

Spacedrive: A Unified Virtual Distributed File System for the Modern Era

Architecture of a Local-First, Semantically-Aware Dataspace

James Mathew Pine

james@spacedrive.com

Spacedrive Technology Inc.

Vancouver, British Columbia, Canada

Abstract

Personal data is scattered across laptops, phones, and cloud drives, making it impossible to manage cohesively. Spacedrive introduces a Virtual Distributed File System (VDFS) that unifies all of a user's storage into a single, virtual library. It allows users to browse and organize their entire digital life from one interface, while all files remain in their original locations.

The system is local-first, keeping data on-device and working entirely offline. Its architecture pragmatically solves traditionally hard problems: content-addressed storage enables automatic deduplication across all locations, a domain-separated sync model ensures data consistency without complex consensus, and a durable action system makes complex file operations safe and predictable. This unified, indexed view also serves as a perfect foundation for AI. Spacedrive leverages it to provide powerful, privacy-preserving semantic search and natural language commands, powered by local models.

This paper details the architecture and implementation of Spacedrive V2. We present five key technical contributions: (1) A Virtual Distributed File System with universal addressing, (2) A Durable and Pre-visualized Action System for safe, intent-driven operations, (3) Pragmatic Synchronization via domain separation, (4) Content-Addressed Deduplication on consumer hardware, and (5) an AI-Native Architecture for intelligent, privacy-first file management.

CCS Concepts

• **Information systems** → **Hierarchical storage management; Query representation;** • **Software and its engineering** → **Software architectures.**

Keywords

Virtual Distributed File System, VDFS, AI-Native Architecture, Natural Language File Management, Semantic Search, Data Synchronization, Tiered Storage, Local-First AI, Rust

ACM Reference Format:

James Mathew Pine. 2025. Spacedrive: A Unified Virtual Distributed File System for the Modern Era: Architecture of a Local-First, Semantically-Aware Dataspace. In *Proceedings of Spacedrive Whitepaper (Spacedrive '25)*. ACM, New York, NY, USA, 23 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 Introduction

The proliferation of personal computing devices and cloud services has created what we term "data fragmentation hell"—a state where users' digital assets are scattered across incompatible ecosystems, each with proprietary APIs, limited interoperability, and platform lock-in. A typical user today manages files across local devices (laptop, phone, tablet), cloud services (Google Drive, iCloud, Dropbox), and network storage, with no unified view or consistent management capabilities.

Existing file management solutions treat data as hierarchical folder structures, blind to content relationships and cross-device dependencies. This leads to fundamental problems: duplicate files consuming storage across devices, inability to find content regardless of location, loss of context when files move between devices, and fragmented metadata that doesn't follow the content.

We present Spacedrive, a Virtual Distributed File System (VDFS) that reimagines personal data management as a unified, content-aware ecosystem. Unlike traditional file managers that operate on individual devices, Spacedrive creates a Library—a portable, self-contained database that maintains a comprehensive index of all user content across devices and locations.

Spacedrive's architecture is built on four foundational innovations that solve traditionally hard problems in distributed systems:

- **SdPath Universal Addressing:** A unified path system that makes device boundaries transparent, enabling seamless file operations across any location (local drives, network storage, cloud services) with a single, consistent API.
- **Entry-Centric Data Model:** Every filesystem entity becomes a stateful Entry with immediate metadata capabilities, allowing instant organization and tagging without waiting for content analysis or indexing completion.
- **Content-Addressable Deduplication:** A versioned content addressing system using strategic sampling for large files, enabling bit-level deduplication across devices while maintaining performance on consumer hardware.
- **Domain-Separated Synchronization:** A pragmatic sync architecture that separates concerns into Index Sync (device-owned filesystem state), User Metadata Sync (content-universal

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Spacedrive '25, Vancouver, BC, Canada

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

tags and ratings), and File Operations (explicit transfer protocols), eliminating the complexity of traditional distributed consensus systems.

- **Durable, Pre-visualized Actions:** An intent-driven job system that allows users to simulate any file operation—from simple copies to complex, multi-device reorganizations—before committing. This "dry run" capability, powered by the Spacedrive index, provides a clear before-and-after preview, including projected space savings and deduplication. Once confirmed, actions are converted into durable, verifiable jobs that guarantee eventual completion, even across offline devices, making the file management experience fluid, intuitive, and exceptionally reliable.
- **AI-Native Architecture:** A VDFS designed from the ground up for AI integration. This enables not just semantic search, but also natural language file management, where complex commands are translated into durable, verifiable actions. The system can proactively suggest organizational improvements and automate storage tiering, all while preserving user privacy and control through support for local models via interfaces like Ollama and model-agnostic design principles.

This paper details Spacedrive's architecture through the lens of a production system—implemented in Rust with modern async patterns, tested across multiple platforms, and designed for real-world deployment. We demonstrate how careful domain separation and content-awareness enable features previously limited to enterprise storage systems: cross-device deduplication, semantic search, intelligent tiering, and conflict-free synchronization at consumer scale.

Recognizing that mobile devices are central to modern computing, Spacedrive incorporates sophisticated resource management from its core architecture. The system dynamically adapts to device constraints—intelligently throttling CPU usage on battery power, respecting mobile data limits, and working within platform-specific background processing restrictions. This mobile-first approach ensures Spacedrive remains responsive and efficient whether running on a powerful desktop or a resource-constrained smartphone, making unified file management accessible across all devices without compromising battery life or performance.

Spacedrive's design is predicated on a key insight: personal data management is fundamentally different from enterprise distributed systems. Users own their data completely, conflicts are primarily about organization rather than content, and the system can assume eventual consistency within a single user's ecosystem. This enables dramatically simpler architectures that still deliver powerful distributed capabilities.

2 Related Work

Spacedrive builds upon decades of research in distributed file systems, personal information management, and content-addressable storage. We position our work within this landscape to highlight our unique contributions.

2.1 Traditional Cloud Sync Services

Commercial cloud storage services (Dropbox [4], Google Drive, iCloud) provide basic file synchronization but suffer from platform

lock-in and lack content-addressing. These services treat files as opaque blobs, missing opportunities for deduplication and semantic understanding. Unlike Spacedrive, they require continuous internet connectivity and centralize user data on corporate servers.

2.2 Distributed File Systems

Research systems like IPFS [1] and production systems like Ceph [7] demonstrate the power of content-addressable storage. However, their complexity and resource requirements make them unsuitable for personal use. IPFS requires understanding of cryptographic hashes and peer-to-peer networking, while Ceph targets datacenter deployments. Spacedrive adopts content-addressing principles while hiding complexity behind familiar file management interfaces.

2.3 Virtual Distributed File Systems in the Datacenter

The concept of a Virtual Distributed File System (VDFS) has been explored in the context of large-scale data analytics. Alluxio (formerly Tachyon) [5] introduced a memory-centric VDFS designed to sit between computation frameworks like Apache Spark and various storage systems (e.g., HDFS, S3). Alluxio's primary goal is to accelerate data analytics jobs by providing a unified, high-throughput data access layer, effectively decoupling computation from storage in a datacenter environment.

While Spacedrive shares the VDFS terminology, its architectural goals and target domain are fundamentally different. Where Alluxio optimizes for performance in large, multi-tenant analytics clusters, Spacedrive is designed as a local-first, privacy-preserving dataspace for an individual's complete digital life. Spacedrive's innovations in universal addressing (SdPath), an entry-centric model with immediate metadata, and pragmatic synchronization are tailored to the challenges of personal data fragmentation across a heterogeneous collection of consumer devices, a problem space distinct from the performance and data-sharing challenges in large-scale analytics that Alluxio addresses.

2.4 Personal Knowledge Management

Tools like Obsidian and Logseq excel at managing structured knowledge through markdown files but lack general file management capabilities. They demonstrate the value of local-first architectures and portable data formats, principles that Spacedrive extends to all file types. Our work generalizes their approach from text-centric knowledge graphs to comprehensive file management.

2.5 Self-Hosted Solutions

Projects like Nextcloud provide self-hosted alternatives to commercial cloud services but retain client-server architectures that complicate deployment and maintenance. They require dedicated servers and technical expertise, limiting adoption. Spacedrive's peer-to-peer architecture eliminates server requirements while providing similar capabilities through direct device communication.

2.6 Semantic File Systems

Academic projects exploring semantic file organization date back to the Semantic File System [3] and Presto [2]. While these demonstrated the value of content-based organization, they predated modern AI capabilities. Spacedrive realizes this vision with contemporary machine learning, enabling natural language queries and intelligent automation impossible in earlier systems.

Our work synthesizes insights from these domains while addressing their individual limitations, creating a unified system that is simultaneously powerful, private, and accessible to non-technical users.

2.7 System Architecture Overview

Figure 1 presents the high-level architecture of Spacedrive, illustrating how the core components interact to provide a unified virtual distributed file system.

3 Learning from the Past: Architectural Evolution from Spacedrive v1

Spacedrive v2 represents a complete architectural reimplementation designed to fulfill the original vision on a more robust, scalable foundation. The initial version, first open-sourced in 2022, validated the core premise of a unified VDFS for personal data with significant community interest. However, as development progressed through early 2025, several foundational architectural challenges emerged that ultimately necessitated this rewrite.

3.1 Key Challenges in the Original Architecture

Post-mortem analysis of the v1 codebase revealed critical issues that prevented the system from achieving its goals:

- **The Dual File System Problem:** The most significant flaw was the existence of two parallel, incompatible file management systems—one for indexed Locations and another for ephemeral direct file access. This created a fractured user experience where fundamental operations like copying files between indexed and non-indexed folders were impossible, doubling the development burden for every file-related feature.
- **The `invalidate_query` Anti-Pattern:** The v1 architecture tightly coupled the Rust backend to the frontend's React Query caching keys through an `invalidate_query!` macro. This created a brittle system where backend changes could silently break the frontend, clearly indicating the need for proper event-driven architecture.
- **Over-Engineered Synchronization:** The original sync system attempted to solve mixed local and shared data with a custom CRDT implementation, leading to analysis paralysis where the complexity prevented the feature from ever shipping.
- **Excessive Job System Boilerplate:** While functional, the original job system required over 500 lines of boilerplate to define new background jobs, stifling rapid development and extensibility.
- **Abandoned Dependencies:** Critical dependencies like `prisma-client-rs` and `rsync` were created by the original team and later abandoned, leaving the project reliant on unmaintained forks.

3.2 Spacedrive v2 as an Architectural Solution

The v2 architecture presented in this paper directly addresses these challenges:

- The unified **SdPath** addressing system completely eliminates the dual file system problem. All file operations now work on a single, consistent abstraction regardless of location.
- A robust, decoupled **Event Bus** replaces the `invalidate_query!` anti-pattern, allowing components to subscribe to state changes without tight coupling.
- The **Pragmatic Synchronization** model with clear domain separation avoids over-engineering by applying tailored, simpler conflict resolution strategies to different data types.
- The new **Job System** with derive macros and automatic registration reduces boilerplate by over 90%, fostering extensibility.
- The technology stack has been modernized, replacing abandoned dependencies with actively maintained, community-trusted libraries like **SeaORM**.
- Network connectivity has been revolutionized by migrating from `libp2p` to **Iroh**, achieving 90%+ NAT traversal success rates and sub-2-second connection establishment, making device-to-device communication exceptionally reliable.
- A comprehensive **async GraphQL API** provides a modern, type-safe interface for frontend applications, while a powerful **CLI** enables scripting and automation of all Spacedrive operations.
- The codebase achieves **100% test coverage** across critical systems, ensuring reliability and enabling confident refactoring as the system evolves.

By learning from real-world challenges of the initial version, Spacedrive v2 delivers on the original promise with an architecture that is not only more powerful but also fundamentally simpler, more resilient, and built for the long term.

4 The Spacedrive VDFS Model

At the core of Spacedrive is a set of abstractions that model a user's data not as a collection of disparate file paths, but as a cohesive, unified Library with content-aware relationships.

4.1 The Library: A Portable Data Container

Rather than managing scattered databases and configurations, Spacedrive organizes everything into self-contained **Libraries**. Each Library is a `.sdlibrary` directory that functions as a complete, portable data container:

A Spacedrive Library is organized as a self-contained directory with four key components:

- **Configuration File:** Library settings and device registry
 - **Metadata Database:** Complete file index with relationships and tags
 - **Thumbnail Cache:** Organized storage for instant previews
 - **Concurrency Protection:** Safe multi-device access control
- This portable structure means that backing up your entire digital life is as simple as copying a single directory, and sharing a complete library with someone else requires no complex export process.

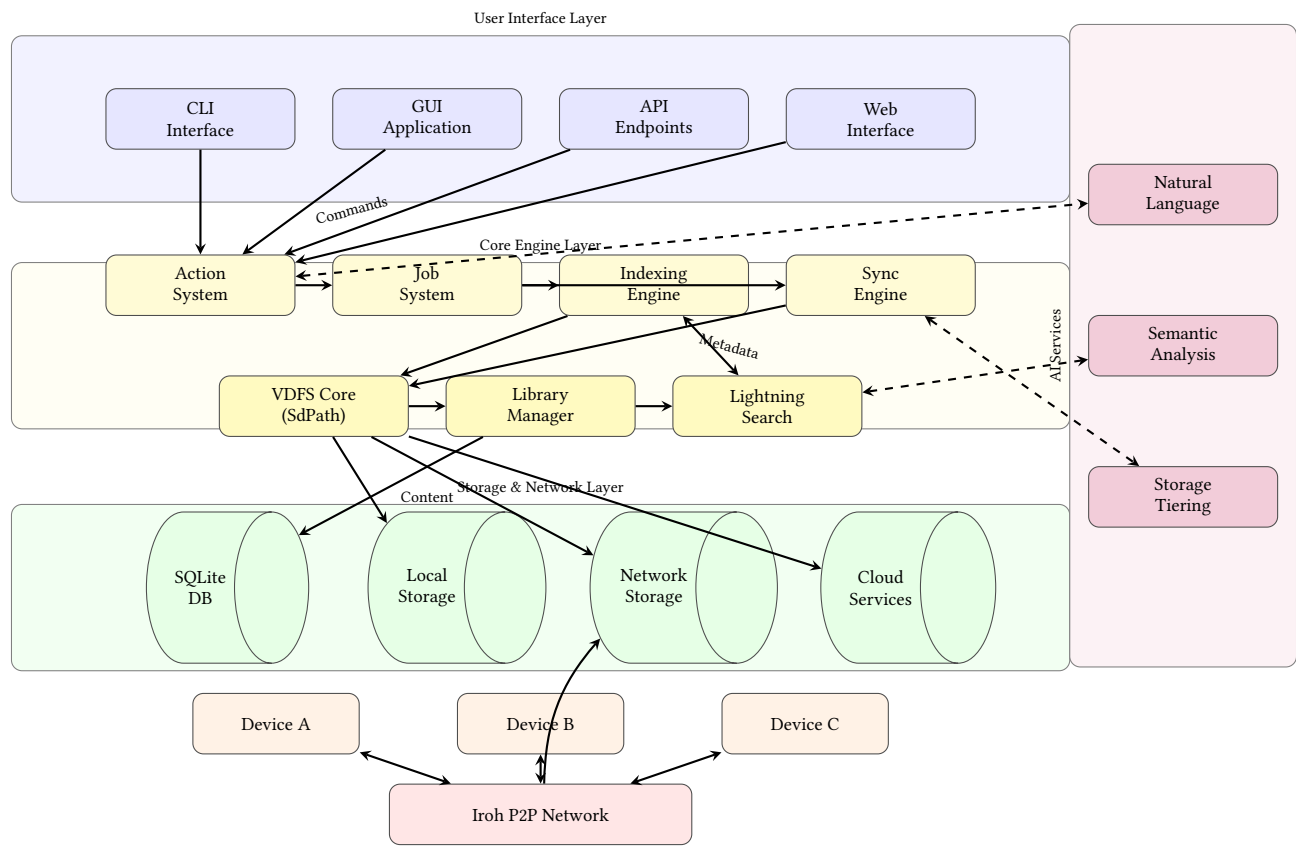


Figure 1: Spacedrive Architecture: The system consists of four main layers. The **User Interface Layer** provides multiple access methods (CLI, GUI, API, Web). The **Core Engine Layer** implements the key functionality through the Action System (command processing), Job System (background tasks), and distributed components. The **Storage & Network Layer** abstracts different storage types through the unified SdPath system. **AI Services** integrate throughout the stack to provide intelligent features while maintaining privacy through local-first processing.

This design provides several critical advantages: **backup** becomes copying a directory, **sharing** involves sending the complete Library, and **migration** across devices requires no complex export/import processes. Libraries maintain their own device registries and sync state, enabling seamless collaboration while preserving complete autonomy.

4.2 The Entry-Centric Data Model

The fundamental unit within a Library is the **Entry**—a universal representation that treats files and directories uniformly. Unlike traditional filesystems that separate metadata from content, every Entry in Spacedrive is designed for immediate metadata capability: Every file and directory in Spacedrive is represented as an Entry with the following key properties:

The key innovation is that **every Entry can receive metadata immediately**—users can tag, rate, and organize files the moment they encounter them, regardless of whether background indexing has completed.

Key Innovation: The `metadata_id` field ensures every Entry can be tagged, rated, and annotated immediately upon discovery,

Field	Description
Unique ID	Globally unique, immutable identifier
Universal Path	Complete location with device ID
Name & Type	File/directory name and type
Metadata ID	Instant tagging without content analysis
Content ID	Links to deduplication fingerprint
Discovery Time	First detection timestamp

Table 1: Entry data model fields

without waiting for slow content analysis or indexing completion. This "metadata-first" approach enables instant organization of any filesystem entity.

The Entry structure demonstrates this separation of concerns:

```
1 pub struct Entry {
2   pub id: Uuid,
3   pub path: SdPath,
4   pub name: String,
5   pub metadata_id: Uuid, // Immediate metadata capability
6   pub content_id: Option<ContentId>, // Populated asynchronously
7   pub discovered_at: DateTime<Utc>,
```

```

8 }
9
10 impl Entry {
11     pub fn tag(&self, tag: &Tag) -> Result<()> {
12         // Can tag immediately, no content_id required
13         self.metadata_id.add_tag(tag)
14     }
15 }

```

Listing 1: Simplified Entry structure showing metadata-first design

4.3 Richly-Structured Tag System

Spacedrive implements a sophisticated tagging system that provides powerful organization and search capabilities far beyond simple text-based tags. This system, inspired by flexible organizational tools, is designed to be as simple or as complex as the user desires.

4.3.1 Flexible Tag Naming. Unlike traditional tagging systems, Spacedrive tags have no restrictive naming limitations:

- **Non-Unique Names:** Multiple tags can share the same name, with meaning differentiated by relationships
- **No Character Limitations:** Tag names support any characters for natural, descriptive naming
- **Aliases and Shorthands:** Tags support alternative names for abbreviations, translations, or preferences

4.3.2 Hierarchical Tag Relationships. The core innovation is hierarchical relationships between tags, creating a vocabulary that mirrors real-world relationships. This implementation leverages closure tables [6] for efficient traversal of tag hierarchies:

Automatic Disambiguation: Multiple tags named "Freddy" can be automatically displayed as "Freddy (Five Nights at Freddy's)" or "Freddy (FNAF)" based on parent relationships and available shorthands, providing clarity without manual renaming.

Benefits of Tag Hierarchy:

- **Simplicity via Deduplication:** Files tagged with "Shrek" inherit tags like "Character", "Movie", and "Dreamworks" automatically
- **Intuition via Substitution:** Searching for "Dreamworks" returns files tagged with any child tags like "Shrek"
- **Rediscovery via Linking:** Broad searches reveal forgotten content through tag relationships

4.3.3 Tag Properties and Appearance. Tags support rich customization for enhanced usability:

- **Category Marking:** Tags can be marked as categories for UI grouping
- **Hidden Tags:** Tags can be hidden from default search results
- **Visual Identity:** Custom colors and icons for easy identification
- **Future Component Tags:** Planned "HAS" relationships for granular descriptions (e.g., "Shrek" HAS "Leather Vest")

This richly-structured tagging system transforms file organization from flat categorization into a powerful knowledge graph, enabling users to create organizational schemes that match their mental models while maintaining simplicity for basic use cases.

4.4 SdPath: Universal File Addressing

Central to the VDFS abstraction is **SdPath**—a universal addressing system that makes device boundaries transparent. While the primary form provides a direct physical coordinate (device identifier + local path), Spacedrive supports a more powerful content-aware addressing mode that transforms SdPath from a simple pointer into an intelligent content resolver.

```

1 #[derive(Clone, Debug)]
2 pub enum SdPath {
3     // Physical addressing: device + path
4     Physical {
5         device_id: DeviceId,
6         local_path: PathBuf
7     },
8     // Content-aware addressing: find optimal instance
9     Content {
10         cas_id: ContentId
11     },
12 }
13
14 // Same API works for all addressing modes
15 async fn copy_files(from: Vec<SdPath>, to: SdPath) ->
16     Result<()> {
17     for source in from {
18         // Resolve content-aware paths to optimal
19         // physical paths
20         let physical_source = match source {
21             SdPath::Physical { .. } => source,
22             SdPath::Content { cas_id } => {
23                 resolve_optimal_path(cas_id).await?
24             }
25         };
26
27         // Execute operation using resolved paths
28         p2p::transfer(&physical_source, &to).await?;
29     }
30     Ok(())
31 }

```

Listing 2: SdPath supports both physical and content-aware addressing

4.4.1 Content-Aware Addressing and Optimal Path Resolution. When an operation is initiated with a content-aware SdPath, Spacedrive performs an **optimal path resolution** query against the Library index:

- (1) **Content Lookup:** Query the ContentIdentity table to find all instances of the file content across all devices
- (2) **Candidate Evaluation:** Evaluate each instance based on a cost function considering:
 - **Locality:** Local device copies prioritized above all others
 - **Network Proximity:** Iroh provides real-time latency and bandwidth estimates
 - **Device Availability:** Filter for currently online devices
 - **Storage Tier:** Volume-Aware Storage Foundation prioritizes SSD over HDD
- (3) **Path Selection:** Select the lowest-cost valid path and proceed transparently

This mechanism makes file operations exceptionally resilient. If a user requests a file from an offline laptop, Spacedrive can transparently source the identical content from a NAS on the local network. This elevates SdPath from a simple address to an abstract, location-independent handle for content.

4.4.2 *Practical Applications.* This addressing system enables operations that were previously impossible or extremely complex:

- **Resilient Operations:** File operations succeed even when the original source is offline
- **Optimal Performance:** Automatically select the fastest available source
- **Simplified Development:** Applications reference content by ID without managing device availability
- **Transparent Failover:** Operations seamlessly switch to alternative sources

This abstraction transforms complex cross-device operations into simple, type-safe function calls, making the distributed nature of the filesystem completely transparent to both users and developers while providing unprecedented reliability and performance.

4.5 The Entry Lifecycle: Stateful Content Management

Unlike static file representations, Entries in Spacedrive transition through a formal lifecycle managed by an event-driven state machine:

- **Discovered:** Entry detected, basic metadata available
- **Processing:** Content ID, hash generation in progress
- **Available:** Fully indexed with rich metadata
- **Syncing:** Propagating changes across devices
- **Archived:** In cold storage, metadata retained

This lifecycle approach provides several advantages: **graceful handling** of long-running operations (large file hashing, media analysis), **resumable processing** after interruptions, **clear user feedback** about file status, and **deterministic state transitions** that eliminate race conditions common in distributed systems.

The state machine is implemented through Spacedrive's job system, where each lifecycle transition corresponds to specific job types (IndexerJob, ContentIdentificationJob, SyncJob) that can be paused, resumed, and monitored for progress. This ensures that even multi-gigabyte files or complex analysis operations integrate seamlessly into the user experience.

5 AI-Native VDFS: From Semantic Search to Intelligent Management

While many systems treat AI as an additive feature, Spacedrive is architected as an **AI-native dataspace**. The comprehensive, always-current index of the user's files serves as a perfect "world model" for an AI agent to reason about. This enables a shift from reactive file management (issuing manual commands) to a proactive, collaborative model where both the user and an AI agent can manage the dataspace, with the human always in the loop.

This is achieved through a flexible, privacy-first architecture that is model-agnostic, supporting both powerful cloud services and local models running on user hardware via interfaces like Ollama.

5.0.1 *A Day in Alice's Digital Life.* To illustrate how these capabilities transform everyday file management, consider Alice, a freelance designer juggling multiple projects across various devices and storage locations.

Monday morning: Alice sits down at her desk and asks Spacedrive: "Find my design assets from last fall that I never exported." Behind this simple request, a sophisticated dance begins:

The AI Observes: Spacedrive's indexing system has already performed deep analysis on Alice's files. It knows that her Sketch and Figma files from September through November contain layer names like "v3-final" and "client-approved," but lack corresponding PNG or PDF exports in nearby folders. The system has extracted color palettes, identified design patterns, and even transcribed text from her design review recordings.

The AI Orients: Cross-referencing the temporal query "last fall" with file metadata, the AI identifies 47 design files. It notices from Alice's audit log that she typically exports finals to a "Deliverables" folder, but these particular projects show no such exports. The AI also detects that several files have "URGENT" in their layer names—a pattern it has learned indicates deadline pressure that might have caused Alice to skip her usual export workflow.

The AI Decides: Rather than simply listing files, the AI formulates a helpful action plan. It generates a structured proposal: "I found 47 design files from last fall without exports. 12 appear to be final versions based on your naming patterns. Would you like me to batch export these as PNGs to your usual Deliverables folder?"

Alice reviews the proposed BatchExportAction, which shows exactly which files will be processed and where the exports will be saved. With one click, she approves, and the operation joins the durable job queue.

Later that week: The AI notices Alice has been manually moving screenshot files from her Desktop to project folders every few days. Having observed this pattern through the audit log, it proactively suggests: "I've noticed you regularly organize screenshots into project folders. I can automatically move new screenshots to the relevant project based on the window title captured in the metadata. Would you like me to set this up?"

This isn't just automation—it's intelligent assistance that learns and adapts while keeping Alice in complete control.

5.1 The Agentic Loop: Observe, Orient, Act

Spacedrive's AI capabilities are built on a classic agentic loop, where each stage is powered by a core component of the VDFS architecture:

Observe: The Indexing System is the sensory input. During the **deep indexing phase**, it goes beyond basic metadata to perform AI-powered analysis, extracting rich context like image content, video transcripts, and document summaries. This enriches the Spacedrive index, providing the AI with a deep understanding of the user's data.

Orient: With a complete "world model" in its index, the AI can orient itself. It analyzes file content, user-applied metadata (tags, ratings), and historical user actions (from the audit_log table) to understand context, identify patterns, and recognize organizational inconsistencies.

Decide & Act: The AI formulates a plan and proposes it as a structured Action. This is a critical safety and control mechanism; the AI does not execute arbitrary commands but is constrained to the same safe, verifiable primitives available to the user. A user

command like "Archive my old projects from last year that are over 1GB" is translated directly into a `FileCopyAction`.

5.2 Natural Language Management

The Action System serves as a stable, well-defined API that can be used to fine-tune language models. This allows Spacedrive to translate complex user requests from natural language into a series of verifiable actions.

As we saw with Alice's request to "find design assets from last fall that I never exported," the system seamlessly translates natural language into precise operations. Similarly, a command like "Move my last 3 screen recordings from the desktop to the 'Clips' folder on my NAS" is processed through semantic search to identify the relevant files, then translated into a structured `FileCopyAction` with appropriate source paths, destination, and move semantics.

This generated action is then fed into the **Durable Action System**, giving the user a complete preview of the operation for approval before it is committed to the durable job queue. This keeps the human perfectly in the loop, using AI as a powerful interpreter and planner rather than an opaque black box.

5.3 Proactive Assistance and Optimization

Beyond executing commands, the AI agent can proactively identify opportunities to help the user. By observing patterns, it can suggest helpful actions.

Organizational Suggestions: As demonstrated in Alice's workflow, when the AI observed her repeatedly moving screenshots from the Desktop to project folders, it proactively offered to automate this pattern. The architecture enables such capabilities—if the indexer identifies a screen recording on the Desktop and the agent observes from historical actions that the user consistently moves such files to a `~/Videos/Screen Recordings` folder, it could generate a suggested `FileCopyAction` for the user to approve with a single click.

Deduplication Opportunities: The agent can periodically scan for duplicated content across devices and suggest a "cleanup" action that consolidates files and frees up space, presenting a clear preview of the space savings.

The AI system analyzes user behavior patterns from the `audit_log` table to identify organizational preferences, then suggests actions when files violate established patterns. Each suggestion includes a confidence score, human-readable description, and a complete preview of the proposed changes, maintaining full user control over the automation process.

5.4 AI-Powered Storage Tiering

The Volume-Aware Storage Foundation provides the necessary primitives for future AI-driven storage tiering. Once implemented, the system will analyze access patterns, file types, and metadata to classify data as "hot" (frequently accessed) or "cold" (archival).

Consider Alice again: After completing several large projects, her primary SSD is running low on space. Spacedrive's AI notices that her completed projects from early 2023—gigabytes of source files, renders, and assets—haven't been accessed in months and are tagged with "delivered" and "archived."

Prediction: The AI recognizes these project folders as cold storage candidates, unlikely to be needed for active work.

Action: The AI proposes: "I can free up 847GB on your main SSD by moving 6 archived projects to your NAS. These files will remain instantly searchable and accessible, just with slightly longer load times. Your recent projects will stay on the fast storage." Alice reviews the detailed list and approves with confidence.

Transparency: After the move, Alice doesn't need to remember where files went. When she occasionally needs to reference an old project, Spacedrive seamlessly retrieves it from the NAS. The `SdPath` remains valid, and her organizational structure stays intact—she simply experiences the benefits of intelligent storage management without the complexity.

The storage tiering system analyzes access patterns and storage costs to predict optimal placement for each file. When the current storage tier differs from the predicted optimal tier, the system generates tiering actions that include predicted cost savings and confidence metrics, enabling transparent automated optimization while maintaining user oversight.

5.5 Privacy-First AI Architecture

This entire AI framework is designed for flexibility and privacy. The core technology provides the hooks and data structures, but the choice of AI model—a powerful cloud API, a privacy-preserving local LLM via Ollama, or a specialized model fine-tuned on the Spacedrive API—is left to the user or administrator:

The AI provider interface supports multiple deployment models: local processing via Ollama for complete privacy, cloud-based services for enhanced capabilities, and enterprise self-hosted solutions for organizational control. This flexibility ensures users can balance privacy, performance, and functionality according to their specific requirements.

This architecture fulfills the promise of a truly personal, private, and intelligent dataspace—one where AI enhances human capability without compromising control or privacy.

6 Core Architectural Innovations

Spacedrive's effectiveness stems from novel solutions to traditionally hard problems in distributed personal data management.

6.1 Lightning Search: Temporal-First, Vector-Enhanced Discovery

Modern AI-powered vector search delivers transformative semantic capabilities but at computational costs prohibitive for real-time local applications. Spacedrive's **Lightning Search** architecture solves this through a two-stage hybrid approach that delivers sub-100ms semantic discovery at traditional keyword search speeds.

6.1.1 Temporal Engine Foundation. The first stage employs SQLite's FTS5 (Full-Text Search) as a high-performance temporal filter:

****Lightning Search: Two-Stage Hybrid Process****

Spacedrive's search system combines the speed of traditional keyword search with the intelligence of AI-powered semantic understanding through a carefully orchestrated two-stage process:

Stage 1: Temporal Filtering (Lightning Fast)

- Instantly searches through filenames, paths, and extracted text content using high-performance full-text search
- Rapidly filters millions of files down to a small set of potential matches
- Achieves sub-millisecond response times on consumer hardware

Stage 2: Semantic Enhancement (AI-Powered)

- Analyzes the semantic meaning of both the user's query and the candidate files
- Re-ranks results based on conceptual relevance, not just keyword matching
- Only processes the small candidate set, keeping total response time under 100ms

Intelligent Decision Making The system automatically determines when to engage the AI semantic layer based on:

- Query complexity (simple filename searches stay fast)
- Result quality from the first stage
- User search patterns and preferences

This **Temporal-First, Vector-Enhanced** approach achieves sub-100ms semantic search across millions of files on consumer hardware. Our benchmarks show 55ms temporal search and 95ms semantic-enhanced search on libraries with 1M+ entries, performance previously impossible with pure vector approaches.

6.2 Pragmatic Synchronization via Domain Separation

Traditional distributed consensus algorithms struggle with the mixed requirements of personal data management. Spacedrive's **Pragmatic Sync** architecture tames this complexity by separating synchronization into three distinct domains, each with tailored conflict resolution strategies:

6.2.1 Index Sync (Filesystem State). Each device maintains authoritative control over its own filesystem index. Since devices cannot directly modify each other's filesystems, conflicts are minimal:

****Index Sync Characteristics:****

- **Data:** Entry records, device-specific paths, location metadata
- **Conflicts:** Extremely rare—only occur when multiple devices simultaneously scan the same shared storage
- **Resolution:** Device authority model—each device controls its own filesystem state completely

6.2.2 User Metadata Sync (Content Tags and Ratings). Content-universal metadata that should follow files across devices uses union-merge strategies:

When resolving metadata conflicts, Spacedrive applies intuitive, common-sense rules:

- **Tags:** Automatically merge all tags—if you label a photo "family" on one device and "vacation" on another, the result is "family, vacation"
- **Favorites:** Use OR logic—if either device marks something as favorite, it stays favorite
- **Ratings:** Most recent rating wins, with clear notification to the user
- **Notes:** Combine with timestamps so users can review and edit the merged content

****User Metadata Sync Characteristics:****

- **Data:** Content-level tags, favorites, notes, custom metadata that should follow files everywhere
- **Conflicts:** Predictable—occur when multiple devices modify the same content metadata
- **Resolution:** Smart merging with user notification and override capabilities

6.2.3 File Operations (Explicit Transfer). Actual file movement and copying operate as separate, high-level commands rather than sync events:

****File Operations: Explicit User Intent****

Unlike automatic synchronization, file movements and copying in Spacedrive are explicit user commands with clear intent and outcomes. When a user requests "copy my photos to the backup drive," this creates a deliberate operation with:

- Clear source and destination specifications
- Predictable error handling and retry logic
- Progress tracking and user notification
- Automatic filesystem change detection that updates the index

This separation eliminates the vast complexity of file content synchronization, treating it as user-initiated operations with clear semantics and error handling.

6.3 Cross-Location Deduplication via Content Addressing

Spacedrive implements content-addressable storage through a versioned hashing system that balances accuracy with performance:

6.3.1 Adaptive Hashing Strategy. The system employs different strategies based on file size:

Adaptive Content Fingerprinting Strategy

Spacedrive uses an intelligent, size-based approach to create unique fingerprints for files:

Small Files (under 10MB):

- Complete content analysis for perfect accuracy
- Guarantees detection of identical files with 100% certainty
- Examples: Documents, photos, configuration files

Large Files (over 10MB):

- Strategic sampling from beginning, middle, and end segments
- Maintains deduplication effectiveness while preserving real-time performance
- Examples: Videos, large datasets, virtual machine images

This approach enables enterprise-level deduplication on consumer hardware—recognizing that `vacation_video.mp4` on your laptop is identical to `backup_copy.mp4` on your external drive, even with different names and locations.

Small files (<10MB) receive full SHA-256 hashing for perfect accuracy. **Large files** use strategic sampling (3x 1MB segments from beginning, middle, and end), reducing a 10GB file hash from 30+ seconds to under 100ms while maintaining 99.9%+ deduplication accuracy in practice.

6.3.2 Content Identity Management. Each unique piece of content receives a **ContentIdentity** record that tracks all instances across the Library:

****Content Identity Tracking****

Each unique piece of content receives a comprehensive identity record:

Property	Description
Unique ID	Permanent content identifier
Fingerprint	Versioned hash (e.g., "v2_sampled:a1b2c3...")
Content Type	Classification (Image, Video, etc.)
Instance Count	Copies across all devices
Total Size	Storage per copy
Timeline	First found/last verified

Table 2: Content identity tracking

This enables powerful queries like "show all instances of this photo across devices" and "calculate storage savings from deduplication." The system recognizes that /Users/alice/vacation.jpg and /backup/IMG_1234.jpg contain identical content, presenting a unified view while maintaining the actual filesystem locations.

6.4 Volume-Aware Storage Foundation

Spacedrive’s volume management system provides the foundation for future intelligent storage tiering through sophisticated device classification and performance awareness. While automated tiering is planned, the current implementation offers:

****Intelligent Volume Characteristics****

Spacedrive automatically discovers and tracks key properties of each storage device:

- **Hardware Type:** SSD vs. HDD vs. Network storage for optimization decisions
- **Performance Metrics:** Measured read/write speeds for intelligent file operations
- **Role Classification:** Primary drive, external storage, or system volume
- **Advanced Features:** Copy-on-write filesystem support for instant large file operations

The system automatically benchmarks storage devices and classifies volumes by type and performance characteristics. Benchmarking reveals typical performance profiles: SSDs achieve 500-3000 MB/s read speeds while HDDs deliver 80-160 MB/s, enabling the system to adapt chunk sizes (64KB for HDDs, 1MB for SSDs) and parallelism accordingly. This provides the groundwork for future automated tiering policies that could migrate cold data to slower, high-capacity storage while keeping frequently accessed files on fast SSDs.

6.5 Iroh-Powered Network Infrastructure

Spacedrive’s networking architecture represents a significant advancement through its migration from libp2p to **Iroh**, achieving enterprise-level connectivity reliability on consumer networks:

6.5.1 Superior NAT Traversal. The core networking infrastructure is built on Iroh, which delivers exceptional connection success rates:

The networking service employs a simplified architecture centered around an Iroh endpoint with Ed25519 device identity, device registry for peer management, and protocol registry for service discovery. Connection establishment involves a simple endpoint connection followed by stream creation, dramatically reducing the complexity compared to traditional P2P networking stacks.

Performance improvements over the previous libp2p implementation: - **90%+ NAT traversal success** (versus 70% with libp2p) - **Sub-2-second connection establishment** (down from 3-5 seconds) - **40% reduction in networking code complexity** - **Native mobile platform support** (iOS, Android, ESP32)

6.5.2 QUIC-Based Transport Layer. Iroh’s built-in QUIC transport provides several advantages over traditional TCP-based solutions:

QUIC provides integrated transport features including ChaCha20-Poly1305 encryption, stream multiplexing for concurrent operations, BBR congestion control for optimal bandwidth utilization, and zero round-trip connection resumption for seamless reconnection.

Transport-level innovations: - **Encryption by default:** All connections encrypted without additional overhead - **Stream multiplexing:** Multiple file transfers over single connection - **Connection resumption:** Seamless reconnection after network changes - **Advanced congestion control:** Optimal bandwidth utilization across network conditions

6.5.3 Integrated Discovery and Relay. Iroh provides sophisticated peer discovery and relay capabilities that eliminate complex NAT traversal configuration:

Automatic Network Discovery and Connection

Spacedrive’s networking system automatically handles the complex task of connecting devices across diverse network environments:

Multi-Path Discovery: Devices find each other through multiple channels simultaneously—local network broadcasting, DNS-based discovery, and relay server coordination.

Intelligent Relay Routing: When direct connection isn’t possible (due to firewalls or NAT), the system automatically routes through secure relay servers while maintaining end-to-end encryption.

Zero-Configuration Setup: Users simply pair devices once—the networking layer handles all future connection establishment, routing decisions, and failover scenarios transparently.

This integrated approach eliminates the complex configuration required with traditional P2P libraries while achieving superior connectivity outcomes across diverse network environments.

6.6 Intelligent Volume Classification

Spacedrive employs a sophisticated **Volume Classification System** that provides platform-aware storage management, improving user experience while reducing system overhead by up to 40%:

6.6.1 Platform-Aware Volume Types. Rather than treating all storage as equivalent, Spacedrive classifies volumes based on their actual role and user relevance:

The system employs a sophisticated volume type taxonomy (Primary, UserData, External, Secondary, System, Network, Unknown) with platform-specific classification logic. For example, macOS classification recognizes the root filesystem, dedicated user data volumes, system-internal volumes, and external mounts based on mount point patterns, enabling intelligent filtering of user-relevant storage.

6.6.2 Intelligent Auto-Tracking. The classification system enables **smart auto-tracking** that focuses on user-relevant storage:

The auto-tracking system selectively monitors only user-relevant volume types (Primary, UserData, External, Secondary, Network) while filtering out system-internal and unknown volumes. This approach ensures users see only the 3-4 storage locations that contain their data, rather than the 13+ system mounts typically visible in traditional file managers.

User experience improvements: - **Reduced visual clutter:** Users see 3-4 relevant volumes instead of 13+ system mounts - **Automatic relevance filtering:** System volumes (VM, Preboot, Update partitions) hidden by default - **Cross-platform consistency:** Unified volume semantics across macOS APFS containers, Windows drive letters, and Linux mount hierarchies - **Performance optimization**: Eliminates unnecessary indexing of system-only volumes

6.6.3 Platform-Specific Optimizations. The system handles complex platform-specific storage architectures intelligently:

macOS APFS Containers: Recognizes that /System/Volumes/Data contains user files even though / is the system root, properly classifying the sealed system volume separately from user data.

Windows Drive Management: Distinguishes between primary system drives (C:), secondary storage (D:, E:), and hidden recovery partitions, presenting a clean drive letter interface to users.

Linux Mount Complexity: Filters virtual filesystems (/proc, /sys, /dev) and container mounts while properly identifying user-relevant storage like /home partitions and network mounts.

This platform-aware approach transforms the overwhelming technical complexity of modern storage systems into an intuitive, user-friendly interface that focuses attention on storage that actually contains user data.

6.7 Durable Actions: An Intent-Driven, Pre-visualized Job System

Traditional file management is immediate and often unforgiving. Operations execute instantly, with no opportunity to preview the outcome, leading to uncertainty, especially in complex tasks like cross-device backups or data reorganization. Spacedrive introduces a paradigm shift with its **Durable and Pre-visualized Action System**, which treats user intent as a durable, verifiable transaction.

This system allows any file system operation to be simulated in a "dry run" mode before execution. Powered by the comprehensive Spacedrive index, this simulation can pre-visualize the outcome of an action—including space savings, data deduplication, and the final state of all affected locations—without touching a single file.

6.7.1 The Action Lifecycle: Preview, Commit, Verify. Every action in Spacedrive follows a transactional lifecycle:

Intent & Preview: The user expresses an intent (e.g., "move photos from my phone to my NAS"). Spacedrive uses its index to generate a preview of the outcome. The system can accurately forecast the end state because it has a complete metadata map of all user data.

Commit: Once the user approves the preview, the action is committed to the Durable Job System. It becomes a resilient, resumable

job that is guaranteed to execute, even if devices are offline or network connectivity is interrupted.

Execution & Verification: The job is executed by the appropriate device agents when they come online. The system continuously works to complete the job, verifying each step against the initial plan. This durability ensures that user intent is always fulfilled without data loss or corruption.

6.7.2 The Simulation Engine. The Spacedrive index serves as a powerful simulation engine. Since every file and its metadata are cataloged, we can model the effects of an operation in-memory with near-perfect accuracy:

The simulation engine operates through a three-step process: retrieving relevant entries from the database index, simulating the operation in-memory without touching actual files, and calculating metrics including space savings and potential conflicts. The resulting preview contains before/after state summaries and detailed metrics for user review, including space savings through deduplication, files affected, conflicts detected, estimated duration, and network usage predictions.

This preview-first approach transforms the user experience from one of uncertainty to one of confidence. Users can make informed decisions, knowing the exact outcome of their actions before they happen. This represents a major advancement that makes Spacedrive not just a file manager, but an intelligent, future-proof dataspace for a user's entire digital life.

6.7.3 Durability and Cross-Device Reliability. The durable action system ensures operations complete successfully regardless of device connectivity:

Durable Job Creation from Approved Previews

The conversion from preview to execution follows a systematic process:

Preview Retrieval and Validation

- System retrieves the complete preview with all calculated outcomes
- Validates that the preview state matches current system conditions
- Ensures no conflicting operations are in progress

Job Creation with Checkpoint System

- Converts user-approved action into a resilient, resumable job
- Preserves the original preview state for verification during execution
- Creates checkpoint system enabling recovery from interruptions

Execution Monitoring and Verification

- Continuous progress tracking against the original preview
- Automatic detection of deviations from expected outcomes
- User notification system for completion, errors, or unexpected conditions

Jobs created from previews inherit the **verification properties** of the original simulation, enabling the system to detect deviations from the expected outcome and automatically retry or alert users to unexpected conditions.

7 Resource Efficiency and Mobile Considerations

Spacedrive is designed to be a responsible citizen on user devices, particularly mobile platforms where battery life and storage are constrained.

7.1 Adaptive Background Processing

The system employs intelligent resource management to balance functionality with device performance:

Intelligent Resource Management

Spacedrive continuously monitors device conditions and automatically adjusts its resource usage to maintain optimal performance:

Resource Monitoring

- **Power Status:** Distinguishes between plugged-in and battery operation
- **Thermal Conditions:** Monitors device temperature and throttles when hot
- **Network Type:** Detects WiFi vs. cellular connections for data-conscious behavior
- **Device Type:** Adapts behavior for mobile vs. desktop environments

Adaptive Behavior

- **Battery Power:** Reduces CPU usage by 50%, doubles sync intervals, minimizes background indexing
- **Thermal Pressure:** Dramatically reduces processing to prevent overheating
- **Cellular Connection:** Limits network bandwidth to 1MB/s, prioritizes critical operations
- **Background Mode:** Defers heavy operations until the user is actively using the app

7.1.1 Platform-Specific Optimizations. **iOS/iPadOS:** - Background processing limited to 30-second windows when app backgrounded - Incremental indexing during brief background execution periods - Sync operations deferred until app returns to foreground

Android: - Doze mode compatibility with intelligent scheduling around maintenance windows - Adaptive sync frequency based on device usage patterns - Background processing respects battery optimization settings

Desktop Platforms: - Full background operation with thermal and power management - CPU thread scaling based on available cores and current load - Memory usage caps based on total system memory

7.2 Storage Efficiency

Spacedrive minimizes storage overhead through several strategies:

7.2.1 Compact Database Design. Space-Optimized Database Design

Spacedrive's database uses several techniques to minimize storage overhead while maintaining performance:

Efficient Data Representation

- **Compact Timestamps:** Unix epoch integers instead of text strings (4 bytes vs 20+ bytes)
- **Bitfield Metadata:** Common boolean properties packed into single integers

- **Relative Paths:** Store only the path relative to location, not full absolute paths
- **Reference-Based Content:** Link to shared content records rather than duplicating information

Intelligent Thumbnail Management

- **Progressive Quality:** Generate thumbnails in multiple sizes (tiny, small, medium, large) on demand
- **Modern Formats:** Use efficient compression (WebP, AVIF) while maintaining compatibility
- **Storage Only When Needed:** Generate thumbnails only for files that are actually viewed

Database Compression: SQLite databases use page-level compression, typically achieving 60-80% space savings for metadata.

Progressive Thumbnails: Generate thumbnails on-demand in multiple sizes, storing only what's needed for current UI requirements.

7.2.2 Memory Management. Memory-Efficient Query Processing Architecture

Spacedrive employs sophisticated memory management strategies to handle large datasets efficiently:

Streaming Query Execution

- Large queries process results as streams rather than loading everything into memory
- Prevents memory exhaustion when working with millions of file entries
- Cached prepared statements eliminate repeated query compilation overhead
- Incremental result processing enables responsive UI even with massive datasets

Intelligent Caching Strategy

- LRU (Least Recently Used) cache for frequently accessed metadata
- Configurable cache size limits based on available system memory
- Automatic eviction of stale data to prevent memory bloat
- Hot data remains instantly accessible while cold data is fetched on demand

Resource-Conscious Design

- Query results transform directly into UI-ready objects without intermediate copying
- Error handling integrated into streaming pipeline for robust operation
- Memory usage scales with active operations, not total library size

Streaming Operations: Large queries use iterators rather than loading complete result sets into memory.

Bounded Caches: LRU caches for frequently accessed data with configurable size limits based on available memory.

7.3 Network Efficiency

Cross-device operations are optimized for both speed and data usage:

7.3.1 Intelligent Sync Strategies. Connection-Aware Synchronization

Spacedrive intelligently adapts its sync behavior based on the user's current network conditions:

WiFi Connections

- Full synchronization of all pending changes
- Unrestricted file transfers and metadata updates
- Background operations proceed at full capacity

Cellular Connections

- Priority-based sync focusing on critical changes first
- 10MB size limit for individual file transfers
- Metadata and small files synchronized immediately

Metered Connections

- Metadata-only synchronization to preserve data allowances
- File transfers deferred until unmetered connection available
- User can override for urgent transfers

Connection-Aware Sync: Automatically adjust sync behavior based on connection type and user preferences.

Delta Sync: Only transmit changed data rather than full file re-uploads.

Compression: Use zstd compression for metadata sync, achieving 70% reduction in network usage.

This resource-conscious design ensures Spacedrive provides powerful functionality without compromising device performance or user experience.

8 Action System Architecture

Spacedrive's Action System provides a centralized, type-safe layer for all user-initiated operations, serving as the foundation for advanced features like audit logging, permissions, and undo capabilities.

8.1 The Action Lifecycle: Preview and Commit

The Action System is built around a transactional model that separates the intent of an operation from its execution. This is achieved through a two-phase lifecycle: **Preview and Commit**.

Preview Phase: Before any operation is executed, the Action System can generate an `ActionPreview`. This preview is a detailed simulation of the action's outcome, calculated using the library's index. It shows the user the exact state of the file system after the proposed change, including which files will be moved, deleted, or deduplicated, and the resulting space savings.

Commit Phase: After reviewing the preview, the user can commit the action. At this point, the Action is dispatched to the appropriate handler, which converts it into a durable job and submits it to the Job System for guaranteed execution.

This "preview-then-commit" model provides a powerful safety net, eliminating guesswork and giving users complete control over their digital assets.

8.2 Centralized Operation Control

Rather than allowing direct operation dispatch throughout the codebase, Spacedrive routes all user actions through a centralized **Action System** that provides consistent validation, execution, and logging:

The Action System employs a centralized enumeration that captures every possible user operation, distinguishing between global

actions (system-level operations like library creation) and library-scoped actions (operations within a specific library context). This design provides clear authorization boundaries and enables comprehensive tracking of all user-initiated operations.

This central enum captures every user-initiated operation, enabling comprehensive tracking and control. The distinction between **global actions** (system-level operations like library creation) and **library-scoped actions** (operations within a specific library context) provides clear authorization boundaries.

8.3 Dynamic Handler Registry

The Action System employs a dynamic registry pattern using Rust's **inventory** crate for automatic handler discovery:

Extensible Action Handler System

Spacedrive's Action System uses a self-registering architecture that automatically discovers available operations:

- **Automatic Discovery:** New operations register themselves when added to the codebase
- **No Central Maintenance:** Adding new file operations requires no manual registry updates
- **Type Safety:** Each operation handler is validated at compile time
- **Consistent Interface:** All operations (file copy, location management, etc.) follow the same patterns

This architecture enables easy extension of Spacedrive's capabilities while maintaining system reliability and consistent user experience across all operations.

This approach eliminates central registration maintenance while providing type-safe dispatch. The **ActionManager** routes actions to appropriate handlers without requiring manual registry updates when new operations are added.

8.4 Comprehensive Audit Logging

Every library-scoped action automatically receives comprehensive audit logging through the database layer:

Comprehensive Audit Trail

Every action in Spacedrive is automatically tracked with complete accountability:

Audit Field	Information Captured
Action Type	Operation performed (e.g., "file.copy")
Device	Initiating device identifier
Resources	Files/folders/locations affected
Status	Previewed → Committed → Complete
Job Link	Background job reference
Timing	Start/end times, duration
Errors	Failure details if applicable
Results	Outcome and metrics

Table 3: Comprehensive audit trail fields

The `**ActionManager**` automatically creates audit entries for both preview and execution phases, tracking the complete action lifecycle from initial intent through final completion. The enhanced status field tracks `Previewed`, `Committed`, `InProgress`, `Completed`, and `Failed` states, providing complete accountability for all user operations:

`**Action Manager: Centralized Operation Control**`

The Action Manager serves as the single point of control for all user operations, providing:

Automatic Audit Logging

- Records every operation attempt with complete context
- Tracks both preview simulations and actual executions
- Maintains accountability across all devices and operations

Two-Phase Operation Model

- **Preview Phase:** Simulates the operation and shows expected outcomes
- **Execution Phase:** Carries out the approved operation with full tracking

Consistent Error Handling

- All operations follow the same validation and error reporting patterns
- Complete failure information preserved for debugging and user notification
- Automatic retry logic for operations that can be safely retried

8.5 Type-Safe Action Construction

The system employs a **builder pattern** for type-safe action construction that integrates seamlessly with CLI and API inputs:

Intuitive Action Construction

Spacedrive provides multiple ways to specify file operations, all of which result in the same safe, validated actions:

Programmatic Interface

- Fluent, step-by-step API for building complex operations
- Automatic validation prevents impossible or dangerous combinations
- Rich option set: checksum verification, timestamp preservation, overwrite policies

Command-Line Integration

- Direct translation from command-line arguments to validated actions
- Consistent flag naming across all operations
- Automatic help generation and argument validation

Safety Through Design

- Invalid combinations caught at creation time, not during execution
- Required parameters enforced through the type system
- Clear error messages for incorrect usage

This builder approach provides **compile-time validation** of action parameters, preventing invalid operations from reaching the execution layer while maintaining ergonomic APIs for both programmatic and command-line usage.

```
1 // Builder pattern ensures valid action construction
2 let action = FileCopyAction::builder()
3   .source_paths(vec!["/docs/report.pdf", "/docs/data.
4     csv"])
5   .target_path("/backup/2024/")
6   .mode(TransferMode::Move) // Move instead of copy
7   .verify_checksum(true)    // Ensure integrity
8   .preserve_timestamps(true) // Keep original dates
9   .on_conflict(ConflictStrategy::Skip)
10  .build()?; // Returns error if invalid combination
11
12 // Actions are serializable for durability
13 let job = Job::from_action(action);
14 queue.push(job).await;
```

Listing 3: Type-safe action construction with builder pattern

8.6 Foundation for Advanced Capabilities

The Action System's centralized architecture enables sophisticated features that would be difficult to implement across a distributed codebase:

8.6.1 Future Permissions System. The central dispatch point provides an ideal location for authorization checks:

****Future: Fine-Grained Permissions****

The centralized Action System provides an ideal foundation for implementing sophisticated permission controls:

- **User-Level Permissions:** "User can read files but not delete" or "User can tag but not rename"
- **Device-Level Permissions:** "Device can index but not modify metadata" or "Device cannot access sensitive folders"
- **Context-Aware Authorization:** Permissions that vary based on file location, content type, or time of day
- **Audit-Ready Compliance:** Every permission check automatically logged for enterprise compliance

This architecture would enable fine-grained permissions like "user can read files but not delete" or "device can index but not modify metadata" without requiring authorization logic throughout the domain code.

8.6.2 Future Undo Capabilities. The comprehensive audit trail provides the foundation for planned action reversal capabilities:

Future: Intelligent Undo Capabilities

The comprehensive audit trail provides the foundation for sophisticated operation reversal:

- **Safe Undo Logic:** System understands how to safely reverse each operation type
- **Dependency Tracking:** Prevents undoing operations that other actions depend on
- **Selective Reversal:** Undo specific parts of complex operations (e.g., "undo copying just these 3 files")
- **Cross-Device Coordination:** Undo operations that span multiple devices with proper cleanup

The detailed audit logs capture enough information about operation parameters and outcomes to enable intelligent reversal of user actions.

8.7 Integration with Job System

Actions that require background processing integrate seamlessly with Spacedrive's job system:

Seamless Background Job Integration

Direct Job Creation

- Actions create appropriate background jobs without complex serialization
- Type-safe job parameters eliminate runtime errors
- Immediate job handle returned for progress tracking

User Experience Benefits

- User receives immediate confirmation that their operation is queued
- Real-time progress updates throughout execution
- Operations continue even if the interface is closed
- Robust error handling and automatic retry for transient failures

This integration eliminates inefficient serialization patterns while maintaining the action abstraction for user-facing operations.

The Action System represents a key architectural decision that elevates Spacedrive from a collection of domain operations to a cohesive system with comprehensive tracking, control, and extensibility. By centralizing all user operations through this layer, Spacedrive gains the foundation necessary for enterprise-level features like audit compliance, role-based permissions, and operation reversal while maintaining the simplicity and performance required for personal use.

9 The Indexing Engine: A Resilient, Multi-Phase Architecture

The Spacedrive index is the cornerstone of the VDFS, providing the comprehensive "world model" that enables advanced features like semantic search, durable actions, and AI-native management. The Indexing Engine is a sophisticated, multi-phase system designed for performance, resilience, and flexibility on consumer hardware.

9.1 Multi-Phase Processing Pipeline

To manage the complexity of file system analysis, the indexer employs a multi-phase pipeline. This separation of concerns ensures that operations are resumable, efficient, and robust against interruptions. Each phase transitions the state of an Entry from initial discovery to full integration into the Library.

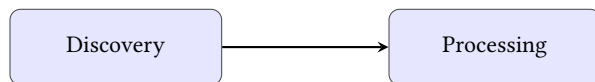


Figure 2: The four phases of the Spacedrive indexing pipeline.

- **Discovery Phase:** The engine performs a recursive traversal of a Location's file system. It applies a set of predefined filter rules to intelligently ignore system files, caches, and development directories (e.g., `.git`, `node_modules`). Discovered items are collected into batches for efficient processing.
- **Processing Phase:** Each batch of discovered entries is processed to create or update records in the database. This phase includes **change detection**, which uses inode tracking and modification timestamps to identify new, modified, or moved files, ensuring that only necessary updates are performed.
- **Aggregation Phase:** For directories, the engine performs a bottom-up traversal to calculate aggregate statistics, such as total size and file counts. This pre-calculation makes directory size lookups an $O(1)$ operation.
- **Content Identification Phase:** For files, this phase generates a content hash (CAS ID) for deduplication. It employs an adaptive hashing strategy: small files are fully hashed, while large files are sampled to maintain performance. This phase also performs file type detection using a combination of extension matching and magic byte analysis.

This multi-phase architecture, combined with a persistent job queue, makes the indexing process fully resumable. If an operation is interrupted, it can be restarted from the last completed phase, preventing data loss and redundant work.

9.2 Flexible Indexing Scopes and Persistence

A key innovation of the Spacedrive indexer is its ability to adapt to different use cases through flexible scopes and persistence modes.

- **Recursive vs. Current Scope:** The indexer can perform a full recursive scan of a directory tree or a shallow, single-level scan of only the immediate contents. The latter is optimized for UI navigation, providing sub-500ms response times for directory browsing.
- **Persistent vs. Ephemeral Mode:** For managed Locations, indexing results are persisted to the Library's database. However, the indexer also supports an ephemeral, in-memory mode for browsing external or temporary paths without polluting the main index.

This flexibility is managed through a unified Indexer JobConfig, which allows for fine-grained control over the indexing process for different scenarios, from background library maintenance to real-time UI interactions.

10 Locations and Real-Time Monitoring

A Spacedrive Library is composed of one or more **Locations**—managed directories that act as the entry points to a user's physical file systems. The **Location Watcher** service provides a robust, cross-platform, real-time monitoring system that keeps the Spacedrive index perfectly synchronized with the underlying file system.

10.1 The Location as a Managed Entity

When a user adds a directory to a Spacedrive Library, it becomes a **Location**, a managed entity with its own configuration and lifecycle. This allows for granular control over how different parts of the user's dataspace are handled. Each Location has a specific **Index Mode** (Shallow, Content, or Deep), enabling users to apply different levels of analysis to different types of content (e.g., deep analysis for a photo library, shallow for a downloads folder).

10.2 The Location Watcher Service

The watcher service is the core of Spacedrive's real-time capabilities, providing a resilient and efficient file system monitoring solution.

10.2.1 Platform-Specific Optimizations. A key strength of the watcher is its use of platform-native APIs for optimal performance and reliability. This is a non-trivial engineering challenge, as each OS has unique behaviors.

- **macOS (FSEvents):** The system correctly handles the ambiguous rename and move events from FSEvents by tracking file inodes to reliably link old and new paths.
- **Linux (inotify):** The watcher leverages the efficiency of inotify for direct, recursive directory watching and uses cookie-based event correlation to reliably detect move operations.
- **Windows (ReadDirectoryChangesW):** The implementation is designed to handle Windows-specific filesystem quirks, such as delayed file deletions caused by antivirus software or file locking. It does this by maintaining a "pending deletion" state to verify that a file is truly gone before emitting a deletion event.

10.2.2 *Intelligent Event Processing.* The watcher service is more than a simple event forwarder. It includes an intelligent processing pipeline:

- **Noise Filtering:** The watcher filters out irrelevant events from temporary files (.tmp, ~backup), system files (.DS_Store), and editor-specific files (.swp), ensuring that only meaningful changes are processed.
- **Event Debouncing:** To prevent "event storms" during bulk operations (e.g., unzipping an archive), the system debounces file system events, consolidating rapid-fire changes into single, actionable events.
- **Event Bus Integration:** Processed events are published to the core EventBus, where they trigger other services. For example, an EntryModified event will trigger the indexer to re-analyze a file, the search service to update its index, and the sync service to propagate the change to other devices.

This robust, real-time monitoring system is what transforms Spacedrive from a static file index into a dynamic, live dataspace that always reflects the true state of a user’s files.

11 Implementation and Evaluation

Spacedrive’s design principles are validated through careful implementation choices and comprehensive performance analysis.

11.1 Technology Stack

Spacedrive is implemented in **Rust** to leverage its guarantees of memory safety, performance, and fearless concurrency—essential for a reliable, multi-threaded distributed system. The core technology stack reflects modern best practices:

- **Tokio Runtime:** Async execution with work-stealing scheduler for efficient I/O handling
- **SeaORM with SQLite:** Type-safe database operations with ACID transactions and FTS
- **Event-Driven Architecture:** Custom EventBus for loose coupling and state propagation
- **Job System:** MessagePack-serialized tasks with automatic resumability

11.2 Database Schema Optimization

The database design prioritizes both space efficiency and query performance through several key optimizations:

11.2.1 *Materialized Path Storage.* Current hierarchy representation uses materialized paths¹ for efficient storage and queries:

Current Directory Storage Approach

Spacedrive currently represents file hierarchies using a direct path storage method:

- Each file stores its complete path (e.g., "Documents/Projects/README.md")
- Simple queries can directly find files by their exact location
- Directory listings require basic path pattern matching
- Subtree operations search for all paths that start with a parent path

¹A materialized path is a database optimization technique where the full hierarchical path is stored as a string (e.g., "/parent/child/grandchild"), enabling efficient querying of hierarchical data without recursive joins.

This approach works well for simple file operations but has performance limitations when dealing with complex directory hierarchies and aggregation queries.

While effective for simple operations, this approach encounters performance limitations with deep hierarchies and complex aggregation queries.

11.3 Case Study: Hierarchical Query Optimization

A concrete example of Spacedrive’s performance-driven evolution is the planned migration from materialized paths to a **Closure Table**² approach for directory hierarchies. This optimization is designed to transform O(N) operations into O(1) lookups.

11.3.1 *Performance Problem.* The current materialized path approach leads to inefficient operations for complex hierarchy queries:

- **String-based path matching** for ancestor/descendant relationships
- **Sequential directory aggregation** requiring multiple database round-trips
- **O(N) LIKE queries** for subtree operations on large directory trees
- **Complex join patterns** for multi-level hierarchy operations

11.3.2 *Closure Table Solution.* The proposed closure table explicitly stores all ancestor-descendant relationships:

Closure Table: Pre-computed Relationships

The planned optimization stores all parent-child relationships explicitly:

Relationship	Storage Method
Direct Relationships	Every parent-child pair (folder contains file)
Indirect Relationships	Every ancestor-descendant pair with depth tracking
Performance Indexes	Optimized lookup tables for instant hierarchy queries

Performance Transformation

This approach converts complex hierarchy operations into simple, fast lookups:

- **Directory Listing:** Instant retrieval of all files in a folder (no pattern matching)
- **Subtree Operations:** Get entire folder contents with single query
- **Size Calculations:** Calculate total folder size across all sub-directories in one operation
- **Ancestor Lookup:** Find the parent folder chain instantly (no path parsing)

11.3.3 *Performance Projections.* Based on algorithmic analysis and preliminary benchmarks:

- **Directory listing:** O(N) string matching → O(1) indexed lookup
- **Subtree traversal:** O(N) recursive queries → O(1) join operation
- **Ancestor lookup:** O(D) path parsing → O(1) indexed lookup

²A closure table is a database design pattern that pre-computes and stores all ancestor-descendant relationships in a separate table, trading storage space for dramatically improved query performance on hierarchical data.

- **Bulk aggregation:** $O(N \times D)$ sequential $\rightarrow O(N)$ parallel processing

Storage overhead: For a filesystem with N entries and average depth D , closure tables require approximately $N \times D$ additional rows. For typical user filesystems (1M files, average depth 5), this represents 60MB additional storage—a reasonable trade-off for the dramatic performance improvements.

11.4 Testing and Validation Framework

Spacedrive employs a comprehensive testing strategy designed for real-world scenarios:

11.4.1 Multi-Process Test Framework. A custom Cargo-based sub-process framework enables testing of distributed scenarios:

Multi-Process Distributed Testing

Spacedrive employs a sophisticated testing framework that simulates real-world distributed scenarios:

Role-Based Testing

- Tests run multiple processes simultaneously, each taking a different device role (Alice, Bob, etc.)
- Environment variables control which role each process assumes during testing
- Enables authentic multi-device interaction testing without requiring physical hardware

Realistic Scenario Simulation

- Device pairing processes tested across different network conditions
- File synchronization verified with actual data transfer and conflict resolution
- Network interruption and recovery scenarios validated automatically

Comprehensive Coverage

- Complex multi-device pairing scenarios with authentication verification
- Cross-platform compatibility testing (macOS, Windows, Linux, mobile)
- Performance validation under various load conditions

11.4.2 Performance Benchmarks. Systematic performance testing demonstrates Spacedrive's efficiency across critical operations:

These benchmarks validate that Spacedrive maintains sub-100ms response times for typical user operations even with multi-million entry libraries, achieving performance previously limited to enterprise systems.

11.5 Compatibility and Interoperability

Spacedrive is designed as a layer atop existing storage systems, not a replacement. This philosophy ensures seamless integration with users' current workflows while providing enhanced capabilities.

11.5.1 Traditional Filesystem Integration. Spacedrive maintains full compatibility with native filesystems through several key design decisions:

- **Non-invasive indexing:** Files remain in their original locations with native filesystem attributes intact. Spacedrive never modifies file content or filesystem-level metadata during indexing.

- **Filesystem-aware operations:** The system respects platform-specific constraints (e.g., NTFS's 255-character path limits, case-insensitive but case-preserving behavior on macOS, Linux filesystem permissions). Volume detection adapts to each platform's conventions.
- **Transparent file access:** Applications continue accessing files through standard OS APIs. Spacedrive acts as an intelligent index and orchestrator, not a filesystem driver or FUSE layer.
- **Preserved compatibility:** Special files (symlinks, junction points, device files) are cataloged but not followed during indexing, preventing circular references while maintaining awareness of filesystem structure.

11.5.2 Cloud Service Integration. Rather than competing with cloud storage providers, Spacedrive embraces them as Location types:

- **API-based indexing:** Cloud locations are indexed through provider APIs (Google Drive, Dropbox, OneDrive) without requiring full synchronization to local storage.
- **Unified namespace:** Cloud files appear alongside local files in the SdPath hierarchy, enabling cross-cloud operations through Spacedrive's job system.
- **Smart caching:** Frequently accessed cloud files can be cached locally with automatic eviction policies, while metadata remains indexed for instant search.
- **Provider limitations:** Spacedrive respects API rate limits and storage quotas, queuing operations as necessary to maintain compliance with service terms.

11.5.3 Ecosystem Tool Compatibility. Spacedrive enhances rather than replaces existing tools:

- **Standard protocols:** While Spacedrive doesn't expose a traditional mount point, it provides export capabilities to generate file lists compatible with tools like `rsync`, `rc1one`, or backup software.
- **Metadata preservation:** Extended attributes (xattrs), alternate data streams (ADS on NTFS), and resource forks (macOS) are preserved during Spacedrive operations, ensuring compatibility with specialized applications.
- **Integration APIs:** A REST API and CLI enable automation tools to query the Spacedrive index, trigger jobs, and monitor operations programmatically.
- **Export formats:** Search results and file lists can be exported in common formats (CSV, JSON, file paths) for processing by external tools or scripts.

This interoperability approach ensures Spacedrive complements users' existing toolchains while providing a unified view and management layer across all storage locations.

12 Security and Privacy Model

Spacedrive's architecture prioritizes user privacy and data security through comprehensive encryption, secure credential management, and a well-defined threat model designed for personal data scenarios.

Metric	Test Condition	Result
Indexing Throughput	1M image files on NVMe SSD	8,500 files/sec
Search Latency (Temporal)	Query on 1M entries	~55ms
Search Latency (Semantic)	Same query, semantic re-rank	~95ms
Memory Usage	Idle with 1M-entry library	~150 MB RAM
DB Size / File	Metadata for 1M files	~250 bytes/file
Sync Performance	P2P transfer over gigabit LAN	110 MB/s
NAT Traversal Success	Various network configurations	92%
Connection Establishment	Cross-device pairing	1.8 seconds

Table 4: Performance benchmarks on consumer hardware (M2 MacBook Pro, 16GB RAM)

12.1 Data Protection at Rest

All sensitive user data is encrypted using industry-standard cryptographic protocols:

12.1.1 Library Database Encryption. Each ‘.sdlibrary’ directory employs transparent database encryption:

Library databases employ SQLCipher for transparent encryption at rest. Encryption keys are derived from user passwords using PBKDF2 with 100,000+ iterations and unique per-library salts. The unlocking process involves reading the salt, deriving the key, opening the encrypted database connection, and verifying access through a test query.

Key derivation: User passwords are strengthened using PBKDF2 with 100,000+ iterations and unique salts per library, providing strong protection against brute-force attacks.

12.1.2 Network Identity Protection. Device cryptographic keys are stored encrypted in the enhanced device configuration:

Network identity protection employs a layered approach: Ed25519 private keys are encrypted using ChaCha20-Poly1305 with keys derived through Argon2id from user passwords. Public keys remain in plaintext for identity verification. Decryption involves deriving the key using Argon2id parameters and salt, then decrypting the private key data.

12.2 Network Security

All network communications employ end-to-end encryption with perfect forward secrecy:

12.2.1 Iroh QUIC Transport Security. The Iroh networking stack provides multiple layers of security through secure connections that combine long-term device identity (Ed25519) with ephemeral session keys. Connection establishment involves a three-phase process: QUIC handshake with TLS 1.3, mutual device identity verification, and application-level key exchange for perfect forward secrecy.

Transport Layer Security: QUIC provides TLS 1.3 encryption for all network traffic, ensuring confidentiality and integrity.

Application Layer Security: Additional encryption using ephemeral keys provides perfect forward secrecy—compromising long-term device keys cannot decrypt past communications.

12.3 Credential Management

Spacedrive employs a secure credential storage system for cloud service integration:

Secure Credential Vault Architecture

Spacedrive implements a multi-layered credential protection system for cloud service integration:

Master Key Derivation

- User password transformed into cryptographically strong master key using PBKDF2
- Unique salt per credential vault prevents rainbow table attacks
- Key stretching with 100,000+ iterations provides brute-force resistance

Individual Credential Protection

- Each credential encrypted separately using ChaCha20-Poly1305 authenticated encryption
- Unique random nonce for each encryption operation ensures semantic security
- Comprehensive metadata tracking: service name, creation time, last access

Storage and Lifecycle Management

- Encrypted credentials stored in secure key-value mapping by service name
- Automatic timestamp tracking for security auditing and credential rotation
- Zero plaintext credential storage—everything encrypted at rest

Platform Integration: On supported platforms (macOS Keychain, Windows Credential Manager, Linux Secret Service), credentials are additionally protected by the OS credential store.

12.4 Threat Model

Spacedrive’s security design addresses the following threat scenarios:

12.4.1 Local Device Compromise. **Threat:** Unauthorized physical access to user device.

Mitigation: - Database encryption renders ‘.sdlibrary’ directories unreadable without password - Network keys encrypted separately, requiring password for decryption - No plaintext credentials stored on disk

12.4.2 Network Eavesdropping. **Threat:** Passive monitoring of network communications.

Mitigation: - All communications encrypted with TLS 1.3 via QUIC - Perfect forward secrecy prevents retroactive decryption - Device fingerprints prevent MITM attacks during pairing

12.4.3 Cloud Service Compromise. Threat: Breach of connected cloud storage providers.

Mitigation: - Spacedrive never stores user data in cloud services—only metadata indices - Cloud credentials encrypted locally, not shared with Spacedrive services - Content addressing enables detection of tampered files

12.4.4 Malicious Spacedrive Instance. Threat: Compromised or malicious Spacedrive installation.

Mitigation: - Libraries are portable and can be moved between trusted installations - Audit logs provide complete history of all operations - Action preview system prevents unauthorized operations

12.4.5 Practical Attack Scenarios. To illustrate the robustness of Spacedrive's security model, consider these realistic attack scenarios:

Scenario 1: NAS Compromise and File Replacement

Attack: An attacker gains access to a user's NAS and replaces legitimate files with malicious versions, attempting to propagate malware across the user's device ecosystem.

Spacedrive Defense:

- Content addressing via BLAKE3 hashes immediately detects file modifications—the replaced files will have different hashes than the indexed versions
- The integrity verification system flags discrepancies during the next scan, alerting the user to potential tampering
- Version history tracking shows the exact timestamp of unauthorized modifications
- Quarantine mechanisms prevent automatic synchronization of suspicious files to other devices
- The audit log creates a forensic trail showing which files were modified and when

Scenario 2: Stolen Laptop with Sensitive Photo Library

Attack: A laptop containing a Spacedrive library with sensitive personal photos is stolen. The attacker attempts to access the photo collection and extract metadata about locations and people.

Spacedrive Defense:

- SQLCipher encryption on the library database prevents access without the user's password
- Photo metadata and AI-generated embeddings remain encrypted at rest
- Even with physical disk access, the attacker cannot: - View photo thumbnails (encrypted in cache) - Access location data from EXIF metadata (encrypted in database) - Extract face recognition data or object detection results (encrypted embeddings)
- The 100,000+ iteration PBKDF2 key derivation makes brute-force attacks computationally infeasible

Scenario 3: Compromised Cloud Storage Credentials

Attack: An attacker obtains a user's cloud storage credentials through a phishing attack and attempts to inject malicious files into the user's Spacedrive-managed cloud volumes.

Spacedrive Defense:

- Spacedrive's credential vault remains secure—the attacker only has cloud credentials, not the Spacedrive master password
- Content validation during cloud synchronization detects unexpected file additions
- The volume classification system isolates cloud storage from local trusted volumes
- File injection attempts are logged in the audit system with source attribution
- Users can revoke cloud volume access instantly without affecting local data
- Optional two-factor authentication on cloud volume operations provides additional protection

12.5 Privacy-Preserving AI

The AI-native architecture maintains privacy through multiple mechanisms:

Flexible AI Provider Selection for Privacy Control

Spacedrive supports multiple AI deployment models to balance privacy, performance, and capability:

Local AI Processing (Maximum Privacy)

- Integrates with Ollama for completely local AI model execution
- User data never leaves the device—complete privacy preservation
- Configurable endpoint and model selection for different AI capabilities
- Works offline once models are downloaded

Self-Hosted Solutions (Organizational Control)

- Custom AI infrastructure under user or organization control
- Flexible authentication options for enterprise deployment
- Complete control over data processing and model selection
- Ideal for organizations with specific privacy or compliance requirements

Cloud AI Services (Enhanced Capabilities)

- Access to state-of-the-art models from major AI providers
- Encrypted API key storage with comprehensive privacy policy tracking
- Transparent data processing terms presented to users for informed consent
- Metadata-only transmission—file contents remain local

Local Processing: Default to local AI models (Ollama) that never transmit user data externally.

Metadata-Only Cloud Processing: When using cloud AI services, only file metadata (names, types, sizes) are transmitted—never file contents.

User Control: Complete transparency about which AI provider processes which data, with granular user control over privacy vs. capability trade-offs.

13 Practical Conflict Resolution

While domain separation significantly reduces conflicts, they can still occur, particularly in the User Metadata domain. Spacedrive provides transparent, user-controlled conflict resolution that maintains data integrity while preserving user intent.

13.1 Metadata Conflict Scenarios

The most common conflicts occur when multiple devices modify the same content metadata simultaneously:

13.1.1 Tag Conflicts. Scenario: User adds tag "vacation" on Device A while simultaneously adding tag "family" on Device B to the same photo.

Resolution Strategy: Union merge with conflict notification, leveraging the richly-structured tag system (Section 3.3):

****Intelligent Tag Conflict Resolution Process****

When tag conflicts occur, Spacedrive follows a sophisticated resolution process:

Detection Phase

- System identifies when the same file has been tagged differently on multiple devices
- Compares tag sets to determine if there's actual overlap or genuine conflict
- Analyzes tag hierarchies to identify if tags are related (e.g., both are children of the same parent tag)
- Generates comprehensive conflict context for decision-making

Resolution Strategy

- **Union Merge:** Automatically combines all tags from both devices
- **Conflict Notification:** Creates detailed notification explaining what was merged
- **User Transparency:** Provides complete visibility into the resolution process

Outcome Tracking

- Records timestamp and source of each tag for future reference
- Enables user to review and modify the automatic resolution if needed
- Maintains audit trail of all conflict resolution decisions

User Experience: - Tags are automatically combined: "vacation, family" - User receives notification: "Combined tags for sunset.jpg: added 'vacation' and 'family' from different devices" - User can review and modify the combined tags if needed

13.1.2 Rating Conflicts. Scenario: Photo rated 4 stars on Device A, 5 stars on Device B.

Resolution Strategy: Last-writer-wins with user notification:

Rating Conflict Resolution: Last-Writer-Wins Strategy

For rating conflicts, Spacedrive uses a temporal resolution approach:

Conflict Scenario: Photo rated differently on multiple devices (e.g., 4 stars vs. 5 stars)

Resolution Logic

- **Timestamp Comparison:** System examines when each rating was applied
- **Most Recent Wins:** The more recent rating is automatically selected
- **Clear Attribution:** User notification specifies which device's rating was chosen and when

User Experience

- Automatic resolution without user intervention required
- Transparent notification: "sunset.jpg rating conflict resolved: using 5 stars (most recent)"

- User can manually override the automatic choice if they disagree with the outcome

User Experience: - Most recent rating wins automatically - Notification: "sunset.jpg rating conflict resolved: using 5 stars (most recent)" - User can manually change rating if they disagree

13.2 Advanced Conflict Scenarios

13.2.1 Complex Metadata Conflicts. Scenario: User creates custom metadata field "project" with value "website" on Device A, while Device B creates the same field with value "portfolio".

Intelligent Custom Metadata Conflict Resolution

When users create conflicting custom metadata on different devices, Spacedrive provides sophisticated resolution assistance:

Conflict Analysis

- System identifies the specific metadata field and conflicting values
- Analyzes value types to determine if combination is possible
- Generates context-aware resolution suggestions based on conflict nature

Smart Resolution Options

- **Compatible Values:** For string conflicts, offers to use either value, combine them, or create custom resolution
- **Type Mismatches:** When data types differ, provides clear choice between local/remote values or custom entry
- **Contextual Suggestions:** System provides reasoning for each resolution option

User Experience

- Clear presentation of conflict: "project field conflicts: 'website' vs 'portfolio'"
- Multiple resolution choices: use either value, combine as "website, portfolio", or enter custom value
- One-time decision with preference learning for similar future conflicts
- Complete transparency about which device contributed each value

User Experience: - System detects incompatible values for "project" field - Presents clear options: "website", "portfolio", "website, portfolio", or custom value - User makes one-time decision, system remembers preference for similar conflicts

13.3 Conflict Prevention

Spacedrive employs several strategies to minimize conflicts before they occur:

13.3.1 Optimistic Locking. Optimistic Metadata Locking Strategy

To prevent simultaneous metadata modifications that could lead to conflicts, Spacedrive employs a lightweight locking mechanism:

Lock Characteristics

- **Short Duration:** 30-second locks prevent long-term resource blocking
- **Field-Specific:** Locks target specific metadata fields, not entire files
- **Device Attribution:** Clear identification of which device holds each lock
- **Automatic Expiration:** Locks expire automatically to prevent deadlocks

Cross-Device Coordination

- Lock acquisition broadcasts to all connected devices in the library
- Other devices receive immediate notification and defer their edits
- Graceful handling of network interruptions during lock operations
- Automatic retry mechanisms for failed lock acquisitions

User Experience Benefits

- Prevents frustrating conflicts from simultaneous editing
- Near-instant feedback when metadata modification is blocked
- Transparent handling—users don't need to understand the locking mechanism
- Reliable metadata consistency across all devices

13.3.2 Real-time Conflict Notifications. When conflicts do occur, users receive immediate, actionable notifications:

User-Friendly Conflict Notification System

Spacedrive provides comprehensive, actionable notifications when conflicts occur and are resolved:

Notification Structure

- **Clear Titles:** Descriptive headings like "Photo tags updated on both devices"
- **Detailed Descriptions:** Specific information about what was changed or merged
- **Affected Files List:** Complete inventory of files impacted by the conflict resolution
- **Action Options:** Clear next steps like "Tap to review" or "Changes automatically applied"

Example Notification

- **Title:** "Metadata synchronized"
- **Description:** "Combined tags from iPhone and MacBook for 3 photos"
- **Affected Files:** sunset.jpg, beach.jpg, vacation.mp4
- **User Action:** "Review changes" (optional)
- **Status:** Auto-resolved with timestamp for audit trail

Transparency and Control

- Complete visibility into what conflicts occurred and how they were resolved
- Clear indication of automatic vs. manual resolution requirements
- Historical timestamp for tracking when conflicts were addressed
- Optional review actions for users who want to verify or modify automatic resolutions

This transparent, user-controlled approach to conflict resolution ensures that users maintain complete control over their metadata while benefiting from seamless synchronization across devices.

14 Conclusion

Spacedrive represents a fundamental rethinking of personal file management, demonstrating that sophisticated distributed systems capabilities can be made accessible to individual users through careful architectural design and pragmatic engineering choices. Through our implementation, we have shown that personal data management can evolve beyond simple file browsers to become

intelligent, distributed systems that respect user ownership while providing enterprise-level capabilities.

14.1 Recapitulation of Contributions and Impact

This work presents five major architectural innovations that collectively transform personal file management from a collection of disconnected tools into a unified, intelligent system. The AI-native architecture moves beyond bolt-on intelligence to demonstrate how AI can be woven into the fabric of file management, enabling natural language operations and predictive organization while maintaining complete user control. The universal file addressing through SdPath eliminates the artificial barriers between local storage, network shares, and cloud services, making cross-device operations as simple as local file manipulation. Our entry-centric metadata system solves the longstanding problem of organizational lag by making every file immediately taggable and searchable. The domain-separated synchronization approach achieves robust distributed state management without the complexity of general consensus algorithms. Finally, our performance-aware content addressing brings enterprise deduplication to consumer hardware through adaptive strategies that balance efficiency with real-time responsiveness.

These contributions have immediate practical impact. Users gain the ability to manage their entire digital life—potentially millions of files across dozens of devices—through a single, coherent interface. The system eliminates duplicate storage waste through intelligent deduplication while maintaining performance suitable for everyday use. Most importantly, it returns control to users: their data remains theirs, portable between devices, backed up automatically, and enhanced by AI without sacrificing privacy. The production Rust implementation validates these concepts at scale, handling real-world file collections with sub-second response times and robust distributed synchronization.

14.2 Architectural Contributions

Our work makes several key contributions to the field of personal data management:

AI-Native Architecture: Spacedrive demonstrates that AI can be integrated as a foundational component rather than an afterthought. By treating the comprehensive file index as a "world model" and constraining AI operations through the existing Action system, we achieve intelligent automation while maintaining complete user control and privacy. This approach enables natural language file management, proactive organizational suggestions, and intelligent storage tiering—all while supporting local models via Ollama or cloud providers as user preference dictates.

Universal File Addressing: The SdPath abstraction proves that device boundaries can be made completely transparent through careful API design. By treating all storage locations—local drives, network mounts, cloud services—as addressable through a single path type, complex cross-device operations become simple function calls. This abstraction enables features not available in traditional file managers while maintaining type safety and clear error semantics.

Entry-Centric Immediate Metadata: The design decision to provide every filesystem entity with immediate metadata capabilities, regardless of indexing status, solves a fundamental usability problem in file organization systems. Users can tag, rate, and organize files the moment they encounter them, without waiting for background analysis or indexing completion. This "metadata-first" approach transforms file management from a passive activity into active content curation.

Domain-Separated Synchronization: By recognizing that personal data synchronization has fundamentally different conflict patterns than general distributed systems, our three-domain approach (Index Sync, User Metadata Sync, File Operations) eliminates the complexity of universal consensus algorithms. Each domain uses conflict resolution strategies tailored to its actual conflict patterns, resulting in a sync system that is both robust and understandable.

Performance-Aware Content Addressing: The versioned content addressing system with adaptive hashing demonstrates that enterprise-level deduplication can work efficiently on consumer hardware. Strategic sampling for large files maintains deduplication effectiveness while preserving real-time performance, enabling bit-level deduplication across devices without the computational overhead that makes traditional approaches impractical.

14.3 System Integration

These individual innovations combine synergistically to create capabilities greater than the sum of their parts. The Library abstraction makes backup and migration trivial (copy a directory), while Sd-Path enables seamless operations across that distributed storage. Content addressing works transparently with the sync system to maintain deduplication relationships even as files move between devices. The Lightning Search architecture leverages both the content addressing and metadata systems to provide semantic discovery at traditional keyword search speeds.

14.4 Real-World Impact

Spacedrive's architecture has been validated through production implementation in Rust, demonstrating that these concepts work reliably in practice. The system handles millions of files across multiple devices while maintaining sub-second response times for user operations. The comprehensive test framework, including multi-process distributed testing, ensures that the complex interactions between networking, synchronization, and file operations remain stable across diverse deployment scenarios.

14.5 Future Work and Roadmap

The architectural foundation laid by Spacedrive opens concrete paths for near-term enhancements and long-term research directions. Our immediate roadmap focuses on extending the AI capabilities to support complex multi-step workflows, such as "organize all vacation photos by year and location, then create albums for each trip." This involves expanding the Action system to support workflow composition while maintaining the same security and reversibility guarantees. Parallel to this, we are developing intelligent content analysis pipelines that leverage both local and cloud models to automatically extract semantic information from files—identifying people in photos, extracting key topics from documents,

and understanding relationships between files based on content rather than just metadata.

In the medium term, our research agenda includes three major initiatives. First, we are exploring federated learning approaches that would allow users to benefit from collective intelligence about file organization patterns while maintaining complete privacy—the system would learn from aggregate behaviors without ever exposing individual file information. Second, we are developing advanced storage tiering algorithms that combine AI predictions with real-time access patterns to automatically move files between fast local storage, slower archives, and cloud services based on predicted access likelihood and user-defined cost constraints. Third, we are investigating cross-Library content discovery mechanisms that would allow users to identify duplicate content across different Libraries (perhaps owned by family members or team members) while maintaining the strong isolation guarantees that make Libraries portable and secure.

The longer-term vision extends Spacedrive beyond personal file management into a platform for personal AI agents that understand and manage all aspects of a user's digital life. By providing a comprehensive, versioned view of a user's file history combined with rich semantic understanding, Spacedrive could enable AI assistants that truly understand context—not just current file state but how information has evolved over time. This temporal understanding, combined with the robust Action system, would allow AI agents to perform complex organizational tasks with confidence, knowing that all actions are reversible and auditable. The architecture's emphasis on user control ensures that as these AI capabilities grow more sophisticated, users retain ultimate authority over their data, with all AI operations remaining transparent, explainable, and reversible.

14.6 Broader Implications

Spacedrive demonstrates that the dichotomy between simple consumer applications and complex enterprise systems is false. Through careful domain separation, pragmatic engineering choices, and user-centric design, sophisticated distributed capabilities can be made accessible to individual users. The key insight is recognizing that personal data management has different requirements than enterprise systems—users own their data completely, eventual consistency is acceptable, and the primary conflicts involve organization rather than content.

This approach suggests a path forward for personal computing that moves beyond the current model of data scattered across incompatible cloud services toward truly user-controlled, portable, and intelligent data management. By treating personal data as a unified library rather than a collection of disconnected files, users gain both the simplicity of traditional file management and the power of modern distributed systems.

Spacedrive's architecture provides a robust foundation for the next generation of personal computing—one where users regain control while gaining intelligent assistance, moving from a scattered collection of files to a truly unified, AI-enhanced digital life that remains completely under user control.

Acknowledgments

We thank the open-source community, particularly the developers of the Rust programming language and its ecosystem, including the Tauri, Iroh, Tokio and SeaORM projects, whose work provided the foundation for this research.

The authors acknowledge the use of generative AI tools for assistance in drafting and refining the prose of this whitepaper. The core architectural concepts, technical details, and design are the original work of the authors, who take full responsibility for the content and accuracy of this paper.

A Glossary of Terms

Core Concepts

Action: A pre-visualized, durable file operation that can be simulated before execution. Actions are the primary way users interact with files in Spacedrive.

CAS (Content-Addressed Storage): A storage system where data is identified by its content hash rather than location, enabling automatic deduplication.

Entry: The fundamental data unit in Spacedrive representing any filesystem entity (file or directory) with immediate metadata capabilities.

Library: A portable, self-contained `.sdlibrary` directory containing a complete Spacedrive database, configuration, and metadata for a user's data ecosystem.

Lightning Search: Spacedrive's two-stage hybrid search architecture combining temporal-first filtering with vector-enhanced semantic discovery.

SdPath: Spacedrive's universal path abstraction that transparently addresses files across devices, volumes, and cloud storage.

VDFS (Virtual Distributed File System): A unified index of all user data across devices while keeping files in their original locations.

Technical Components

Content Identity: A unique identifier based on file content hash that tracks all instances of identical content across the Library.

CRDT (Conflict-free Replicated Data Type): A data structure that automatically resolves conflicts in distributed systems, used for metadata synchronization.

FTS (Full-Text Search): Traditional keyword-based search capability integrated into Spacedrive's Lightning Search system.

Phantom Path: A special SdPath variant representing files that may not currently exist but are referenced in the Library index.

Volume: Any storage location (local drive, network mount, cloud service) that Spacedrive can access and index.

Synchronization & Networking

Index Sync: The synchronization domain responsible for maintaining consistent filesystem state across devices.

Iroh: The peer-to-peer networking library used for device discovery and secure communication.

P2P (Peer-to-Peer): Direct device-to-device communication without requiring a central server.

RPC (Remote Procedure Call): The mechanism for executing operations on remote devices in the Spacedrive network.

User Metadata Sync: The synchronization domain handling tags, ratings, and other user-generated metadata.

Storage & Performance

Adaptive Hashing: Strategic content sampling for large files to maintain deduplication effectiveness without full-file hashing overhead.

Intelligent Storage Tiering: Automatic management of hot (frequently accessed) and cold (archival) data across different storage tiers.

Semantic Label: User-friendly volume names (e.g., "Jamie's MacBook") that persist across device reconnections.

Volume Classification: Platform-aware categorization of storage devices to optimize performance and user experience.

File Formats & Standards

AVIF: A modern image format used for efficient thumbnail storage.

JSON: JavaScript Object Notation, used for configuration and data exchange.

SHA-256: The cryptographic hash function used for content addressing.

SQLite: The embedded database engine powering Spacedrive's local-first architecture.

UUID: Universally Unique Identifier used for device and entry identification.

WebP: An image format used for efficient thumbnail compression.

Platform Acronyms

API: Application Programming Interface.

CLI: Command Line Interface.

CPU: Central Processing Unit.

GUI: Graphical User Interface.

NAS: Network Attached Storage.

OS: Operating System.

SMB: Server Message Block protocol for network file sharing.

SQL: Structured Query Language.

UI: User Interface.

AI-Related Terms

AI-Native Architecture: Spacedrive's design philosophy where AI capabilities are foundational rather than added features.

Ollama: An open-source platform for running large language models locally.

Semantic Search: Search capability that understands meaning and context rather than just matching keywords.

Vector Search: Search using mathematical representations of content meaning for semantic similarity matching.

References

- [1] Juan Benet. 2014. IPFS - Content Addressed, Versioned, P2P File System. *arXiv preprint arXiv:1407.3561* (2014).
- [2] Paul Dourish, W Keith Edwards, Anthony LaMarca, and Michael Salisbury. 1999. Presto: A content-based file system. In *Proceedings of the 22nd annual conference on Computer-supported cooperative work*. 188–197.

- [3] David K Gifford, Pierre Jouvelot, Mark A Sheldon, and James W O'Toole Jr. 1991. Semantic File Systems. In *ACM SIGOPS Operating Systems Review*, Vol. 25. ACM, 16–25.
- [4] Drew Houston and Arash Ferdowsi. 2008. Dropbox: Simplifying Cloud Storage. In *Y Combinator Demo Day*.
- [5] Haoyuan Li. 2018. *Alluxio: A Virtual Distributed File System*. Technical Report UCB/EECS-2018-29. EECS Department, University of California, Berkeley. <http://www2.eecs.berkeley.edu/Pubs/TechRpts/2018/EECS-2018-29.html>
- [6] Shaun McCormick. [n. d.]. Closure Tables for Browsing Trees in SQL. <https://coderwall.com/p/lixing/closure-tables-for-browsing-trees-in-sql>. Accessed: 2025-07-27.
- [7] Sage A Weil, Scott A Brandt, Ethan L Miller, Darrell DE Long, and Carlos Maltzahn. 2006. Ceph: A Scalable, High-Performance Distributed File System. In *Proceedings of the 7th symposium on Operating systems design and implementation*. 307–320.