# Git

Version control systems (VCSs) - tools used to track changes to source code - help maintain a history of changes - facilitate collaboration. - track changes to a folder and its contents in a series of snapshots, where each snapshot encapsulates the entire state of files/folders within a top-level directory. - store metadata: who created each snapshot, messages associated with each snapshot

Modern VCSs let you easily answer questions like:

- Who wrote this module?
- When was this particular line of this particular file edited? By whom? Why was it edited?
- Over the last 1000 revisions, when/why did a particular unit test stop working?

While other VCSs exist, **Git** is the de facto standard for version control. This XKCD comic captures Git's reputation:

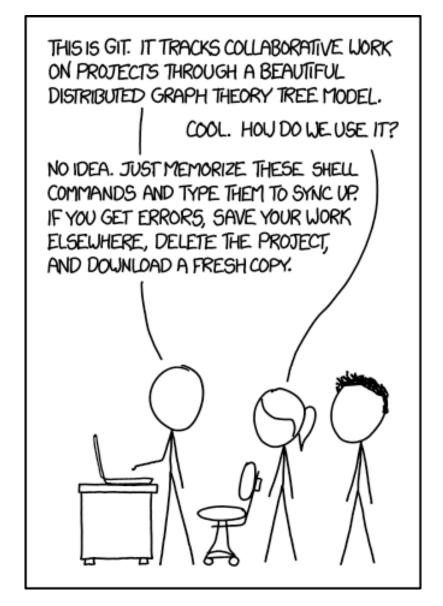


Figure 1: xkcd 1597

While Git admittedly has an ugly interface, its underlying design and ideas are beautiful. While an ugly interface has to be *memorized*, a beautiful design can be *understood*. For this reason, we give a bottom-up explanation of Git, starting with its data model and later covering the command-line interface. Once the data model is understood, the commands can be better understood, in terms of how they manipulate the underlying data model.

# Git's data model

Git has a well thought-out model that enables all the nice features of version control, like maintaining history, supporting branches, and enabling collaboration.

# **Snapshots**

Git models the history of a collection of files and folders within some top-level directory as a series of snapshots. In Git terminology, a file is called a "blob", and it's just a bunch of bytes. A directory is called a "tree", and it maps names to blobs or trees (so directories can contain other directories). A snapshot is the top-level tree that is being tracked. For example, we might have a tree as follows:

```
1 <root> (tree)
2 |
3 +- foo (tree)
4 | |
5 | + bar.txt (blob, contents = "hello world")
6 |
7 +- baz.txt (blob, contents = "git is wonderful")
```

The top-level tree contains two elements, a tree "foo" (that itself contains one element, a blob "bar.txt"), and a blob "baz.txt".

# **Modeling history: relating snapshots**

How should a version control system relate snapshots? One simple model would be to have a linear history. A history would be a list of snapshots in time-order. For many reasons, Git doesn't use a simple model like this.

In Git, a history is a *directed acyclic graph* (DAG) of snapshots. All this means is that each snapshot in Git refers to a set of "parents", the snapshots that preceded it. It's a set of parents rather than a single parent (as would be the case in a linear history) because a snapshot might descend from multiple parents, for example due to combining (merging) two parallel branches of development.

Git calls these snapshots "commit"s. Visualizing a commit history might look something like this:

In the ASCII art above, the os correspond to individual commits (snapshots). The arrows point to the parent of each commit (it's a "comes before" relation, not "comes after"). After the third commit, the

history branches into two separate branches. This might correspond to, for example, two separate features being developed in parallel, independently from each other. In the future, kkthese branches may be merged to create a new snapshot that incorporates both of the features, producing a new history that looks like this, with the newly created merge commit shown in bold:

Commits in Git are immutable. This doesn't mean that mistakes can't be corrected, however; it's just that "edits" to the commit history are actually creating entirely new commits, and references (see below) are updated to point to the new ones.

# Data model, as pseudocode

Git's data model written down in pseudocode:

```
1 // a file is a bunch of bytes
2 type blob = array<byte>
3
4 // a directory contains named files and directories
5 type tree = map<string, tree | file>
6
7 // a commit has parents, metadata, and the top-level tree
8 type commit = struct {
9    parent: array<commit>
10    author: string
11    message: string
12    snapshot: tree
13 }
```

It's a clean, simple model of history.

### **Objects and content-addressing**

An "object" is a blob, tree, or commit:

```
1 type object = blob | tree | commit
```

In Git data store, all objects are content-addressed by their SHA-1 hash.

```
1 objects = map<string, object>
2
3 def store(object):
4  id = shal(object)
```

```
5   objects[id] = object
6
7   def load(id):
8    return objects[id]
```

Blobs, trees, and commits are unified in this way: they are all objects. When they reference other objects, they don't actually *contain* them in their on-disk representation, but have a reference to them by their hash.

For example, the tree for the example directory structure above (visualized using git cat-file -p 698281bc680d1995c5f4caaf3359721a5a58d48d), looks like this:

```
1 100644 blob 4448adbf7ecd394f42ae135bbeed9676e894af85 baz.txt 2 040000 tree c68d233a33c5c06e0340e4c224f0afca87c8ce87 foo
```

The tree itself contains pointers to its contents, baz.txt (a blob) and foo (a tree). If we look at the contents addressed by the hash corresponding to baz.txt with git cat-file -p 4448 adbf7ecd394f42ae135bbeed9676e894af85, we get the following:

```
1 git is wonderful
```

#### References

Now, all snapshots can be identified by their SHA-1 hash. That's inconvenient, because humans aren't good at remembering strings of 40 hexadecimal characters.

Git's solution to this problem is human-readable names for SHA-1 hashes, called "references". References are pointers to commits. Unlike objects, which are immutable, references are mutable (can be updated to point to a new commit). For example, the master reference usually points to the latest commit in the main branch of development.

```
references = map<string, string>
2
3
  def update_reference(name, id):
4
       references[name] = id
5
  def read_reference(name):
7
       return references[name]
8
9 def load_reference(name_or_id):
       if name_or_id in references:
10
           return load(references[name_or_id])
12
13
           return load(name_or_id)
```

With this, Git can use human-readable names like "master" to refer to a particular snapshot in the history, instead of a long hexadecimal string.

One detail is that we often want a notion of "where we currently are" in the history, so that when we take a new snapshot, we know what it is relative to (how we set the parents field of the commit). In Git, that "where we currently are" is a special reference called "HEAD".

### **Repositories**

Finally, we can define what (roughly) is a Git repository: it is the data objects and references.

On disk, all Git stores is objects and references: that's all there is to Git's data model. All git commands map to some manipulation of the commit DAG by adding objects and adding/updating references.

Whenever you're typing in any command, think about what manipulation the command is making to the underlying graph data structure. Conversely, if you're trying to make a particular kind of change to the commit DAG, e.g. "discard uncommitted changes and make the 'master' ref point to commit 5d83f9e", there's probably a command to do it (e.g. in this case, git checkout master; git reset --hard 5d83f9e).

# **Staging area**

This is another concept that's orthogonal to the data model, but it's a part of the interface to create commits.

One way you might imagine implementing snapshotting as described above is have a "create snapshot" command that creates a new snapshot based on the *current state* of the working directory. Some version control tools work like this, but not Git. We want clean snapshots, and it might not always be ideal to make a snapshot from the current state. For example, imagine a scenario where you've implemented two separate features, and you want to create two separate commits, where the first introduces the first feature, and the next introduces the second feature. Or imagine a scenario where you have debugging print statements added all over your code, along with a bugfix; you want to commit the bugfix while discarding all the print statements.

Git accommodates such scenarios by allowing you to specify which modifications should be included in the next snapshot through a mechanism called the "staging area".

# Git command-line interface

To avoid duplicating information, we're not going to explain the commands below in detail. See the highly recommended Pro Git for more information, or watch the lecture video.

#### **Basics**

The git init command initializes a new Git repository, with repository metadata being stored in the .git directory:

```
1 $ mkdir myproject
2 $ cd myproject
3 $ git init
4 Initialized empty Git repository in /home/missing-semester/myproject/.
        git/
5 $ git status
6 On branch master
7
8 No commits yet
9
10 nothing to commit (create/copy files and use "git add" to track)
```

How do we interpret this output? "No commits yet" basically means our version history is empty. Let's fix that.

```
1 $ echo "hello, git" > hello.txt
2 $ git add hello.txt
3 $ git status
4 On branch master
6 No commits yet
8 Changes to be committed:
9
   (use "git rm --cached <file>..." to unstage)
           new file: hello.txt
11
12
13 $ git commit -m 'Initial commit'
14 [master (root-commit) 4515d17] Initial commit
15
   1 file changed, 1 insertion(+)
   create mode 100644 hello.txt
```

With this, we've git added a file to the staging area, and then git committed that change, adding a simple commit message "Initial commit". If we didn't specify a -m option, Git would open our text editor to allow us type a commit message.

Now that we have a non-empty version history, we can visualize the history. Visualizing the history as

a DAG can be especially helpful in understanding the current status of the repo and connecting it with your understanding of the Git data model.

The git log command visualizes history. By default, it shows a flattened version, which hides the graph structure. If you use a command like git log --all --graph --decorate, it will show you the full version history of the repository, visualized in graph form.

```
1 $ git log --all --graph --decorate
2 * commit 4515d17a167bdef0a91ee7d50d75b12c9c2652aa (HEAD -> master)
3 Author: Missing Semester <missing-semester@mit.edu>
4 Date: Tue Jan 21 22:18:36 2020 -0500
5
6 Initial commit
```

This doesn't look all that graph-like, because it only contains a single node. Let's make some more changes, author a new commit, and visualize the history once more.

```
1 $ echo "another line" >> hello.txt
2 $ git status
3 On branch master
4 Changes not staged for commit:
    (use "git add <file>..." to update what will be committed)
5
     (use "git checkout -- <file>..." to discard changes in working
6
        directory)
7
8
           modified:
                       hello.txt
10 no changes added to commit (use "git add" and/or "git commit -a")
11 $ git add hello.txt
12 $ git status
13 On branch master
14 Changes to be committed:
     (use "git reset HEAD <file>..." to unstage)
16
           modified:
                     hello.txt
18
19 $ git commit -m 'Add a line'
20 [master 35f60a8] Add a line
21
    1 file changed, 1 insertion(+)
```

Now, if we visualize the history again, we'll see some of the graph structure:

```
1 * commit 35f60a825be0106036dd2fbc7657598eb7b04c67 (HEAD -> master)
2 | Author: Missing Semester <missing-semester@mit.edu>
3 | Date: Tue Jan 21 22:26:20 2020 -0500
4 |
5 | Add a line
6 |
7 * commit 4515d17a167bdef0a91ee7d50d75b12c9c2652aa
Author: Anish Athalye <me@anishathalye.com>
```

```
9 Date: Tue Jan 21 22:18:36 2020 -0500
10
11 Initial commit
```

Also, note that it shows the current HEAD, along with the current branch (master).

We can look at old versions using the git checkout command.

```
1 $ git checkout 4515d17 # previous commit hash; yours will be different
2 Note: checking out '4515d17'.
   You are in 'detached HEAD' state. You can look around, make
      experimental
  changes and commit them, and you can discard any commits you make in
6 state without impacting any branches by performing another checkout.
8 If you want to create a new branch to retain commits you create, you
  do so (now or later) by using -b with the checkout command again.
      Example:
   git checkout -b <new-branch-name>
11
12
13 HEAD is now at 4515d17 Initial commit
14 $ cat hello.txt
15 hello, git
16 $ git checkout master
17 Previous HEAD position was 4515d17 Initial commit
18 Switched to branch 'master'
19 $ cat hello.txt
20 hello, git
21 another line
```

Git can show you how files have evolved (differences, or diffs) using the git diff command:

```
1 $ git diff 4515d17 hello.txt
2 diff --git c/hello.txt w/hello.txt
3 index 94bab17..f0013b2 100644
4 --- c/hello.txt
5 +++ w/hello.txt
6 @@ -1 +1,2 @@
7 hello, git
8 +another line
```

- git help <command>: get help for a git command
- git init: creates a new git repo, with data stored in the .git directory
- git status: tells you what's going on
- git add <filename>: adds files to staging area

- git commit: creates a new commit
  - Write good commit messages!
- git log: shows a flattened log of history
- git log --all --graph --decorate: visualizes history as a DAG
- git diff <filename>: show differences since the last commit
- git diff <revision> <filename>: shows differences in a file between snapshots
- git checkout <revision>: updates HEAD and current branch

# **Branching and merging**

Branching allows you to "fork" version history. It can be helpful for working on independent features or bug fixes in parallel. The git branch command can be used to create new branches; git checkout -b <br/> tranch name> creates and branch and checks it out.

Merging is the opposite of branching: it allows you to combine forked version histories, e.g. merging a feature branch back into master. The git merge command is used for merging.

- git branch: shows branches
- git branch <name>: creates a branch
- git checkout -b <name>: creates a branch and switches to it
  - same as git branch <name>; git checkout <name>
- git merge <revision>: merges into current branch
- git mergetool: use a fancy tool to help resolve merge conflicts
- git rebase: rebase set of patches onto a new base

#### Remotes

- git remote: list remotes
- git remote add <name> <url>: add a remote</pr>
- git push <remote> <local branch>:<remote branch>: send objects to remote, and update remote reference
- git branch --set-upstream-to=<remote>/<remote branch>: set up correspondence between local and remote branch
- git fetch: retrieve objects/references from a remote
- git pull: same as git fetch; git merge
- git clone: download repository from remote

#### Undo

- git commit --amend: edit a commit's contents/message
- git reset HEAD <file>: unstage a file
- git checkout -- <file>: discard changes
- git reset: moves branch reference backwards in time to an older commit
  - moves a branch backwards in time as if the commit had never taken place
  - git reset HEAD~1: moves HEAD back by one commit
- git revert: for reversing changes already propagated to the remote. Adds a new commit of changes

#### **Relative Refs**

Relative refs are useful to refer to commits without having to type out the hash - ^: move upwards by one commit - ^<num>: move upwards to parent reference e.g. for a merge commit with multiple parents (default: 1) - ~<num>: move upwards by num commits - e.g. - git checkout master ^: checks out first parent of master - master ^ ^: grandparent of master - HEAD~4: 4th ancestor of HEAD - chaining e.g. git checkout HEAD~^2~2: move to parent commit, move to second co-parent of commit, move two commits up - Applications - Branch forcing: directly reassign a branch to a commit - git branch - f master HEAD~3 forces master branch to 3 parents behind HEAD

### Moving work around

- git cherry-pick <commit1> <commit2> <...>: copies a series of commits below current location, i.e. add new commit corresponding to commit1, then commit2, etc.
- git rebase -i interactive rebase. Opens a text editor showing which commits are about to be copied below the rebase target. You can then
  - reorder commits by changing order
  - omit commits by setting pick off
  - squash commits (i.e. combine multiple commits)

### Finding your way around

- git tag <name> to apply name to a specific commit e.g. version number
- git describe describes where you are relative to nearest anchor (tagged commit)

- output: <tag>\_<numCommits>\_g<hash>
  - \* tag: tag of nearest anchor
  - \* numCommits: number of commits away to anchor
  - \* hash: hash of commit being described (not the anchor)

# **Advanced Git**

- git config: Git is highly customizable
- git clone --shallow: clone without entire version history
- git add -p: interactive staging
- git rebase -i: interactive rebasing
- git blame: show who last edited which line
- git stash: temporarily remove modifications to working directory
- git bisect: binary search history (e.g. for regressions)
- .gitignore: specify intentionally untracked files to ignore

# Miscellaneous

- **GUIs**: There are many GUI clients out there for Git. We personally don't use them and use the command-line interface instead.
- **Shell integration**: It's super handy to have a Git status as part of your shell prompt (zsh, ). Often included in frameworks like Oh My Zsh.
- **Editor integration**: Similarly to the above, handy integrations with many features. fugitive.vim is the standard one for Vim.
- **Workflows**: we taught you the data model, plus some basic commands; we didn't tell you what practices to follow when working on big projects (and there are many different approaches).
- **GitHub**: Git is not GitHub. GitHub has a specific way of contributing code to other projects, called pull requests.
- Other Git providers: GitHub is not special: there are many Git repository hosts, like GitLab and BitBucket.

### Resources

• Pro Git is **highly recommended reading**. Going through Chapters 1–5 should teach you most of what you need to use Git proficiently, now that you understand the data model. The later

- chapters have some interesting, advanced material.
- Oh Shit, Git!?! is a short guide on how to recover from some common Git mistakes.
- Git for Computer Scientists is a short explanation of Git's data model, with less pseudocode and more fancy diagrams than these lecture notes.
- Git from the Bottom Up is a detailed explanation of Git's implementation details beyond just the data model, for the curious.
- How to explain git in simple words
- Learn Git Branching is a browser-based game that teaches you Git.

# **Exercises**

- 1. If you don't have any past experience with Git, either try reading the first couple chapters of Pro Git or go through a tutorial like Learn Git Branching. As you're working through it, relate Git commands to the data model.
- 2. Clone the repository for the class website. git clone https://github.com/missing-semester/missing-semester
  - 1. Explore the version history by visualizing it as a graph. git log --graph
  - 2. Who was the last person to modify README.md? (Hint: use git log with an argument) git log -n 1 README.md
  - 3. What was the commit message associated with the last modification to the collections : line of \_config.yml? (Hint: use git blame and git show)

```
1 git blame -L /collections:/ _config.yml
2 git show a88b4eac
```

- 3. One common mistake when learning Git is to commit large files that should not be managed by Git or adding sensitive information. Try adding a file to a repository, making some commits and then deleting that file from history (you may want to look at this).
- 4. Clone some repository from GitHub, and modify one of its existing files. What happens when you do git stash? What do you see when running git log --all --oneline? Run git stash pop to undo what you did with git stash. In what scenario might this be useful?
- 5. Like many command line tools, Git provides a configuration file (or dotfile) called ~/. gitconfig. Create an alias in ~/.gitconfig so that when you run git graph, you get the output of git log --all --graph --decorate --oneline.
- 6. You can define global ignore patterns in ~/.gitignore\_global after running git config --global core.excludesfile ~/.gitignore\_global. Do this, and set up your global gitignore file to ignore OS-specific or editor-specific temporary files, like .DS\_Store.

7. Clone the repository for the class website, find a typo or some other improvement you can make, and submit a pull request on GitHub.