

G

gitit.in

A COMPREHENSIVE TECHNICAL ANALYSIS OF ARCHITECTURE,

*"From ancient counting stones to quantum algorithms—
every data structure tells the story of human ingenuity."*

LIVING FIRST EDITION

Updated December 30, 2025

© 2025 Mahdi

CREATIVE COMMONS • OPEN SOURCE

LICENSE & DISTRIBUTION

GIT INTERNALS: VERSION CONTROL SYSTEMS, SOFTWARE ARCHITECTURE, ALGORITHM ANALYSIS

A Living Architecture of Computing

Git Internals is released under the **Creative Commons Attribution-ShareAlike 4.0 International License** (CC BY-SA 4.0).

FORMAL LICENSE TERMS

Copyright © 2025 Mahdi

This work is licensed under the Creative Commons Attribution-ShareAlike 4.0 International License.

License URL: <https://creativecommons.org/licenses/by-sa/4.0/>

You are free to:

- **Share** — copy and redistribute the material in any medium or format for any purpose, even commercially.
- **Adapt** — remix, transform, and build upon the material for any purpose, even commercially.

Under the following terms:

- **Attribution** — You must give appropriate credit to Mahdi, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
- **ShareAlike** — If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
- **No additional restrictions** — You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

DISTRIBUTION & SOURCE ACCESS

Repository: The complete source code (LaTeX, diagrams, examples) is available at:
<https://github.com/m-mdy-m/algorithms-data-structures/tree/main/books/books>

Preferred Citation Format:

Mahdi. (2025). *Git Internals*. Retrieved from
<https://github.com/m-mdy-m/algorithms-data-structures/tree/main/books/books>

Version Control: This is a living document. Check the repository for the most current version and revision history.

WARRANTIES & DISCLAIMERS

No Warranty: This work is provided "AS IS" without warranty of any kind, either expressed or implied, including but not limited to the implied warranties of merchantability and fitness for a particular purpose.

Limitation of Liability: In no event shall Mahdi be liable for any direct, indirect, incidental, special, exemplary, or consequential damages arising from the use of this work.

Educational Purpose: This work is intended for educational and research purposes. Practical implementation of algorithms and techniques should be thoroughly tested and validated for production use.

TECHNICAL SPECIFICATIONS

Typeset with: L^AT_EX using Charter and Palatino font families

Graphics: TikZ and custom illustrations

Standards: Follows academic publishing conventions

Encoding: UTF-8 with full Unicode support

Format: Available in PDF, and LaTeX source formats

License last updated: December 30, 2025

For questions about licensing, contact: bitsgenix@gmail.com

Contents

Title Page	i
Contents	iii
List of Figures	v
List of Tables	v
Preface	vi
Acknowledgments	ix
Introduction	x
I Foundations: Version Control and Historical Context	1
1 The Version Control Problem	3
1.1 Collaboration Without Chaos	3
1.2 History and Accountability	3
1.3 Experimentation Without Fear	3
1.4 Formal Problem Statement	3
2 A Brief History of Version Control	4
2.1 Generation 0: Manual Version Control (1960s-1970s)	5
2.2 Generation 1: Local Version Control (1972-1982)	5
2.2.1 SCCS: Source Code Control System	5
2.2.2 RCS: Revision Control System	5
2.3 Generation 2: Centralized Version Control (1986-2000)	5
2.3.1 CVS: Concurrent Versions System	5
2.3.2 Subversion: Fixing CVS	5
2.3.3 Perforce and Commercial Solutions	5
2.4 Generation 3: Distributed Version Control (2000-Present)	5
2.4.1 BitKeeper and the Linux Kernel	5
2.4.2 The Birth of Git	5
2.4.3 Mercurial, Bazaar, and Alternatives	5
3 Centralized vs Distributed: The Paradigm Shift	6
3.1 The Centralized Model	6

3.1.1	Architecture and Workflow	6
3.1.2	Advantages	6
3.1.3	Fundamental Limitations	6
3.2	The Distributed Model	6
3.2.1	Architecture and Workflow	6
3.2.2	The Clone as First-Class Citizen	6
3.2.3	Advantages and New Possibilities	6
3.2.4	The Cost of Distribution	6
3.3	Mathematical Analysis of Collaboration Models	6
4	The Linux Kernel Crisis and Git's Birth	7
4.1	The BitKeeper Era (2002-2005)	7
4.2	The Crisis of April 2005	7
4.3	Design Goals and Non-Goals	7
4.4	The First Implementation (April-July 2005)	7
4.5	Linus's Design Philosophy	7
5	Core Terminology and Concepts	8
5.1	VCS vs SCM vs DVCS	8
5.2	Repository, Working Tree, and Index	8
5.3	Commits, Trees, and Blobs	8
5.4	Branches, Tags, and References	8
5.5	Merge, Rebase, and Cherry-Pick	8
5.6	Remote, Clone, and Fork	8
II	Core Architecture: Source Code Organization	9
6	Source Tree Organization	11
7	The Build System: Makefile Analysis	12
8	Module Architecture	13
9	Code Conventions and Style	14
10	The Entry Point: cmd_main Deep Dive	15
10.1	The Question of main vs cmd_main	15
10.2	String Vector Initialization	16
10.2.1	What is strvec?	16

10.2.2	Why STRVEC_INIT Macro?	16
10.2.3	Why Not Use malloc?	16
10.3	State Variables	17
10.3.1	The cmd Pointer	17
10.3.2	The done_help Flag	17
10.4	Command Name Extraction	18
10.4.1	Why Check !cmd?	18
10.4.2	The find_last_dir_sep Mystery	18
10.4.3	Why cmd = slash + 1?	19
III	Data Model: The Heart of Git	20
11	Content-Addressable Storage	22
12	The Four Object Types	23
13	The DAG: Mathematical Foundations	24
14	SHA-1: Cryptographic Hash Functions	25
15	Object Storage and Retrieval	26
16	Pack Files and Delta Compression	27
17	The .git Directory: Repository Anatomy	28

List of Figures

List of Tables

Preface

IT STARTED with a simple question: "How does Git actually work?"

Not the surface-level answer—"it's a distributed version control system." The real answer. The one that requires reading thousands of lines of C code, understanding content-addressable storage, analyzing directed acyclic graphs, and appreciating decades of software engineering wisdom.

This paper is my journey through that question.

The Problem with Most Git Resources

Most Git tutorials teach you commands: `git add`, `git commit`, `git push`. They're useful, but they don't explain *why*. Why does Git use SHA-1 hashes? Why is everything a DAG? Why are branches so cheap? Why does rebasing work the way it does?

Books like *Pro Git* give excellent overviews. But they don't show you the source code. They don't prove the complexity bounds. They don't analyze the design patterns or compare implementation alternatives.

Academic papers on version control systems exist, but they're often theoretical or focused on narrow aspects. None provide a comprehensive, ground-up analysis of Git's actual implementation.

That gap is what this paper fills.

What Makes This Different

This is not a user manual. This is not a quick-start guide. This is a **technical deep-dive into Git's internals**, written for:

- Software engineers who want to understand how sophisticated systems are built
- Computer science students learning about data structures, algorithms, and software architecture
- Anyone who's used Git and wondered "how does this actually work?"

Every claim is backed by:

- Source code analysis from the actual Git repository
- Mathematical complexity analysis with proofs
- Comparative studies with other version control systems
- Historical context and design rationale

The Journey

I started by reading Git's `git.c`—the entry point. Line by line. Function by function. Understanding why `cmd_main` exists instead of `main`. Why string vectors are initialized with macros. Why `find_last_dir_sep` is platform-dependent.

Then I dove into the data model. Blobs, trees, commits, tags. Content-addressable storage. The DAG structure. Pack files and delta compression. The mathematical foundations that make it all work.

Next came the command system. How 150+ commands are registered and dispatched. How options are parsed. How aliases are expanded. The design patterns that keep it maintainable.

Then algorithms: object lookup, tree traversal, merge strategies, diff algorithms. Complexity analysis. Trade-offs. Why certain choices were made.

Finally, the bigger picture: how Git compares to CVS, SVN, Mercurial. Why Git won. What makes it special. Where it's going.

How to Read This Paper

This paper is structured in Parts, each building on previous ones:

Part I establishes foundations—what version control is, why Git was created, terminology, and historical context.

Part II analyzes the core architecture—source code organization, build system, module structure.

Part III deeply examines the data model—objects, DAG, content-addressable storage, with full mathematical treatment.

Part IV dissects the command system—entry point, dispatch, option parsing, every line explained.

Part V provides algorithmic analysis—complexity proofs, trade-off analysis, performance characteristics.

Part VI identifies design patterns—which patterns Git uses and why they matter.

Part VII covers advanced features—branching, merging, rebasing, hooks, remotes.

Part VIII offers comparative analysis—Git vs. everything else, and why Git succeeded.

You can read linearly or jump to sections that interest you. Prerequisites are clearly marked. Code is thoroughly commented. Mathematics is rigorous but explained.

A Living Document

This paper continues to evolve. Git itself evolves—new features, performance improvements, design changes. My understanding deepens. Readers provide feedback.

The version you’re reading is a snapshot. Check the repository for updates, corrections, and additions.

Acknowledgments

To Linus Torvalds and the Git development community—for building something worth studying.

To the authors of *Pro Git*, *Version Control with Git*, and countless blog posts that helped me understand pieces of the puzzle.

To every developer who’s contributed to Git’s codebase, documentation, and tests. Your work made this analysis possible.

Let’s Begin

Enough preamble. Time to understand Git from the ground up.

Your Name
December 30, 2025

Acknowledgments

I would like to express my gratitude to everyone who supported me during the creation of this book. Special thanks to the open-source community for their invaluable resources and to all those who reviewed early drafts and provided feedback.

Introduction

VERSION CONTROL is fundamental to modern software development. Yet most developers use Git as a black box—memorizing commands without understanding the elegant engineering beneath.

This paper changes that.

The Central Question

How does Git actually work?

Not "what commands do I type" but "what happens when I type them." Not "how do I use Git" but "how is Git built."

This question leads to deeper ones:

- Why content-addressable storage?
- Why directed acyclic graphs?
- Why are branches cheap?
- How does merge work?
- What makes Git fast?
- How is the source code organized?
- What design patterns are used?
- How does it compare to alternatives?

Scope and Approach

This paper provides:

Complete Source Code Analysis—We read Git's C implementation line-by-line, explaining every design decision.

Mathematical Rigor—Every algorithm gets complexity analysis. Every data structure gets formal specification.

Engineering Insights—We identify design patterns, architectural principles, and software engineering practices.

Comparative Context—We compare Git to CVS, SVN, Mercurial, understanding what makes Git different.

Historical Perspective—We trace version control evolution from SCCS to Git, understanding why it exists.

What You'll Learn

By the end of this paper, you'll understand:

1. **The Data Model**—How Git represents history as a DAG of content-addressed objects
2. **The Architecture**—How Git's codebase is organized and why
3. **The Algorithms**—How merge, diff, and pack algorithms work, with complexity proofs
4. **The Command System**—How 150+ commands are implemented and dispatched
5. **The Design Patterns**—Which patterns make Git maintainable and extensible
6. **The Performance**—Why Git is fast and where bottlenecks exist
7. **The Philosophy**—The design principles and trade-offs that shaped Git

Prerequisites and Reading Strategy

Required Background:

- Basic programming knowledge (any language)
- Willingness to read C code
- Comfort with algorithmic thinking

Helpful but Not Required:

- Data structures (trees, graphs, hash tables)
- Algorithm analysis (Big-O notation)
- Basic Git usage
- Software architecture concepts

Reading Strategies:

Linear Reading—Start at Part I, work through sequentially. Best for comprehensive understanding.

Topic-Focused—Jump to sections of interest (e.g., Part III for data model, Part V for algorithms).

Code-First—Start with Part II and IV (architecture and command system), then fill in theory.

Every section is as self-contained as possible. Cross-references guide you to prerequisites.

Notation and Conventions

Throughout this paper:

- Monospace text indicates code, commands, or filenames
- **Bold monospace** indicates git commands
- *Italic monospace* indicates files
- *Italic monospace with //* indicates directories
- Mathematical notation follows standard conventions
- Complexity is expressed in Big-O, Big-Omega, Big-Theta notation

Why This Matters

Understanding Git internals makes you a better developer:

Debugging—When Git behaves unexpectedly, you'll know why.

Performance—You'll write commit histories that work with Git's model, not against it.

Architecture—Git demonstrates principles applicable to any large software system.

Confidence—Understanding *why* removes the fear of "don't do X or bad things happen."

Beyond Git, this paper teaches:

- How to analyze large C codebases
- How to apply algorithm analysis to real systems
- How to identify and understand design patterns in practice
- How to evaluate software architecture trade-offs

The Structure Ahead

Part I: Foundations Historical context, terminology, motivation

Part II: Architecture Source code organization, build system, modules

Part III: Data Model Objects, DAG, content-addressable storage

Part IV: Command System Entry point, dispatch, option parsing

Part V: Algorithms Complexity analysis, merge, diff, pack algorithms

Part VI: Design Patterns Patterns in Git's implementation

Part VII: Advanced Features Branching, merging, hooks, remotes

Part VIII: Comparative Analysis Git vs. other VCS

A Philosophical Note

Git is not just a tool—it's a **masterclass in software engineering**. It demonstrates:

- How to design for distribution and offline work
- How to make common operations fast
- How to maintain a decade-old codebase
- How to balance flexibility and simplicity
- How to build extensible systems

More fundamentally, Git shows how **choosing the right data model** solves problems that plague other systems. The DAG isn't a detail—it's the insight that makes everything else possible.

Ready?

Let's dive into the source code.

Let's understand the algorithms.

Let's see why Git became the standard.

Turn the page.

"The best way to understand a system is to build it.

The second-best way is to take it apart."

— UNKNOWN

Part I

Foundations: Version Control and Historical Context

BEFORE DIVING into *Git's implementation*, we need context. Why does version control exist? What problems does it solve? How did we get from punch cards to distributed workflows? This part establishes the foundation for everything that follows.

What You'll Learn:

- **The Problem Space:** What version control solves and why it matters
- **Historical Evolution:** From SCCS to *Git*, understanding the journey
- **Terminology:** VCS vs SCM, centralized vs distributed, and precise definitions
- **Git's Origin Story:** The Linux kernel crisis and Linus's response
- **Design Philosophy:** The principles that shaped *Git*

“Those who cannot remember the past are condemned to repeat it.”

— GEORGE SANTAYANA

Chapter 1

The Version Control Problem

Why do we need version control at all? This chapter examines the fundamental problems that version control systems solve, from coordinating multiple developers to maintaining project history. We establish precise definitions and explore the problem space before examining solutions.

1.1 Collaboration Without Chaos

1.2 History and Accountability

1.3 Experimentation Without Fear

1.4 Formal Problem Statement

Chapter 2

A Brief History of Version Control

Version control didn't appear fully-formed. It evolved through decades of iteration, each generation learning from the previous one's limitations. This chapter traces that evolution, from 1970s mainframe systems to modern distributed tools.

2.1 Generation 0: Manual Version Control (1960s-1970s)

2.2 Generation 1: Local Version Control (1972-1982)

2.2.1 SCCS: Source Code Control System

2.2.2 RCS: Revision Control System

2.3 Generation 2: Centralized Version Control (1986-2000)

2.3.1 CVS: Concurrent Versions System

2.3.2 Subversion: Fixing CVS

2.3.3 Perforce and Commercial Solutions

2.4 Generation 3: Distributed Version Control (2000-Present)

2.4.1 BitKeeper and the Linux Kernel

2.4.2 The Birth of Git

2.4.3 Mercurial, Bazaar, and Alternatives

Chapter 3

Centralized vs Distributed: The Paradigm Shift

The jump from centralized to distributed version control wasn't just a feature addition—it was a fundamental reconceptualization of what version control means. This chapter examines both paradigms in depth, understanding the trade-offs and why distribution matters.

3.1 The Centralized Model

3.1.1 Architecture and Workflow

3.1.2 Advantages

3.1.3 Fundamental Limitations

3.2 The Distributed Model

3.2.1 Architecture and Workflow

3.2.2 The Clone as First-Class Citizen

3.2.3 Advantages and New Possibilities

3.2.4 The Cost of Distribution

3.3 Mathematical Analysis of Collaboration Models

Chapter 4

The Linux Kernel Crisis and Git's Birth

Git wasn't created in a vacuum. It emerged from a specific crisis: the Linux kernel project needed a version control system, and nothing adequate existed. This chapter tells that story, examining the requirements that shaped Git's design.

4.1 The BitKeeper Era (2002-2005)

4.2 The Crisis of April 2005

4.3 Design Goals and Non-Goals

4.4 The First Implementation (April-July 2005)

4.5 Linus's Design Philosophy

Chapter 5

Core Terminology and Concepts

Precise terminology prevents confusion. This chapter defines key terms used throughout the paper, establishing a common vocabulary for discussing version control systems.

5.1 VCS vs SCM vs DVCS

5.2 Repository, Working Tree, and Index

5.3 Commits, Trees, and Blobs

5.4 Branches, Tags, and References

5.5 Merge, Rebase, and Cherry-Pick

5.6 Remote, Clone, and Fork

Part II

Core Architecture: Source Code Organization

NOW WE ENTER *Git's codebase*. Before analyzing specific algorithms or data structures, we need to understand how the code is organized. Where do files live? How is functionality partitioned? What's the build system? This part provides the map you'll need for everything that follows.

What You'll Learn:

- **Directory Structure:** How Git's source tree is organized
- **Build System:** How Git compiles and what happens during build
- **Module Organization:** How functionality is partitioned
- **Code Conventions:** C style, memory management, error handling
- **The Entry Point:** Where execution begins and why

"Programs must be written for people to read, and only incidentally for machines to execute." — HAROLD ABELSON

Chapter 6

Source Tree Organization

Chapter 7

The Build System: Makefile Analysis

Chapter 8

Module Architecture

Chapter 9

Code Conventions and Style

Chapter 10

The Entry Point: `cmd_main` Deep Dive

Finally, we read the code. This chapter analyzes Git's entry point line by line, understanding every design decision. Why `cmd_main` instead of `main`? Why are string vectors initialized with macros? Why is `find_last_dir_sep` platform-dependent? We answer everything.

10.1 The Question of `main` vs `cmd_main`

Why This Matters

Most C programs start with `int main(int argc, char **argv)`. Git doesn't. It uses `int cmd_main(int argc, const char **argv)`. Why this unusual choice?

```
1 int cmd_main(int argc, const char **argv)
2 {
3     struct strvec args = STRVEC_INIT;
4     const char *cmd;
5     int done_help = 0;
6     // ...
7 }
```

Listing 10.1: The actual entry point in git.c

The answer lies in Git's build system and testing infrastructure:

1. **Wrapper Programs:** Git has multiple entry points (`git`, `git-upload-pack`, etc.). The `main` function is in a thin wrapper that calls `cmd_main`.
2. **Testing:** Tests can call `cmd_main` directly without going through `main`, allowing better test isolation.

3. **Platform Compatibility:** Some platforms require special `main` handling (Windows Unicode, etc.). Separating concerns keeps `cmd_main` portable.

10.2 String Vector Initialization

Look at the first line of `cmd_main`:

```
1 struct strvec args = STRVEC_INIT;
```

10.2.1 What is `strvec`?

A **string vector** is Git's dynamic array for string storage:

```
1 struct strvec {
2     const char **v; // Array of string pointers
3     size_t nr;     // Number of strings
4     size_t alloc;  // Allocated capacity
5 };
```

Listing 10.2: `strvec` structure definition

10.2.2 Why `STRVEC_INIT` Macro?

```
1 #define STRVEC_INIT { NULL, 0, 0 }
```

Listing 10.3: `STRVEC_INIT` definition

This macro initializes:

- `v = NULL`: No allocation yet (lazy initialization)
- `nr = 0`: Empty array
- `alloc = 0`: No capacity allocated

Key Idea

Lazy Initialization Pattern: Don't allocate until needed. If the program errors out early (e.g., invalid arguments), we haven't wasted memory allocating arrays we never use.

10.2.3 Why Not Use `malloc`?

Stack allocation + macro initialization is faster and simpler than:

```

1 struct strvec *args = malloc(sizeof(struct strvec));
2 args->v = NULL;
3 args->nr = 0;
4 args->alloc = 0;

```

The macro approach:

- Avoids heap allocation
- Simplifies error handling (no need to free if we error out)
- Makes intent clear: "initialize this empty"

10.3 State Variables

```

1 const char *cmd;
2 int done_help = 0;

```

10.3.1 The cmd Pointer

Stores the command name extracted from `argv[0]`. Declared `const char *` because:

- We're pointing to existing memory (in `argv`)
- We won't modify the string itself
- Compiler can optimize knowing this

10.3.2 The done_help Flag

This is a **state machine flag**. It tracks whether we've already attempted to provide help to the user.

Flow:

1. User types invalid command
2. Git tries to suggest similar commands (via `help_unknown_cmd`)
3. If that fails, set `done_help = 1` to prevent infinite loop
4. If we still can't find a command, error out

Important

Without `done_help`, Git could loop infinitely trying to help with a malformed command.

10.4 Command Name Extraction

```

1 cmd = argv[0];
2 if (!cmd)
3     cmd = "git-help";
4 else {
5     const char *slash = find_last_dir_sep(cmd);
6     if (slash)
7         cmd = slash + 1;
8 }
```

Listing 10.4: Extracting command from argv[0]

10.4.1 Why Check !cmd?

POSIX allows argv[0] to be NULL (though rare). If it happens, default to showing help.

10.4.2 The `find_last_dir_sep` Mystery

Why is this a function (or macro) instead of inline `strrchr`?

```

1 #ifdef WIN32
2     // Windows: both / and \ are path separators
3     #define find_last_dir_sep(path) /* complex logic */
4 #else
5     // Unix: only / is path separator
6     #define find_last_dir_sep(path) strrchr(path, '/')
7 #endif
```

Listing 10.5: Platform-specific path handling

Cross-Platform Abstraction

Git runs on Windows, macOS, Linux, BSDs, etc. Path separators differ (\ vs /).

Abstracting this into a macro/function:

- Centralizes platform-specific logic
- Makes code readable (intent is clear)
- Prevents bugs from hardcoding separators

10.4.3 Why `cmd = slash + 1?`

If `argv[0]` is `"/usr/bin/git"`, we want just "git".

```
1 path: /usr/bin/git
2      ^      ^
3      |      |
4 path  slash
```

Listing 10.6: Pointer arithmetic explanation

`slash + 1` gives pointer to 'g', skipping the separator.

(Continuing with more sections...)

Part III

Data Model: The Heart of Git

THIS IS WHERE *Git* gets interesting. *The data model is Git's core insight—the idea that makes everything else possible. This part develops the mathematical foundations of content-addressable storage, directed acyclic graphs, and object representations.*

What You'll Learn:

- **Content-Addressable Storage:** Why hashing content is revolutionary
- **The Four Object Types:** Blob, tree, commit, tag—structure and purpose
- **DAG Theory:** Formal graph-theoretic foundations
- **SHA-1 Mathematics:** Collision probability and cryptographic properties
- **Object Storage:** How objects are stored and retrieved efficiently

"Show me your flowcharts and conceal your tables, and I shall continue to be mystified. Show me your tables, and I won't usually need your flowcharts; they'll be obvious." — FRED BROOKS

Chapter 11

Content-Addressable Storage

Chapter 12

The Four Object Types

Chapter 13

The DAG: Mathematical Foundations

Chapter 14

SHA-1: Cryptographic Hash Functions

Chapter 15

Object Storage and Retrieval

Chapter 16

Pack Files and Delta Compression

Chapter 17

The .git Directory: Repository Anatomy