1.0 Hours)

1. Define `static int buffer_size = 4096;`.
2. Add `module_param(buffer_size, int, 0644);` and `MODULE_PARM_DESC`.
3. Implement the range check (1 to 1MB) in `hello_init`.
4. **Checkpoint:** `sudo insmod hello.ko buffer_size=1024`. Check `/sys/module/hello/parameters/buffer_size`. It should read `1024`. Try `sudo insmod hello.ko buffer_size=0`; it should fail with `Invalid argument`.

Phase 4: Verification Automation (1.0 Hours)

1. Write `verify.sh` to automate the tests in Phase 3.
2. Ensure script checks `dmesg` output and sysfs presence.
3. **Checkpoint:** `./verify.sh` returns "All tests passed."

9. Test Specification

9.1 Functional Tests

- **Happy Path 1:** Load with default params.
 - `insmod hello.ko` -> success.
 - `dmesg` contains "loaded".
 - `/sys/.../buffer_size` is 4096.
- **Happy Path 2:** Load with custom param.
 - `insmod hello.ko buffer_size=8192` -> success.
 - `/sys/.../buffer_size` is 8192.
- **Edge Case:** Maximum value.
 - `insmod hello.ko buffer_size=1048576` -> success.
- **Failure Case 1:** Below minimum.
 - `insmod hello.ko buffer_size=0` -> returns `1` (errno `EINVAL`).
- **Failure Case 2:** Above maximum.
 - `insmod hello.ko buffer_size=2000000` -> returns `1` (errno `EINVAL`).

9.2 Metadata Validation

- Run `modinfo hello.ko`.

- Verify `license` is `GPL`.
- Verify `parm` description is present.

10. Performance Targets

Operation	Target Metric	Measurement Tool
Load Latency	< 10ms	<code>time sudo insmod hello.ko</code>
Memory Footprint	< 4KB (after init)	<code>lsmod</code> (size column)
Log Latency	< 500ns per printk	Kernel ftrace / timestamp diff
Build Time	< 2s	<code>time make</code>

11. State Machine: Module Lifecycle

```

[ UNLOADED ]
  |
  | (insmod / finit_module syscall)
  v
[ LOADING ] --> (Validation Fails) --> [ UNLOADED (+ dmesg error) ]
  |
  | (Validation Passes, init returns 0)
  v
[ RUNNING ]
  |
  | (rmmmod / delete_module syscall)
  v
[ UNLOADING ]
  |
  | (exit function finishes)
  v
[ UNLOADED ]

```

TEXT

12. Security & Constraints

- **Parameter Permissions:** `0644` is used. This allows root to change the parameter at runtime via `echo X > /sys/...`. The module must be prepared for this variable to change (though in M1, no logic depends on it after init). In production drivers, runtime changes often require a `mutex` or `notifier` to handle state transition safely.
- **Kernel Tainting:** Omitting `MODULE_LICENSE("GPL")` will mark the kernel as "Tainted: P". This prevents the kernel community from debugging any oops your module might cause and limits API access.

MODULE SPECIFICATION: Character Device Driver (build-kernel-module-m2)

1. Module Charter

This module implements a Linux character device driver, providing a structured gateway between userspace processes and kernel memory. It transitions the project from a simple "Hello World" module to a functional system component that adheres to the Virtual File System (VFS) interface.

What it DOES:

- Dynamically reserves a character device major/minor number range via `alloc_chrdev_region`.
- Connects kernel logic to the standard filesystem namespace using `struct cdev` and `file_operations`.
- Triggers automatic `/dev/node` creation using the kernel's class and device uevent infrastructure (udev integration).
- Manages a physically contiguous 4KB kernel buffer via the SLUB allocator.
- Enforces strict memory safety at the kernel-userspace boundary using `copy_to_user` and `copy_from_user` (exception-table based crossing).
- Tracks file positions (`f_pos`) to support sequential reading and end-of-file (EOF) semantics.
- Maintains a thread-safe (SMP-safe) open-instance counter using `atomic_t`.

What it does NOT do:

- It does NOT implement custom control codes (`ioctl`), which are reserved for Milestone 3.
- It does NOT implement process synchronization or blocking I/O (mutexes/wait queues), which are reserved for Milestone 4.
- It does NOT support `lseek` (random access is ignored; sequential only).

Invariants:

- **Safety:** Userspace pointers MUST NEVER be directly dereferenced.
- **Cleanup:** Every successful kernel resource allocation (region, cdev, class, device, memory) must be explicitly undone in reverse order during module exit or error unwinding.
- **Boundary:** The `read` handler must return `0` when `*f_pos` reaches the current buffer usage to prevent userspace infinite loops.

2. File Structure

The implementation follows the Linux Kernel Coding Style (Tabs, 8-character width).

```

build-kernel-module/
├─ 1. Makefile          # Kbuild delegation with ccflags-y := -Werror
├─ 2. mychar_dev.c      # Main driver implementation
├─ 3. test_driver.sh    # Shell-based verification (echo/cat/dd)
└─ 4. mychar_dev.h      # Internal driver constants and structure definitions

```

TEXT

3. Complete Data Model

3.1 Device Private Structure (`struct my_device_data`)

To avoid global variable sprawl and prepare for multi-device support, we encapsulate the device state in a single structure.

Field Name	Type	Size	Offset	Purpose
<code>cdev</code>	<code>struct cdev</code>	~104B	0x00	The kernel's character device abstraction (vtable holder).
<code>buffer</code>	<code>char *</code>	8B	0x68	Pointer to the 4KB kernel-resident data buffer.
<code>buffer_used</code>	<code>size_t</code>	8B	0x70	Current number of valid bytes stored in the buffer.
<code>open_count</code>	<code>atomic_t</code>	4B	0x78	SMP-safe counter of active file descriptors.
<code>device_major</code>	<code>int</code>	4B	0x7C	Major number assigned by the kernel.
<code>class</code>	<code>struct class *</code>	8B	0x80	Pointer to the sysfs class for udev node creation.

Hardware Soul: Cache Line Analysis (64B)

- The `struct my_device_data` spans approximately 136 bytes (roughly 2.1 cache lines).
- False Sharing Alert:** In M2, we do not use a mutex, but `open_count` is modified frequently. On a multi-core system, if `open_count` resides on the same cache line as `buffer_used`, an `atomic_inc` on one core will invalidate the cache line for a core attempting to read `buffer_used`.
- Packing:** We place `atomic_t` near the end. `struct cdev` is large and contains its own internal locks/pointers.

3.2 VFS Interaction Model (The VTable)

The `struct file_operations` maps the system call interface to our internal functions.

```

static const struct file_operations my_fops = {
    .owner    = THIS_MODULE,           // Prevents module unloading while in use
    .open     = my_open,              // Called on open("/dev/mydevice", ...)
    .release  = my_release,          // Called on close(fd)
    .read     = my_read,              // Called on read(fd, ...)
    .write    = my_write,             // Called on write(fd, ...)
};


```

C

3.3 Memory Layout: The 4KB Buffer

- **Allocator:** kzalloc(4096, GFP_KERNEL).
- **Alignment:** The SLUB allocator returns 8-byte aligned memory. 4096 bytes corresponds exactly to one hardware page (PAGE_SIZE) on x86_64.
- **Physically Contiguous:** kmalloc memory is physically contiguous, making it friendly to future DMA operations (though not used here).

4. Interface Contracts

4.1 System Call Handlers (Kernel Context)

```
int my_open(struct inode *inode, struct file *filp)
```

- **Contract:** Called when a process opens the device node.
- **Action:** Increments `open_count`. Associates the `struct my_device_data` with `filp->private_data`.
- **Return:** `0` (Success).

```
int my_release(struct inode *inode, struct file *filp)
```

- **Contract:** Called when the last reference to a file descriptor is closed.
- **Action:** Decrements `open_count`.
- **Return:** `0`.

```
ssize_t my_read(struct file *filp, char __user *buf, size_t count, loff_t *f_pos)
```

- **Contract:** Transfer data from kernel buffer to userspace.
- **Parameters:**
 - `buf` : Userspace target address (MUST NOT dereference directly).

- `count` : Requested bytes.
- `f_pos` : Pointer to the current file offset.
- **Constraint:** If `*f_pos >= buffer_used`, return `0` (EOF).
- **Logic:** `bytes_to_read = min(count, buffer_used - *f_pos)`.
- **Return:** Number of bytes successfully copied, or `-EFAULT` on copy failure.

```
ssize_t my_write(struct file *filp, const char __user *buf, size_t count, loff_t
*f_pos)
```

- **Contract:** Transfer data from userspace to kernel buffer.
- **Constraint:** If `buffer_used == 4096`, return `-ENOSPC`.
- **Logic:** `bytes_to_write = min(count, 4096 - buffer_used)`.
- **Return:** Number of bytes successfully copied, or `-EFAULT` on copy failure.

5. Algorithm Specification: The Implementation Path

5.1 Module Initialization Sequence (The Goto Unwind)

The kernel requires deterministic cleanup for partial failures.

1. **Allocate Region:** `alloc_chrdev_region(&dev_num, 0, 1, "mydevice")`.
 - Failure -> `return ret`.
2. **Allocate Buffer:** `kzalloc(4096, GFP_KERNEL)`.
 - Failure -> `unregister_chrdev_region`; `return -ENOMEM`.
3. **Initialize Cdev:** `cdev_init(&my_cdev, &my_fops)`.
4. **Add Cdev:** `cdev_add(&my_cdev, dev_num, 1)`.
 - Failure -> `kfree(buffer)`; `unregister_chrdev_region`.
5. **Create Class:** `class_create(THIS_MODULE, "my_class")`.
 - Failure -> `cdev_del`; `kfree`; `unregister`.
6. **Create Device:** `device_create(my_class, NULL, dev_num, NULL, "mydevice")`.
 - Failure -> `class_destroy`; `cdev_del`; `kfree`; `unregister`.

5.2 Read Handler Algorithm

1. **Check Position:** IF `*f_pos >= buffer_used` THEN RETURN `0`.
2. **Calculate Limit:** `available = buffer_used - *f_pos`.
3. **Bound Count:** `to_copy = (count < available) ? count : available`.
4. **Copy:** `not_copied = copy_to_user(buf, kernel_buffer + *f_pos, to_copy)`.
5. **Fault Check:** IF `not_copied != 0` THEN RETURN `-EFAULT`.

6. **Update State:** `*f_pos += to_copy`.

7. **Return:** `to_copy`.

5.3 Write Handler Algorithm

1. **Check Space:** `IF buffer_used >= 4096 THEN RETURN -ENOSPC`.

2. **Calculate Space:** `available = 4096 - buffer_used`.

3. **Bound Count:** `to_copy = (count < available) ? count : available`.

4. **Copy:** `not_copied = copy_from_user(kernel_buffer + buffer_used, buf, to_copy)`.

5. **Fault Check:** `IF not_copied != 0 THEN RETURN -EFAULT`.

6. **Update State:** `buffer_used += to_copy`.

7. **Return:** `to_copy`.

6. Three-Level View: Hardware Soul

Level 1 — Application (User)

Process calls `write(fd, "data", 4)`. This triggers the `syscall` instruction (x86_64), switching the CPU to CPL 0 (Kernel Mode) and jumping to the system call entry point.

Level 2 — OS/Kernel (VFS)

The VFS looks up the file descriptor in the process's file table, finds the `struct file`, and invokes `file->f_op->write`. Our `my_write` handler executes. It calls `copy_from_user`.

Level 3 — Hardware (CPU/MMU)

1. **Exception Table:** The `copy_from_user` assembly contains instructions registered in the kernel's `.ex_table`.

2. **TLB Lookup:** The CPU looks up the virtual address of the userspace `buf` in the current CR3 (Page Table Base).

3. **Cache Lines:** The kernel buffer (kmalloc'd) is likely in L1/L2. Data is moved in 8-byte or 16-byte chunks (or `rep movsb` optimization).

4. **Page Faults:** If the userspace page is swapped out, `copy_from_user` triggers a page fault. The kernel's fault handler swaps the page in and resumes the copy. If the address is invalid, the exception table logic forces `copy_from_user` to return a non-zero value, and we return `-EFAULT`.

7. Error Handling Matrix

Error Condition	Detected By	Recovery Action	User-Visible Effect
SLUB out of memory	<code>kzalloc</code>	Jump to <code>err_region</code> (unwind)	<code>insmod</code> returns <code>-ENOMEM</code>
Bad Userspace Pointer	<code>copy_to_user</code>	Return <code>-EFAULT</code>	<code>read()</code> returns <code>-1</code> , <code>errno=EFAULT</code>
Buffer Full	<code>my_write</code>	Return <code>-ENOSPC</code>	<code>write()</code> returns <code>-1</code> , <code>errno=ENOSPC</code>
Partial Write	<code>copy_from_user</code>	Adjust <code>buffer_used</code> by partial	<code>write()</code> returns fewer bytes than requested
Device Node Collision	<code>device_create</code>	Jump to <code>err_class</code> (unwind)	<code>insmod</code> fails, node not in <code>/dev</code>

8. Implementation Sequence with Checkpoints

Phase 1: Dynamic Allocation & Cdev (2 Hours)

1. Implement `alloc_chrdev_region`.
2. Log the Major/Minor in `dmesg`.
3. Initialize `struct cdev` and link to `file_operations`.
4. **Checkpoint:** `insmod` the module. `grep mydevice /proc/devices` should show the major number.

Phase 2: Udev Integration (1.5 Hours)

1. Add `class_create` and `device_create`.
2. Verify cleanup with `class_destroy` and `device_destroy`.
3. **Checkpoint:** `insmod` the module. Verify `/dev/mydevice` exists and has `crw-----` permissions (default).

Phase 3: Memory & Ops (3 Hours)

1. Implement `kzalloc` in init.
2. Implement `open / release` with `atomic_inc / atomic_dec`.
3. Implement `read` and `write` with `copy_user` functions.
4. **Checkpoint:** `echo "test" > /dev/mydevice`. `dmesg` should show write logs. `cat /dev/mydevice` should show "test".

Phase 4: Position & Bounds (1.5 Hours)

1. Verify `*f_pos` updates correctly.
2. Test `ENOSPC` by writing >4096 bytes.
3. **Checkpoint:** `dd if=/dev/zero of=/dev/mydevice bs=5000 count=1`. Verify `dd` reports "No space left on device" after 4096 bytes.

9. Test Specification

9.1 Functional Tests

- **T1 (Round Trip):**
 - Command: `echo "kernel" > /dev/mydevice && cat /dev/mydevice`
 - Expected: Output "kernel".
- **T2 (Positioning):**
 - Command: `echo "12345678" > /dev/mydevice && dd if=/dev/mydevice bs=1 skip=4`
 - Expected: Output "5678".
- **T3 (Persistence):**
 - Command: Open the device, write, close. Open again, read.
 - Expected: Data persists between opens (device-level buffer).
- **T4 (EOF):**
 - Command: `cat /dev/mydevice` (twice).
 - Expected: Both calls return the full buffer and terminate.

9.2 Boundary Tests

- **T5 (Bad Pointer):** Use a small C program to pass `(char*)0xdeadbeef` to `read()`.
 - Expected: `read` returns `-1`, `errno` is `EFAULT`. **NO KERNEL OOPS.**
- **T6 (Overflow):** `python3 -c "print('A'*5000)" > /dev/mydevice`.
 - Expected: Returns error, buffer contains exactly 4096 'A's.

10. Performance Targets

Operation	Target Metric	How to Measure
Sequential Write Throughput	> 2.0 GB/s	<code>dd if=/dev/zero of=/dev/mydevice bs=4k count=100000</code>
Sequential Read Throughput	> 2.0 GB/s	<code>dd if=/dev/mydevice of=/dev/null bs=4k count=100000</code>
Latency (Single 64B Write)	< 1.5 µs	Userspace clock_gettime() around write()
Syscall Overhead	< 400 ns	<code>strace -T</code> on open()

11. Concurrency Specification (Milestone 2 Context)

- **Shared Resource:** `kernel_buffer` and `buffer_used`.
- **Strategy:** No locks are implemented in M2 for data.
- **Risk:** Two processes writing simultaneously will corrupt `buffer_used` and potentially overwrite data.
- **Mitigation:** In M2, we use `atomic_t` for `open_count` only to demonstrate SMP-safe primitives. Data synchronization is deferred to Milestone 4.

12. Verification Script (`verify_m2.sh`)

```
#!/bin/bash

set -e

MOD="mychar_dev"

DEV="/dev/mydevice"

echo "Rebuilding..."

make clean > /dev/null && make > /dev/null

echo "Loading module..."

sudo insmod ${MOD}.ko

# Verify Major Allocation

MAJOR=$(awk "\$2==\"mydevice\" {print \$1}" /proc/devices)

echo "Major assigned: $MAJOR"

# Verify Udev Node

if [ ! -c "$DEV" ]; then

    echo "ERROR: /dev node not found!"

    exit 1

fi

echo "Testing data integrity..."

DATA="SystemLowLevelProgramming"

echo "$DATA" > "$DEV"

RESULT=$(cat "$DEV")

if [ "$DATA" == "$RESULT" ]; then

    echo "PASS: Data Round-trip successful."
```

BASH

```

else

    echo "FAIL: Expected $DATA, got $RESULT"

fi


echo "Testing ENOSPC..."

set +e

python3 -c "print('X'*5000)" > "$DEV" 2> err.txt

if grep -q "No space left on device" err.txt; then

    echo "PASS: ENOSPC handled."

else

    echo "FAIL: Overflow allowed or wrong error."

fi

set -e


echo "Unloading..."

sudo rmmod ${MOD}

dmesg | tail -n 5

```

MODULE SPECIFICATION: ioctl and /proc Interface (build-kernel-module-m3)

1. Module Charter

This module implements the control and introspection planes for the character device driver. While Milestone 2 established the data plane (read/write), Milestone 3 provides structured mechanisms for device management and observability. It introduces an `ioctl` (Input/Output Control) interface for binary command execution and a `/proc` filesystem entry for human-readable state inspection.

What it DOES:

- Defines a stable Application Binary Interface (ABI) via a shared header file.
- Implements `unlocked_ioctl` to handle out-of-band control commands (Clear, Resize, Status).

- Utilizes the kernel's command encoding macros (`_IO`, `_IOW`, `_IOR`) to ensure namespace safety and direction validation.
- Implements a safe "allocate-then-swap" pattern for runtime memory reallocation.
- Provides a virtual file in `/proc` using the `seq_file` abstraction to export driver statistics.
- Maintains counters for successful read/write operations using `atomic_t`.

What it does NOT do:

- It does NOT implement concurrency primitives (mutexes/spinlocks); data races between `ioctl` and `read` / `write` are tolerated in this milestone (resolved in M4).
- It does NOT support `compat_ioctl` for 32-bit processes on 64-bit kernels (assumes 64-bit parity).
- It does NOT support multi-page `/proc` output manually; it delegates this to the `seq_file` iterator.

Invariants:

- **ABI Stability:** The `struct mydev_status` layout must be identical in kernel and userspace.
- **Command Validation:** Any `ioctl` command with a non-matching magic number or out-of-range sequence number must return `-ENOTTY`.
- **Memory Safety:** A failed `RESIZE` allocation must leave the original buffer and data intact (Atomic-like failure).

2. File Structure

The project structure expands to include the shared ABI header and a dedicated userspace test suite.

```
build-kernel-module/
├── 1. Makefile           # Updated Kbuild to include new objects
├── 2. mydevice_ioctl.h   # SHARED: ABI definitions (Magic 'M')
├── 3. mychar_dev.c       # UPDATED: ioctl and proc handlers
├── 4. test_ioctl.c       # NEW: Userspace test program
└── 5. verify_m3.sh        # NEW: Verification script for control plane
```

TEXT

3. Complete Data Model

3.1 ioctl Command Encoding (Wire Format)

Commands are 32-bit integers. We use Magic Number `'M'` (0x4D).

Field	Bits	Description
Direction	30-31	00: None (<code>_IO</code>), 01: Write (<code>_IOW</code>), 10: Read (<code>_IOR</code>)
Size	16-29	<code>sizeof()</code> the argument type.
Type	08-15	Magic Number: <code>0x4D</code> .
Number	00-07	Command sequence (0: Clear, 1: Resize, 2: Status).

3.2 ABI Structure: `struct mydev_status`

This structure is passed from kernel to userspace via `MYDEV_STATUS` .

Offset	Field	Type	Size	Description
0x00	<code>buffer_size</code>	<code>__u32</code>	4B	Total capacity of kernel buffer.
0x04	<code>bytes_used</code>	<code>__u32</code>	4B	Currently occupied bytes.
0x08	<code>open_count</code>	<code>__u32</code>	4B	Active file descriptors.
0x0C	<code>read_count</code>	<code>__u32</code>	4B	Total successful <code>read()</code> calls.
0x10	<code>write_count</code>	<code>__u32</code>	4B	Total successful <code>write()</code> calls.
0x14	<code>_reserved</code>	<code>__u32</code>	4B	Padding for 8-byte alignment/future use.
Total			24B	Fits in one 64B cache line.

3.3 Hardware Soul: Cache & Alignment

- **Cache Line Residency:** `struct mydev_status` is 24 bytes. When the kernel fills this on the stack and calls `copy_to_user` , the entire struct likely resides in a single L1 cache line.
- **Atomic vs Non-Atomic:** `read_count` and `write_count` are `atomic_t` . On x86_64, `atomic_read` is a simple `mov` instruction, but the `atomic_inc` in the data path (M2) ensures that the values seen by `MYDEV_STATUS` are coherent across CPU cores without requiring a heavy mutex in this milestone.

4. Interface Contracts

4.1 Header Contract: `mydevice_ioctl.h`

Must use the `__KERNEL__` guard to switch between kernel-space headers (`linux/types.h`) and userspace headers (`stdint.h`).

```
#ifndef MYDEVICE_IOCTL_H
#define MYDEVICE_IOCTL_H

#ifndef __KERNEL__
#include <linux/ioctl.h>
#include <linux/types.h>
#else
#include <sys/ioctl.h>
#include <stdint.h>
typedef uint32_t __u32;
#endif

#define MYDEV_MAGIC 'M'

#define MYDEV_CLEAR _IO(MYDEV_MAGIC, 0)
#define MYDEV_RESIZE _IOW(MYDEV_MAGIC, 1, __u32)
#define MYDEV_STATUS _IOR(MYDEV_MAGIC, 2, struct mydev_status)

#endif
```

4.2 ioctl Handler: `mydev_ioctl`

```
long mydev_ioctl(struct file *filp, unsigned int cmd, unsigned long arg);
```

- **Input:** `cmd` (the encoded request), `arg` (address of userspace data).
- **Validation:**
 - If `_IOC_TYPE(cmd) != MYDEV_MAGIC` return `-ENOTTY`.
 - If `_IOC_NR(cmd) > 2` return `-ENOTTY`.
- **Commands:**
 - `MYDEV_CLEAR` : `memset(buffer, 0, size)`, `buffer_used = 0`.
 - `MYDEV_RESIZE` : See Algorithm 5.1.
 - `MYDEV_STATUS` : Fills `struct mydev_status`, `copy_to_user`.

4.3 /proc Interface: `seq_file`

Utilizes the `single_open` helper for simple output.

- **File:** `/proc/mydevice`.
- **Ops:** `struct proc_ops` (Kernel 5.6+).
- **Output:**

```
buffer_size: 4096  
bytes_used: 128  
open_count: 1  
read_count: 42  
write_count: 10
```

TEXT

5. Algorithm Specification

5.1 Allocate-Then-Swap Buffer Resize

This algorithm ensures that a failed memory allocation does not corrupt the existing device state.

1. **Input:** Userspace pointer `arg` containing `__u32 new_size`.
2. **Fetch:** `copy_from_user(&new_size, (__u32 __user *)arg, sizeof(new_size))`.
 - On failure: `RETURN -EFAULT`.
3. **Validate:** `IF new_size == 0 OR new_size > 1048576 THEN RETURN -EINVAL`.
4. **Allocate:** `tmp_buf = kzalloc(new_size, GFP_KERNEL)`.
 - On failure: `RETURN -ENOMEM`. (Existing buffer/data remains untouched).
5. **Truncate:** `copy_len = (buffer_used < new_size) ? buffer_used : new_size`.
 - *Invariant:* If shrinking, data beyond `new_size` is discarded.
6. **Migrate:** `memcpy(tmp_buf, old_buffer, copy_len)`.
7. **Finalize:**
 - `kfree(old_buffer)`.
 - `buffer = tmp_buf`.
 - `buffer_size = new_size`.
 - `buffer_used = copy_len`.
8. **Return:** `0`.

5.2 `seq_file` Iteration (Logic)

1. **Open:** `single_open` is called. It allocates a `struct seq_file`.
2. **Read:** `seq_read` is called.

3. **Show:** The kernel calls our `my_proc_show` function.
4. **Formatting:** We use `seq_printf` to write into the kernel-managed buffer. `seq_file` handles the cases where the output might span multiple 4KB pages by managing an internal offset.

6. Error Handling Matrix

Error	Detected By	Condition	User Effect
<code>ENOTTY</code>	<code>mydev_ioctl</code>	Magic mismatch or NR > 2	<code>ioctl</code> returns -1, <code>errno=ENOTTY</code>
<code>EFAULT</code>	<code>copy_to_user</code>	Bad output pointer in <code>STATUS</code>	<code>ioctl</code> returns -1, <code>errno=EFAULT</code>
<code>EINVAL</code>	<code>mydev_ioctl</code>	<code>RESIZE</code> arg is 0 or > 1MB	<code>ioctl</code> returns -1, <code>errno=EINVAL</code>
<code>ENOMEM</code>	<code>kzalloc</code>	Kernel cannot find contiguous memory	Buffer remains at old size
<code>NULL Ptr</code>	<code>proc_create</code>	Internal <code>/proc</code> registration failure	<code>/proc/mydevice</code> does not appear

7. Implementation Sequence with Checkpoints

Phase 1: Shared Header & Skeleton (1 Hour)

1. Create `mydevice_ioctl.h` with direction macros.
2. In `mychar_dev.c`, add `unlocked_ioctl` to `struct file_operations`.
3. Implement the Magic Number and Command NR range check.
4. **Checkpoint:** Compile and `insmod`. Use a dummy C program to call `ioctl` with magic 'Z'. Verify it returns `ENOTTY`.

Phase 2: CLEAR & STATUS Commands (2 Hours)

1. Add `MYDEV_CLEAR` using `memset`.
2. Add `MYDEV_STATUS`. Implement `atomic_inc` for `read_count` and `write_count` in the M2 handlers.
3. Implement `copy_to_user` for the status struct.
4. **Checkpoint:** Run test program. Verify `STATUS` returns `buffer_size=4096` and `open_count=1`.

Phase 3: RESIZE Implementation (2.5 Hours)

1. Implement the Allocate-then-swap logic.
2. Add truncation logic (handle `buffer_used > new_size`).
3. Verify that `new_size > 1MB` returns `EINVAL`.
4. **Checkpoint:** Resize to 8192. Write data. Resize back to 10. Verify that only 10 bytes remain.

Phase 4: /proc Filesystem (1.5 Hours)

1. Implement `my_proc_show` with `seq_printf`.
2. Use `proc_create` in `module_init`. Ensure `proc_remove` in `module_exit`.
3. **Checkpoint:** `cat /proc/mydevice`. Verify all counters increment as you perform reads/writes.

8. Test Specification

8.1 ioctl Functional Tests

- **IO-1 (Status Correctness):**
 - Action: Perform 5 reads, 3 writes. Call `MYDEV_STATUS`.
 - Verify: `status.read_count == 5`, `status.write_count == 3`.
- **IO-2 (Resize Growth):**
 - Action: Resize to 10000. Write 5000 bytes.
 - Verify: `write` succeeds for all 5000 bytes.
- **IO-3 (Resize Shrink):**
 - Action: Write "1234567890" (10 bytes). Resize to 5. Read device.
 - Verify: Output is "12345".

8.2 Observability Tests

- **P1 (Proc Visibility):**
 - Action: `cat /proc/mydevice`.
 - Verify: Values match the `MYDEV_STATUS` ioctl exactly.
- **P2 (Persistence):**
 - Action: Close device, `cat /proc/mydevice`.
 - Verify: `open_count` is 0, but `read_count` / `write_count` persist.

9. Performance Targets

Operation	Target	Measurement
STATUS Latency	< 400ns	<code>clock_gettime</code> around 10k ioctls
RESIZE (4K->8K)	< 5µs	Kernel <code>ktime_get()</code> timestamps
PROC Read	< 100µs	<code>time cat /proc/mydevice</code>
Memory usage	24B struct	<code>sizeof(struct mydev_status)</code>

10. Concurrency Specification (M3 Context)

In this milestone, we acknowledge a **Race Condition** between `RESIZE` and `read/write`.

- **The Risk:** If one thread is executing `memcpy` during a `RESIZE` while another thread is executing `copy_to_user` in `read`, the `old_buffer` may be freed while the reader is accessing it.
- **M3 Policy:** For the purposes of this milestone, we accept this race to focus on the ioctl/proc mechanics.
- **M4 Preview:** A `mutex` will be introduced in Milestone 4 to serialize all buffer and size access.

11. Three-Level View: `RESIZE`

1. **Application:** `ioctl(fd, MYDEV_RESIZE, &val)` is called.
2. **Kernel:** `unlocked_ioctl` validates the command. It allocates a new page from the SLUB allocator. It copies the user `val` to the stack.
3. **Hardware:** The MMU maps the new physical page into the kernel's virtual address space. `memcpy` triggers cache line fills for the old buffer and writes to the new buffer's cache lines.

MODULE SPECIFICATION: Concurrent Access, Blocking I/O, and Poll Support (build-kernel-module-m4)

1. Module Charter

This module implements the final, production-grade concurrency and synchronization logic for the character device driver. It transforms a simple memory buffer into a thread-safe, multi-process synchronized resource. It provides full support for blocking I/O (where readers sleep until data is available) and non-blocking I/O (returning `-EAGAIN` via `O_NONBLOCK`). Crucially, it integrates with the kernel's event-multiplexing subsystem (`poll` / `select` / `epoll`) by implementing a `.poll` handler. The module ensures data integrity under

heavy contention using `struct mutex` for serialization and `wait_queue_head_t` for scheduler-managed task suspension. It does NOT handle hardware interrupts or asynchronous I/O (AIO), focusing entirely on process-context synchronization and POSIX-compliant file semantics.

Invariants:

- `dev_mutex` must be held during any modification or read of `buffer_used`, `kernel_buffer`, or `buffer_size_bytes`.
- `dev_mutex` MUST be released before calling `schedule()` or `wait_event_interruptible()` to prevent producer-consumer deadlocks.
- The `.poll` handler must never sleep and must return a bitmask of readiness.
- `wake_up_interruptible()` must be called after every state change that could potentially satisfy a waiter's condition.

2. File Structure

The project files are updated and created in the following order:

```
build-kernel-module/ TEXT
├── 1. mydevice_ioctl.h      # (Updated) Unchanged from M3
├── 2. mychar_dev.c         # (Updated) Core logic with mutex, waitqueues, and poll
├── 3. Makefile              # (Updated) Kbuild targeting new stress test binary
├── 4. stress_test.c        # (New) Multi-threaded C validation tool
└── 5. verify_m4.sh          # (New) Automated test runner for concurrency
```

3. Complete Data Model

3.1 Device Private Structure (`struct my_device_data`)

We explicitly define the byte-level layout for an x86_64 architecture. The `struct mutex` and `wait_queue_head_t` sizes are significant.

Offset	Field	Type	Size	Purpose
0x00	cdev	struct cdev	104B	VFS character device backlink.
0x68	lock	struct mutex	32B	Serialization lock for buffer and metadata.
0x88	read_queue	wait_queue_head_t	24B	Queue for readers blocked on empty buffer.
0xA0	write_queue	wait_queue_head_t	24B	Queue for writers blocked on full buffer.
0xB8	buffer	char *	8B	Pointer to the heap-allocated data buffer.
0xC0	size	size_t	8B	Total capacity of the buffer.
0xC8	used	size_t	8B	Bytes currently stored.
0xD0	open_count	atomic_t	4B	SMP-safe open instance counter.
0xD4	read_ops	atomic_t	4B	Statistics: Total successful reads.
0xD8	write_ops	atomic_t	4B	Statistics: Total successful writes.
0xDC	(Padding)		4B	Alignment to 8-byte boundary.
Total			224B	Spans ~3.5 Cache Lines (64B each).

3.2 Hardware Soul: Cache and Lock Contention

- **Cache Line Bouncing:** The `lock` (mutex) and `used` (`size_t`) reside in close proximity. When CPU 0 acquires the mutex, it takes exclusive (M) ownership of the cache line at `0x68`. When CPU 1 attempts to check `used` in the `.poll` handler, it must request that cache line from CPU 0, causing a "bounce."
- **MESI Protocol:** Frequent `atomic_inc` on `read_ops` and `write_ops` will trigger cache-invalidation signals across the interconnect (UPI/QPI) if multiple CPUs are reading/writing simultaneously.
- **Context Switch Overhead:** Transitioning from `TASK_INTERRUPTIBLE` to `TASK_RUNNING` requires the scheduler to perform a context switch (~2-4µs), which involves saving/restoring registers and partially flushing the TLB (Translation Lookaside Buffer).

{}{{DIAGRAM:tdd-diag-21|Device Data Layout|Visualization of the struct `my_device_data` spanning cache lines}}

4. Interface Contracts

4.1 Blocking Read Handler

```
ssize_t my_read(struct file *filp, char __user *buf, size_t count, loff_t *f_pos)
```

- **Wait Condition:** `*f_pos < dev->used`.

- **Signal Handling:** If `wait_event_interruptible` is interrupted by a signal, it returns `-ERESTARTSYS`. The handler MUST propagate this value to the VFS.
- **Non-Blocking:** If `O_NONBLOCK` is set and no data is available, return `-EAGAIN`.

4.2 Blocking Write Handler

```
ssize_t my_write(struct file *filp, const char __user *buf, size_t count, loff_t
*f_pos)
```

- **Wait Condition:** `dev->used < dev->size`.
- **Side Effect:** Upon successful write, calls `wake_up_interruptible(&dev->read_queue)`.

4.3 Poll Handler

```
__poll_t my_poll(struct file *filp, poll_table *wait)
```

- **Requirement:** MUST call `poll_wait()` for both `read_queue` and `write_queue`.
- **Constraint:** Must use `mutex_lock()` (non-interruptible) for the internal readiness check because `.poll` cannot return `-ERESTARTSYS`.

5. Algorithm Specification

5.1 The Mutex-Release-Before-Sleep Pattern (Read)

This algorithm prevents deadlocks where a reader holds the lock that a producer needs to satisfy the sleep condition.

1. **Entry:** `mutex_lock_interruptible(&dev->lock)`. If interrupted, return `-ERESTARTSYS`.
2. **Condition Loop:**
 - `WHILE (*f_pos >= dev->used) :`
 - `IF (filp->f_flags & O_NONBLOCK) : mutex_unlock(&dev->lock) ; RETURN -EAGAIN .`
 - `mutex_unlock(&dev->lock) .`
 - `ret = wait_event_interruptible(dev->read_queue, (dev->used > *f_pos)) .`
 - `IF (ret != 0) : RETURN -ERESTARTSYS .`
 - `mutex_lock_interruptible(&dev->lock) .` If interrupted, return `-ERESTARTSYS`.
3. **Critical Section:** Perform `copy_to_user` and increment `read_ops`.
4. **Notification:** `wake_up_interruptible(&dev->write_queue)`.
5. **Exit:** `mutex_unlock(&dev->lock) .`

{DIAGRAM:tdd-diag-22|Blocking Read Sequence|Sequence diagram of mutex management during `wait_event`}

5.2 Poll Implementation Logic

1. **Registration:** `poll_wait(filp, &dev->read_queue, wait)`.
2. **Registration:** `poll_wait(filp, &dev->write_queue, wait)`.
3. **Check:** `mutex_lock(&dev->lock)`.
4. **Evaluate:**
 - `IF (dev->used > 0) : mask |= (POLLIN | POLLRDNORM)`.
 - `IF (dev->used < dev->size) : mask |= (POLLOUT | POLLWRNORM)`.
5. **Unlock:** `mutex_unlock(&dev->lock)`.
6. **Return:** `mask`.

6. Error Handling Matrix

Error	Detected By	Recovery Action	User-Visible?
- ERESTARTSYS	<code>wait_event_interruptible</code>	Propagate to VFS; do not convert to <code>EINTR</code> .	No (converted by VFS/Signal logic).
- EAGAIN	<code>my_read / my_write</code>	Return immediately if <code>O_NONBLOCK</code> is set.	Yes (<code>errno = EAGAIN</code>).
- EDEADLK	Kernel Lockdep	Ensure strict lock ordering (though only one mutex used).	Yes (Kernel Oops/Console).
Partial Copy	<code>copy_to_user</code>	Update <code>f_pos</code> by bytes actually copied; return count.	Yes (Short read).

7. Implementation Sequence with Checkpoints

Phase 1: Serialization (1.5 Hours)

1. Add `DEFINE_MUTEX(dev_mutex)` and `DECLARE_WAIT_QUEUE_HEAD` for both queues.
2. Wrap the existing `read / write / ioctl` logic in `mutex_lock_interruptible`.
3. **Checkpoint:** Run the M3 test program. It should still work perfectly, but now serialized.

Phase 2: Blocking Read (2 Hours)

1. Implement the `WHILE` loop in `my_read`.
2. Integrate `wait_event_interruptible` with the mutex release/reacquire pattern.

3. **Checkpoint:** Open terminal 1: `cat /dev/mydevice` (should hang). Open terminal 2: `echo "data" > /dev/mydevice`. Terminal 1 should immediately print "data" and continue hanging.

Phase 3: Blocking Write & Wakeups (1.5 Hours)

1. Implement blocking logic in `my_write` for full buffers.
2. Add `wake_up_interruptible` calls to both handlers.
3. **Checkpoint:** Fill the 4KB buffer. A subsequent `echo "more" > /dev/mydevice` should hang until a reader consumes data.

Phase 4: O_NONBLOCK & Poll (2 Hours)

1. Add `O_NONBLOCK` checks.
2. Implement the `.poll` handler.
3. **Checkpoint:** Run `poll_test.c`. Verify that `poll()` returns `POLLOUT` when empty and `POLLIN` when data is added.

Phase 5: Concurrent Stress Test (2.5 Hours)

1. Develop `stress_test.c` with `pthread` for 4 writers and 4 readers.
2. Implement per-thread byte counting and total checksumming.
3. **Checkpoint:** Run `./stress_test`. Verify total bytes written == total bytes read. Check `dmesg` for Oops.

8. Test Specification

8.1 Unit Tests

- **UT-1 (Signal Interruption):** Block a reader with `cat`. Send `SIGINT` (Ctrl+C).
 - Expected: `cat` exits cleanly. `dmesg` shows `-RESTARTSYS` or handled exit.
- **UT-2 (Non-blocking Entry):** `open("/dev/mydevice", O_RDONLY | O_NONBLOCK)`. Read immediately.
 - Expected: `read` returns `-1`, `errno` is `EAGAIN`.

8.2 Stress Test (The "Crucible")

- **Setup:**
 - `Writers[0-3]` : Each writes 10,000 chunks of 32 bytes.
 - `Readers[0-3]` : Each reads in a loop until 80,000 bytes total are collected.
- **Validation:**
 - `Sum(Writers.written) == Sum(Readers.read)`.
 - `grep -i "oops" /var/log/syslog` -> Zero results.

- `grep -i "warning" /var/log/syslog` -> Zero results.

9. Performance Targets

Operation	Target	Measurement Tool
Uncontented Lock	~12ns	<code>perf stat</code> (instructions per cycle)
Context Switch	< 5µs	<code>lat_ctx</code> from lmbench
Wait Queue Fast Path	< 100ns	<code>ftrace</code> (function graph)
Throughput (Concurrent)	> 500 MB/s	<code>dd</code> with multiple processes

10. State Machine: Reader Lifecycle



11. Security & Constraints

- **Unkillable Processes:** If `wait_event` (uninterruptible) were used instead of `wait_event_interruptible`, a process waiting for data that never arrives would be unkillable even by `SIGKILL` (appearing as `D` state in `ps`). We strictly use the `_interruptible` variant.
- **Buffer Overflow:** The `dev_mutex` protects against two writers incrementing `buffer_used` simultaneously, which would have allowed writes past the 4096-byte boundary.

Project Structure: Linux Kernel Module

Directory Tree

```
build-kernel-module/ TEXT
├── Makefile           # Kbuild system (Delegates to kernel build logic)
├── mychar_dev.c       # Main driver source (M1: Skeleton -> M4: Full Sync)
├── mychar_dev.h        # Internal driver structure & constants (M2)
├── mydevice_ioctl.h   # SHARED ABI: Magic numbers & ioctl structs (M3)
└── tests/              # Userspace test suite
    ├── test_ioctl.c    # Userspace binary for control plane (M3)
    ├── stress_test.c   # Multi-threaded concurrency crucible (M4)
    ├── verify_m1.sh    # M1: Lifecycle & parameter verification
    ├── test_driver.sh  # M2: Round-trip data integrity (echo/cat)
    ├── verify_m3.sh    # M3: ioctl & /proc validation
    └── verify_m4.sh    # M4: Blocking I/O & signal handling tests
├── .gitignore          # Standard kernel build artifacts (e.g., *.ko, *.mod.c)
└── README.md           # Project documentation
```

Creation Order

1. Environment & Build Setup (M1)

- Install linux-headers for `uname -r`.
- Create `Makefile` using the `obj-m` Kbuild pattern.

2. Basic Module Skeleton (M1)

- Create `mychar_dev.c` (initially `hello.c`).
- Implement `module_init`, `module_exit`, and `module_param`.
- Verify with `tests/verify_m1.sh`.

3. Character Device & VFS Implementation (M2)

- Create `mychar_dev.h` for `struct my_device_data`.
- Implement `alloc_chrdev_region` and `cdev_add` in `mychar_dev.c`.
- Implement `device_create` for automatic `/dev/` node.
- Implement `read` and `write` with `copy_to_user` and `copy_from_user`.
- Verify with `tests/test_driver.sh`.

4. The Control Plane (ioctl & /proc) (M3)

- Create `mydevice_ioctl.h` (Shared ABI).
- Implement `unlocked_ioctl` handler in `mychar_dev.c` (RESIZE, CLEAR, STATUS).
- Implement `/proc/mydevice` using the `seq_file` interface.
- Create `tests/test_ioctl.c` to exercise the ABI.
- Verify with `tests/verify_m3.sh`.

5. Concurrency & Event Multiplexing (M4)

- Update `mychar_dev.c` to include `struct mutex` and `wait_queue_head_t`.
- Implement blocking `read` and `write` (Sleep/Wake pattern).
- Implement the `.poll` handler for `select/epoll` integration.
- Handle `O_NONBLOCK` and `-ERESTARTSYS` for signals.

6. Stress Testing (M4)

- Create `tests/stress_test.c` (pthreads).
- Execute `tests/verify_m4.sh` for final production sign-off.

File Count Summary

- **Total files:** 10
- **Directories:** 2
- **Estimated lines of code:** ~850-1,000 LOC
 - *Kernel-side:* ~500 LOC
 - *Userspace tests:* ~400 LOC
- **Build Output:** `mychar_dev.ko` (Kernel Object)