

G

gitit.in

ENGINEERING A VERSION CONTROL SYSTEM

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

LIVING FIRST EDITION

Updated December 31, 2025

© 2025 Mahdi

CREATIVE COMMONS • OPEN SOURCE

LICENSE & DISTRIBUTION

GIT INTERNALS: SOFTWARE ENGINEERING, ARCHITECTURE, DESIGN PATTERNS

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 31, 2025

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

Contents

Title Page	i
Contents	iii
Preface	vi
Introduction	ix
I Context	1
1 Version Control Evolution	3
2 Git's Origin Story	4
3 Distributed vs Centralized	5
4 Git's Philosophy	6
II Architecture	7
5 Codebase Structure	9
6 Command Dispatch System	10
7 Platform Abstraction	11
8 The Entry Point and Initialization	12
III Memory and Resource Management	13
9 Memory Allocation Patterns	15
10 Memory Pools and Arenas	16
11 Reference Counting	17
12 Buffer Management	18
13 Object Cache	19
14 Memory Mapping	20
IV Data Structures and Algorithms	21
15 Hash Tables	23
16 Dynamic Arrays	24
17 Linked Lists	25

18	Trees and Tries	26
19	Graphs and DAGs	27
20	Bitmaps and Bit Manipulation	28
21	Bloom Filters	29
22	Union-Find (Disjoint Set)	30
V	The Object Model	31
23	Content-Addressable Storage	33
24	Object Types and Structure	34
25	Object Serialization	35
26	Loose Objects	36
27	References System	37
28	The Index (Staging Area)	38
29	RefLog	39
VI	Pack Files and Compression	40
30	Pack File Format	42
31	Pack Index Format	43
32	Delta Compression	44
33	Deltification Process	45
34	Multi-Pack Index (MIDX)	46
35	Cruft Packs	47
36	Pack Protocol	48
37	Pack Window Management	49
VII	String Handling and Parsing	50
38	String Buffer (strbuf)	52
39	String Lists and Vectors	53
40	Path Handling	54
41	Config File Parsing	55
42	Command Line Parsing	56

43	Object Format Parsing	57
44	Diff Format Parsing	58
VIII	File System and I/O	59
45	File I/O Patterns	61
46	Directory Traversal	62
47	File System Watching	63
48	Working Tree State	64
49	Lock Files	65
50	Temporary Files and Cleanup	66
IX	Core Algorithms	67
51	Graph Traversal	69
52	Merge Base Computation	70
53	Three-Way Merge	71
54	Merge Strategies	72
55	Diff Algorithms	73
56	Rename Detection	74
57	Tree Comparison	75
58	Object Lookup	76
59	Reachability Algorithms	77
60	History Simplification	78
61	Blame Algorithm	79
62	Packfile Generation	80
X	Advanced Engineering Concepts	81
63	Immutability and Its Consequences	83
64	Copy-on-Write Patterns	84
65	Lazy Loading and Evaluation	85
66	State Machines	86
67	Callback and Iterator Patterns	87

68	Bitmap Indices	88
69	Prefix Compression	89
70	Sparse Data Structures	90
71	Atomic Operations	91
72	Error Handling Strategies	92
73	Extension and Plugin Architecture	93
XI	Performance and Optimization	94
74	Operation Complexity	96
75	Caching Strategies	97
76	Parallelization	98
77	I/O Optimization	99
78	Network Optimization	100
79	Geometric Repacking	101
80	Commit Graph Optimization	102
81	Profiling and Benchmarking	103
XII	Security and Integrity	104
82	Cryptographic Hashing	106
83	Object Integrity	107
84	Input Validation	108
85	GPG Integration	109
86	Attack Surfaces	110
87	Fuzzing and Testing	111
XIII	Advanced Features	112
88	Rebasing Internals	114
89	Submodules	115
90	Worktrees	116
91	Partial Clone	117
92	Sparse Checkout	118

93	LFS Integration	119
94	Garbage Collection	120
XIV	Comparisons and Context	121
95	Git vs Mercurial	123
96	Git vs Subversion	124
97	Git vs Perforce	125
98	Modern Alternatives	126
	Acknowledgments	127

Preface

I HAD JUST FINISHED a project and wanted to start something new. I was working on gix—a Git wrapper I’d been thinking about—and realized I didn’t actually understand Git well enough.

Sure, I could use it. I knew the commands. But how does it work? Why is it designed this way? What makes it different from other version control systems?

I read some articles. Looked at Pro Git’s internals chapter. Found a small booklet on Git internals. But I wanted more. I wanted to see the actual code. The decisions. The trade-offs.

So I cloned the Git repository and started reading.

Why I’m Writing This

At first, I was just making notes for myself. Trying to understand the codebase. Mapping out how things connect.

Then I thought maybe I could write an article about it. Share what I learned.

But the more I read, the more I found. This wasn’t article material. There was too much to say.

So here we are. A book about Git’s internals from an engineering perspective.

What This Book Actually Is

This is me reading Git’s source code and explaining what I find.

Not line by line—that would be boring and pointless. But the important parts. The clever solutions. The patterns that keep showing up. The reasons things are built the way they are.

I care about software engineering. How systems are designed. Why certain approaches work. What problems they solve. Git is interesting because it does things differently and those differences matter.

Who This Is For

You if you've used Git but want to understand how it works underneath.
You if you're interested in system design and want to see a real, mature codebase.
You if you like reading code and learning from it.
You don't need to be a C expert. I'll explain the tricky bits. You don't need to know Git internals already. That's what this book is about.

What This Isn't

This isn't a Git tutorial. I'm not teaching you how to use Git. There are better resources for that.
This isn't comprehensive. Git's codebase is massive. I'm focusing on the core ideas and the interesting engineering choices.
This isn't perfect. I'm still learning this stuff. I'll probably get some things wrong. That's okay.

How I'm Approaching This

I started with Git's first commit. The original implementation by Linus Torvalds. It's small and shows the core model clearly.
Then I worked forward, looking at how the codebase evolved. What got added. What got changed. Why.
I'm trying to explain not just what the code does but why it's designed that way. What problem was being solved? What alternatives existed? What trade-offs were made?

This Is Ongoing

I'm still reading the code. Still learning. This book reflects what I understand right now.
Maybe in six months I'll understand more and want to revise some sections. That's how learning works.
Consider this a snapshot of the journey, not the final destination.

Why Bother?

Because Git is everywhere and most people don't understand it.
Because understanding your tools makes you better at using them.
Because Git's design has lessons that apply beyond version control.

Because it's interesting.

Let's start.

Mahdi

2025

Introduction

EVERY DESIGN decision has a reason. Every data structure solves a problem.
Every algorithm makes a trade-off.

When you look at mature software, you're seeing years of evolution. Hundreds of developers. Thousands of decisions. Some good, some less good, all made for specific reasons at specific times.

Git is like that. It's been around since 2005. It's used by millions. It's been optimized, refactored, extended. And understanding it teaches you about software engineering.

What You're Going To Learn

This book walks through Git's internals from an engineering perspective.

Not how to use Git. How Git is built.

We'll look at the data model and why it's clever. The codebase structure and how it evolved. The algorithms and their complexity. The design patterns that keep appearing. The performance optimizations. How it all fits together.

And more importantly—why. Why these choices? What problems do they solve? What are the trade-offs?

What This Book Isn't

This is not a Git tutorial. If you need to learn Git commands, read Pro Git or the official documentation.

This is not complete reference. Git has over 2000 files and hundreds of commands. I'm covering the important parts, the core insights, the interesting engineering.

This is not academic. I'll talk about algorithms and complexity, but I'm not writing proofs or formal specifications. This is practical engineering.

This is not flawless. I'll make mistakes. I'll miss things. That's fine. The goal is understanding, not perfection.

Why Git?

Because you probably use it every day. Most developers do.

Because it's well-designed at its core. The data model is elegant even if the command interface is messy.

Because it's mature. Twenty years of evolution. You can learn from that evolution—what worked, what didn't, what got changed and why.

Because it does things differently. Most version control systems work one way. Git works another way. Understanding why teaches you about design choices.

Prerequisites

You should know basic programming. Any language is fine. Concepts like variables, functions, data structures.

You should know basic Git usage. Add, commit, push, merge. Not expert level, just the basics.

You should know basic data structures. Trees, graphs, hash tables. Nothing fancy.

You don't need to know C deeply. I'll explain the important parts.

You don't need to know advanced Git internals. That's what this book teaches.

You don't need algorithm analysis background. I'll teach what matters.

How To Read This

Start with Part I if you want historical context. Skip it if you already know about version control evolution and why Git was created.

Parts II and III are essential. They cover architecture and the data model. Read those.

After that, jump around based on your interests. Want to understand merging? Go to Part VI. Curious about performance? Part IV. Interested in patterns? Part V.

Code examples are simplified for clarity but accurate in concept. When I show code, I'll tell you if I've simplified it.

The Core Ideas

Most version control systems think in terms of changes. They track what changed between versions. Git doesn't.

Git thinks in terms of snapshots. It stores the complete state of your project at each commit. Not the difference—the whole thing.

That one decision changes everything. It makes branching cheap. It makes merging easier. It makes Git fast. But it also creates challenges. How do you store thousands of snapshots efficiently? How do you transfer them over networks? How do you handle conflicts?

Git solves these problems in specific ways. We're going to look at how and why.

Content-Addressable Storage

Everything in Git is identified by its content hash. Files aren't named by their path. They're named by their SHA-1 hash.

If two files have the same content, they have the same hash. Git only stores them once. Automatic deduplication.

This seems weird at first but it's powerful. Git can tell if two things are identical by comparing hashes. It can detect corruption by verifying hashes. It can reference any object by its hash.

The entire system builds on this idea.

Starting Point

We'll start with foundations. What is version control? Why does Git exist? What problems does it solve?

Then architecture. How is the codebase organized? Where are things located?

Then the data model. This is the heart of Git. Objects, hashes, the DAG structure.

Then we go deeper. Storage, algorithms, patterns, performance.

By the end, you'll understand how Git works. Not just what commands do, but how they're implemented and why.

A Word on Engineering

Good engineering isn't about being clever. It's about solving problems effectively. Sometimes the best solution is simple. Sometimes it's complex because the problem is complex. Sometimes it's ugly because reality is messy.

Git has all of these. Simple elegant parts. Complex necessary parts. Ugly historical parts.

Learning to distinguish between them—to see which complexity is essential and which is accidental—that's what studying real systems teaches you.

Let's start.

Part I

Context

Before diving into Git's engineering, you need to understand why it exists and what problems it solves.

This part is quick. Just enough background to make sense of the decisions that follow.

Chapter 1

Version Control Evolution

Chapter 2

Git's Origin Story

Chapter 3

Distributed vs Centralized

Chapter 4

Git's Philosophy

Part II

Architecture

How is a 20-year-old codebase organized? Where does functionality live? How do pieces connect? This part maps the terrain before we dive into details.

Chapter 5

Codebase Structure

Chapter 6

Command Dispatch System

Chapter 7

Platform Abstraction

Chapter 8

The Entry Point and Initialization

Part III

Memory and Resource Management

Manual memory management in C. How Git allocates, tracks, and frees memory without leaking or crashing. The unglamorous foundation that everything else relies on.

Chapter 9

Memory Allocation Patterns

Chapter 10

Memory Pools and Arenas

Chapter 11

Reference Counting

Chapter 12

Buffer Management

Chapter 13

Object Cache

Chapter 14

Memory Mapping

Part IV

Data Structures and Algorithms

The building blocks. Hash tables, lists, trees, graphs. How Git implements them, optimizes them, and uses them everywhere.

Chapter 15

Hash Tables

Chapter 16

Dynamic Arrays

Chapter 17

Linked Lists

Chapter 18

Trees and Tries

Chapter 19

Graphs and DAGs

Chapter 20

Bitmaps and Bit Manipulation

Chapter 21

Bloom Filters

Chapter 22

Union-Find (Disjoint Set)

Part V

The Object Model

Git's core insight. Content-addressable storage. Four object types. How they relate. Why this model is powerful. This is the heart of everything.

Chapter 23

Content-Addressable Storage

Chapter 24

Object Types and Structure

Chapter 25

Object Serialization

Chapter 26

Loose Objects

Chapter 27

References System

Chapter 28

The Index (Staging Area)

Chapter 29

RefLog

Part VI

Pack Files and Compression

How Git stores thousands of objects efficiently. Delta compression. Pack formats. Transfer protocols. The engineering that makes Git fast.

Chapter 30

Pack File Format

Chapter 31

Pack Index Format

Chapter 32

Delta Compression

Chapter 33

Deltification Process

Chapter 34

Multi-Pack Index (MIDX)

Chapter 35

Cruft Packs

Chapter 36

Pack Protocol

Chapter 37

Pack Window Management

Part VII

String Handling and Parsing

C doesn't have good string handling. Git built its own. How it manages strings, parses formats, and handles text safely.

Chapter 38

String Buffer (strbuf)

Chapter 39

String Lists and Vectors

Chapter 40

Path Handling

Chapter 41

Config File Parsing

Chapter 42

Command Line Parsing

Chapter 43

Object Format Parsing

Chapter 44

Diff Format Parsing

Part VIII

File System and I/O

Git interacts with the file system constantly. Reading, writing, watching for changes. How it does this efficiently and correctly.

Chapter 45

File I/O Patterns

Chapter 46

Directory Traversal

Chapter 47

File System Watching

Chapter 48

Working Tree State

Chapter 49

Lock Files

Chapter 50

Temporary Files and Cleanup

Part IX

Core Algorithms

The complex parts. Graph algorithms. Diff. Merge. Tree comparison. Where algorithmic choices determine performance and correctness.

Chapter 51

Graph Traversal

Chapter 52

Merge Base Computation

Chapter 53

Three-Way Merge

Chapter 54

Merge Strategies

Chapter 55

Diff Algorithms

Chapter 56

Rename Detection

Chapter 57

Tree Comparison

Chapter 58

Object Lookup

Chapter 59

Reachability Algorithms

Chapter 60

History Simplification

Chapter 61

Blame Algorithm

Chapter 62

Packfile Generation

Part X

Advanced Engineering Concepts

Patterns, tricks, and techniques that appear throughout the codebase. How Git handles complexity, maintains performance, and stays extensible.

Chapter 63

Immutability and Its Consequences

Chapter 64

Copy-on-Write Patterns

Chapter 65

Lazy Loading and Evaluation

Chapter 66

State Machines

Chapter 67

Callback and Iterator Patterns

Chapter 68

Bitmap Indices

Chapter 69

Prefix Compression

Chapter 70

Sparse Data Structures

Chapter 71

Atomic Operations

Chapter 72

Error Handling Strategies

Chapter 73

Extension and Plugin Architecture

Part XI

Performance and Optimization

Git is fast. Not by accident. Specific optimizations, careful profiling, constant attention to performance. How it stays fast as repositories grow.

Chapter 74

Operation Complexity

Chapter 75

Caching Strategies

Chapter 76

Parallelization

Chapter 77

I/O Optimization

Chapter 78

Network Optimization

Chapter 79

Geometric Repacking

Chapter 80

Commit Graph Optimization

Chapter 81

Profiling and Benchmarking

Part XII

Security and Integrity

Git's "trust but verify" model. Cryptographic hashing. Attack surfaces. How Git protects data integrity and handles malicious input.

Chapter 82

Cryptographic Hashing

Chapter 83

Object Integrity

Chapter 84

Input Validation

Chapter 85

GPG Integration

Chapter 86

Attack Surfaces

Chapter 87

Fuzzing and Testing

Part XIII

Advanced Features

Features built on the core. More complex because they combine multiple primitives. Understanding implementation reveals trade-offs.

Chapter 88

Rebasing Internals

Chapter 89

Submodules

Chapter 90

Worktrees

Chapter 91

Partial Clone

Chapter 92

Sparse Checkout

Chapter 93

LFS Integration

Chapter 94

Garbage Collection

Part XIV

Comparisons and Context

Git made specific trade-offs. Understanding what other systems chose differently teaches you about design space and constraints.

Chapter 95

Git vs Mercurial

Chapter 96

Git vs Subversion

Chapter 97

Git vs Perforce

Chapter 98

Modern Alternatives

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.