

Making Analytics Reusable:

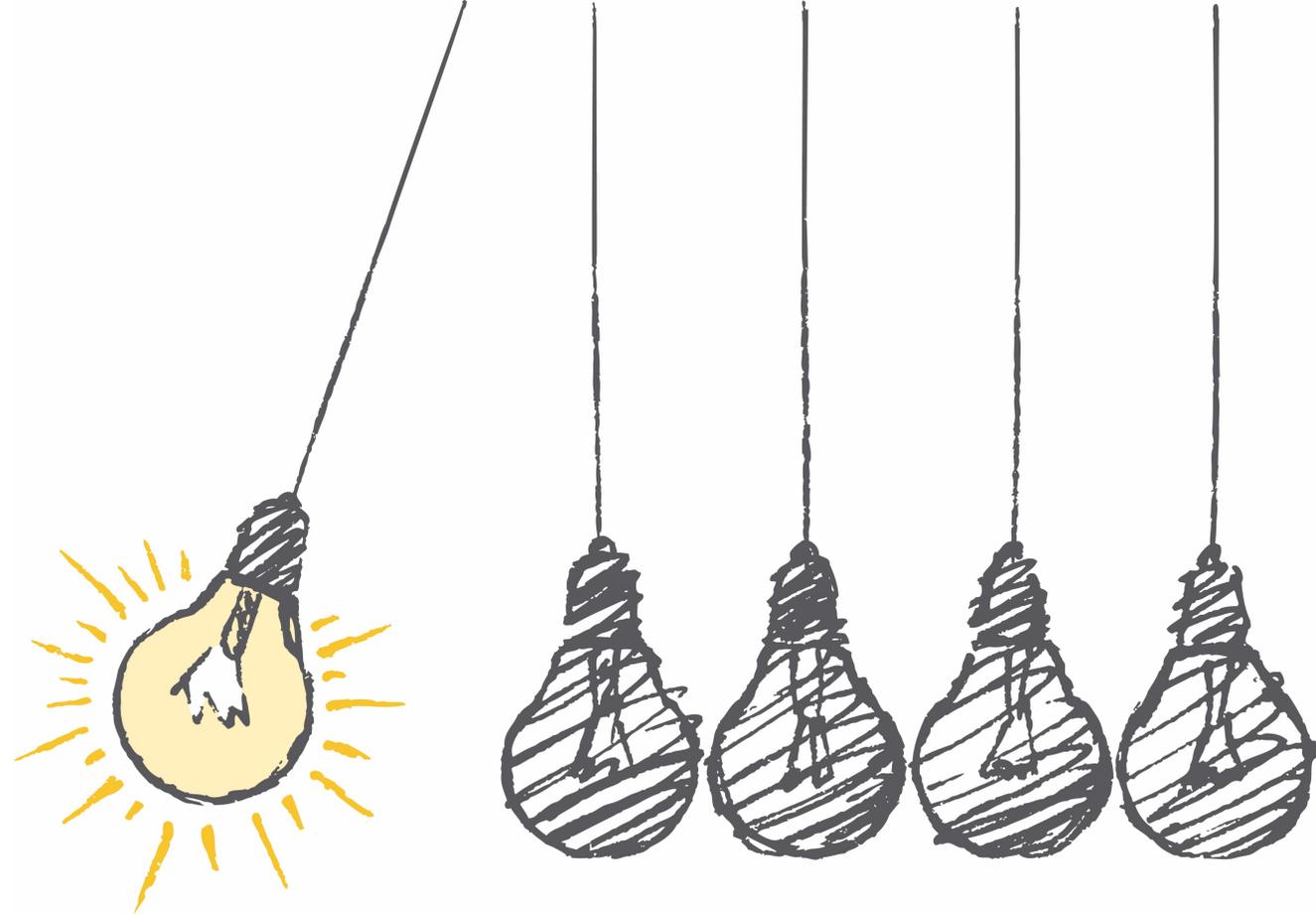
Using Git and GitHub

Prepared by **DIME Analytics**

dimeanalytics@worldbank.org

Presented by **Benjamin Daniels**

bdaniels@worldbank.org



Norad

What is Git?

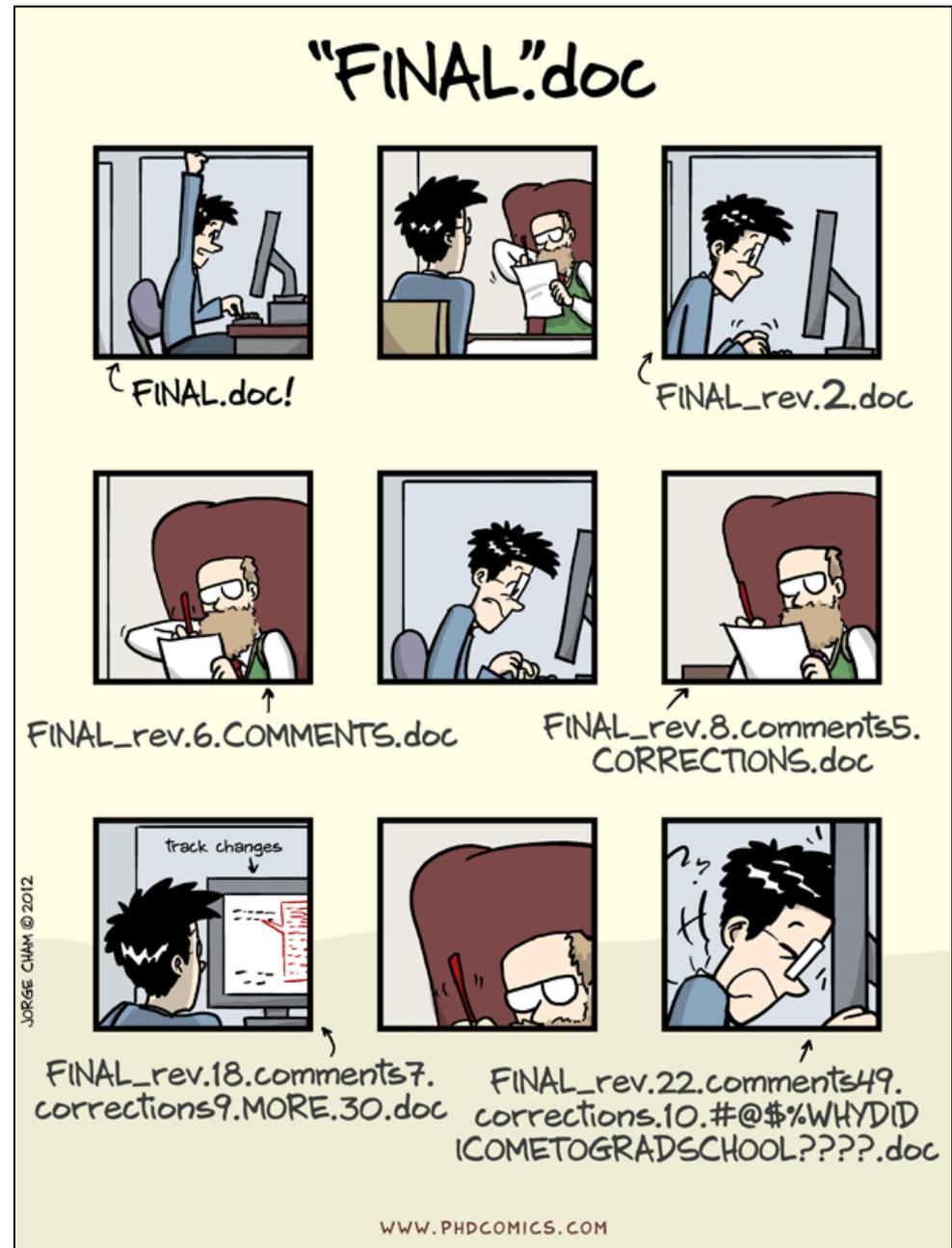
What is GitHub?

First, Git is not GitHub.
Second, GitHub is not Git.

Git is a **version control software**. It organizes *versions* of every file in a project in an orderly way.

GitHub is a “**development platform**”. It manages *contributions* to projects in an orderly way.

Together, Git and GitHub allow you to avoid the “FINAL”.doc problem.

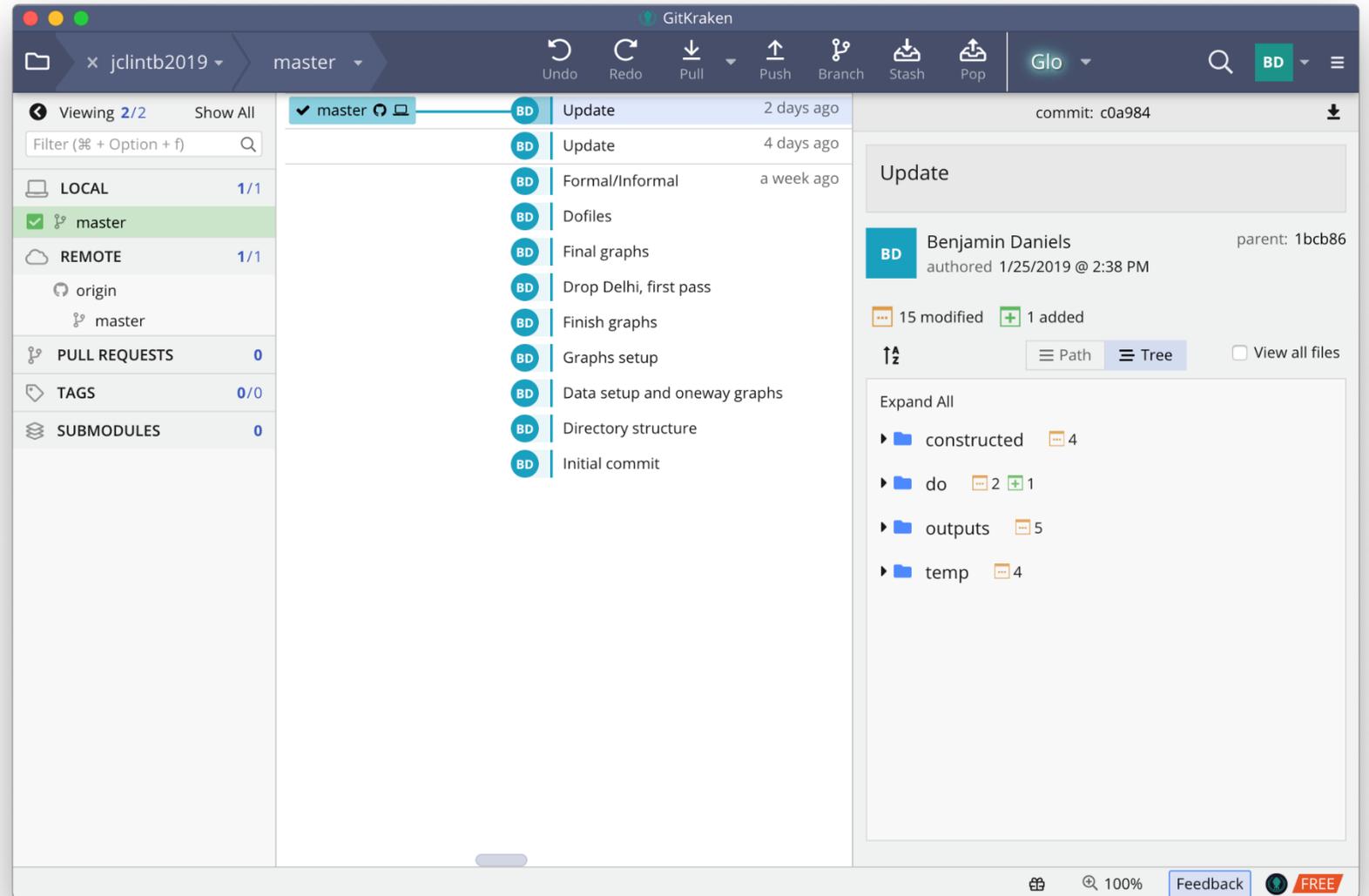


Git: The basics

Git is a software that runs on your computer. It stores “version histories” of your files in a way that is really useful, but that you mainly don’t need to understand.

The complete history of a project is called a “repository”. A repository can be viewed in a desktop “client”, such as [GitKraken](#), shown here.

A very simple repository would look much like this one: it has a series of versions called *commits* that make up its history. Each commit is a record of the changes between itself and the previous commit.

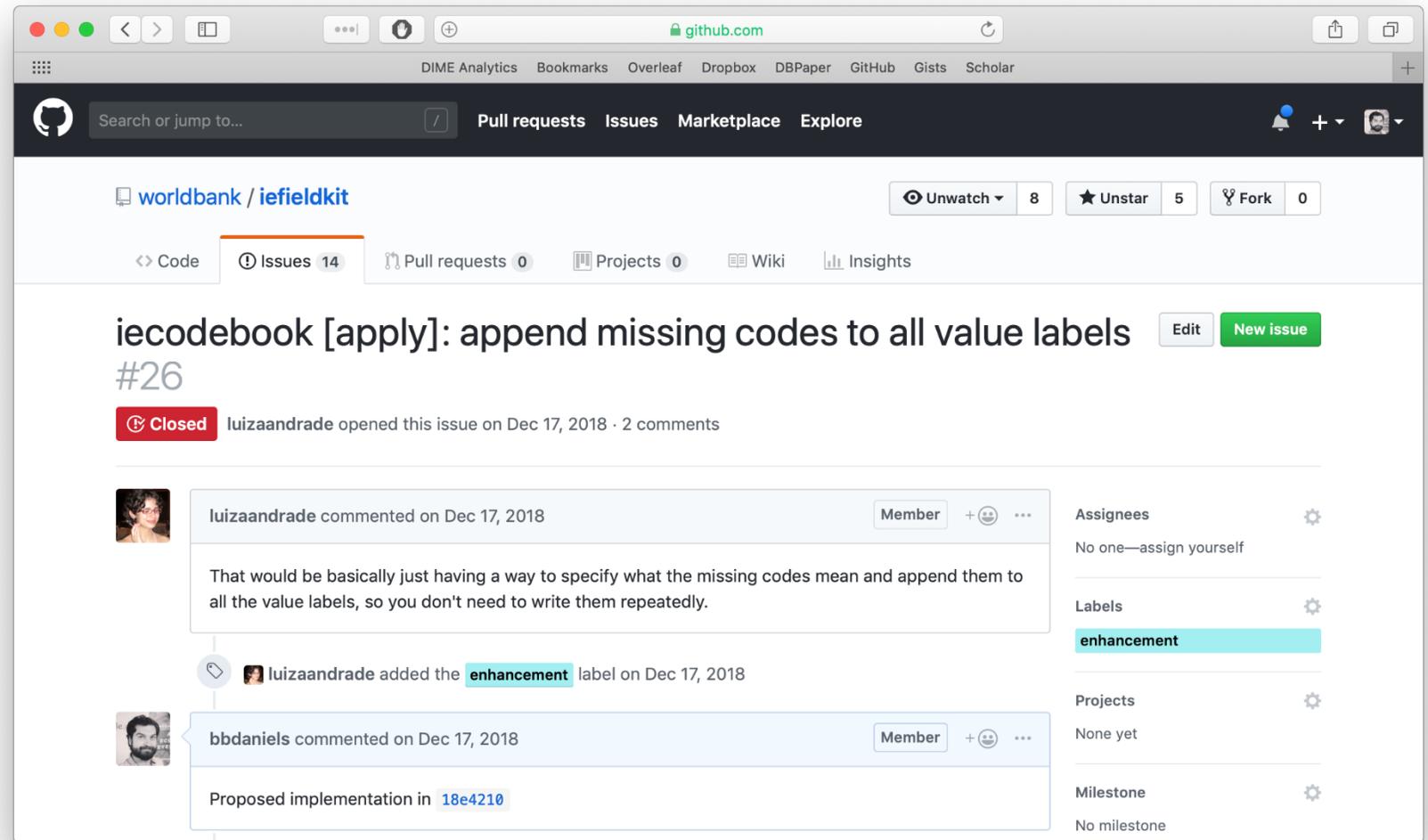


GitHub: The basics

GitHub is a website that you can use to manage collaboration with other people.

This has two main forms: first, managing and tracking tasks (known generally as “issues”), and second, allowing people to submit new materials to the repository.

A very simple issue looks like this one: A problem is posed, and after some discussion a solution is submitted as a “pull request” to the repository.



The screenshot shows a GitHub issue page for the repository 'worldbank / iefieldkit'. The issue title is 'iecodebook [apply]: append missing codes to all value labels' with a green 'New issue' button. The issue is marked as 'Closed' and was opened by 'luizaandrade' on Dec 17, 2018. The issue number is #26. The issue description is: 'That would be basically just having a way to specify what the missing codes mean and append them to all the value labels, so you don't need to write them repeatedly.' A comment by 'bbdaniels' on Dec 17, 2018, proposes an implementation in pull request '18e4210'. The right sidebar shows settings for Assignees (No one—assign yourself), Labels (enhancement), Projects (None yet), and Milestone (No milestone).

Git & GitHub: The basics

Mastering Git with GitHub gives you an awesome tool for managing your own team's workflow, as well as for making and releasing big projects.

Having both your data, code, and *code history* available for all authors and the public makes an incredibly useful tool as a product of your work.

Plus, Git protects all your files much better than other software.

(Shown here: The World Bank Atlas of Sustainable Development Goals.)

worldbank / sdgatl2018

9 commits 1 branch 0 releases 1 contributor

Branch: master New pull request Create new file Upload files Find file Clone or download

File	Commit Message
docs	First clean commit
inputs	Final tweak to data license info
.gitignore	First clean commit
INSTALL.md	Update INSTALL.md
README.md	Merge branch 'master' of github.com:worldbank/sdgatl2018
about_the_atlas.R	First clean commit
make.R	First clean commit
sdg1.R	First clean commit
sdg10.R	Add notes on data terms.
sdg11.R	First clean commit
sdg12.R	First clean commit
sdg13.R	First clean commit
sdg14.R	First clean commit
sdg15.R	Add notes on data terms.
sdg16.R	First clean commit
sdg17.R	First clean commit
sdg2.R	First clean commit
sdg3.R	Add notes on data terms.
sdg4.R	First clean commit
sdg5.R	First clean commit
sdg6.R	First clean commit

The foundation for any evidence is trust: trust that data have been collected, managed, and analyzed responsibly and trust that they have been faithfully presented. The *Atlas* is the first World Bank publication that sets out to be computationally **reproducible**—the majority of its charts and maps are produced with published code, directly from public data sources such as the World Bank's Open Data platform.

The *Atlas* distills the World Bank's knowledge of data related to the SDGs. I hope it inspires you to explore these issues further so that we can collectively accelerate progress toward achieving the SDGs.

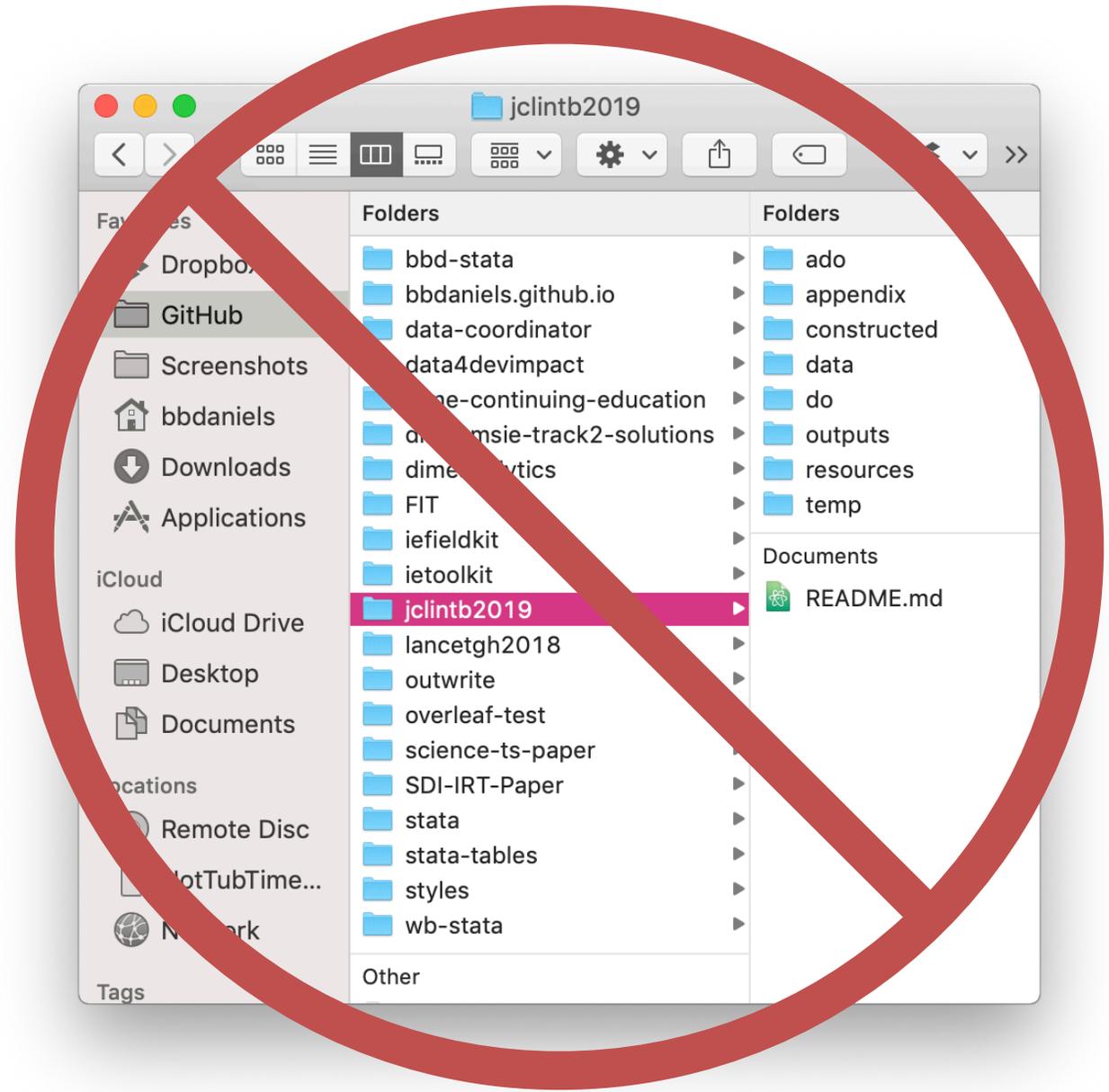
Shanta Devarajan
Senior Director, Development Economics and
Acting Chief Economist
World Bank Group

The files are not the repository

Your Git repositories will probably mainly be used to manage a lot of files. But it is important to remember that although the terms are often used semi-interchangeably, the “repository” is *not* the files themselves.

The files live in what we will call the “working directory” – the project folder (“directory”) where you are working.

A Git repository manages the contents of the working directory.



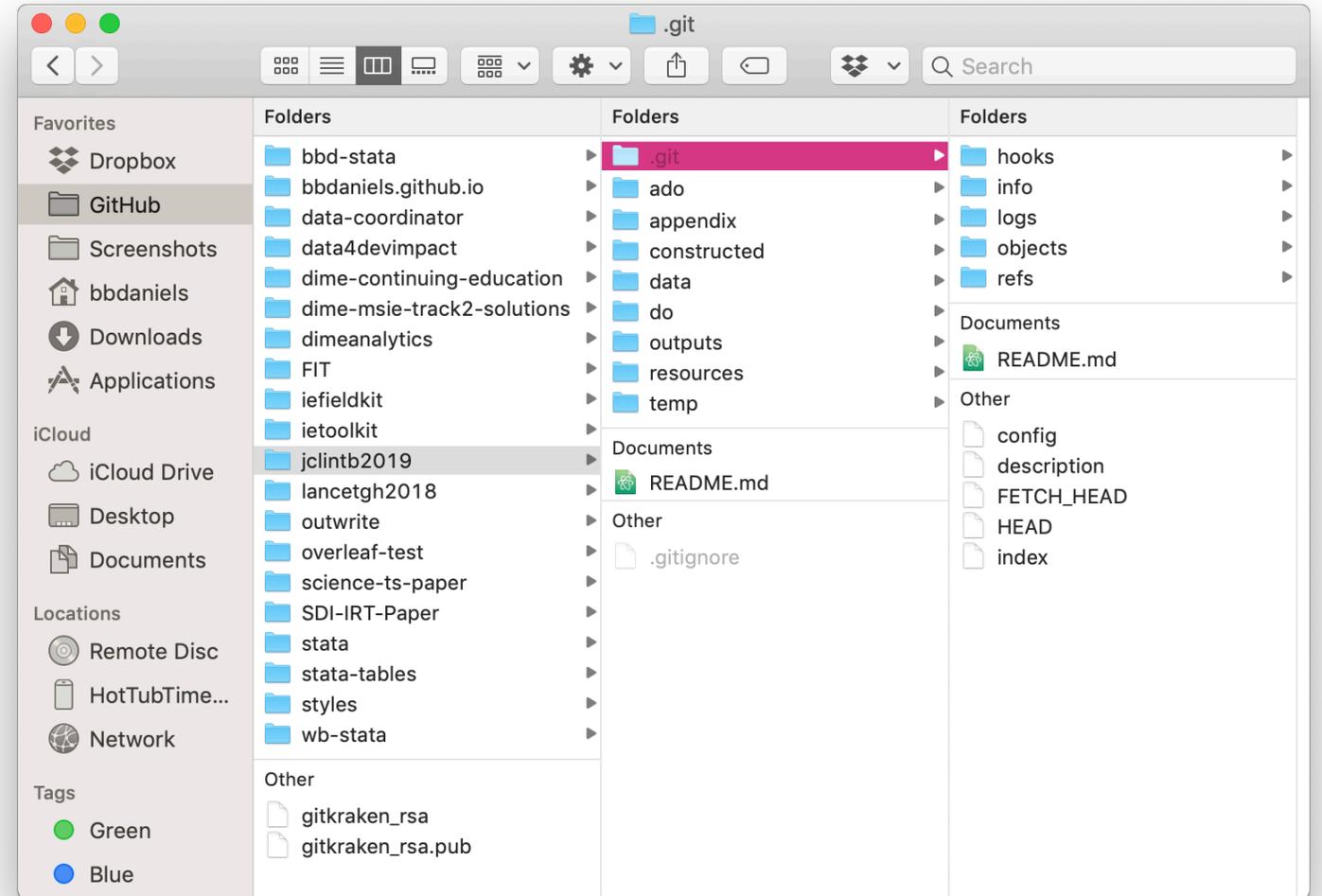
What is the repository then?

You *really* don't need to remember this, but:

The Git repository is a bunch of hidden files *inside* the working directory. This is where the magic happens.

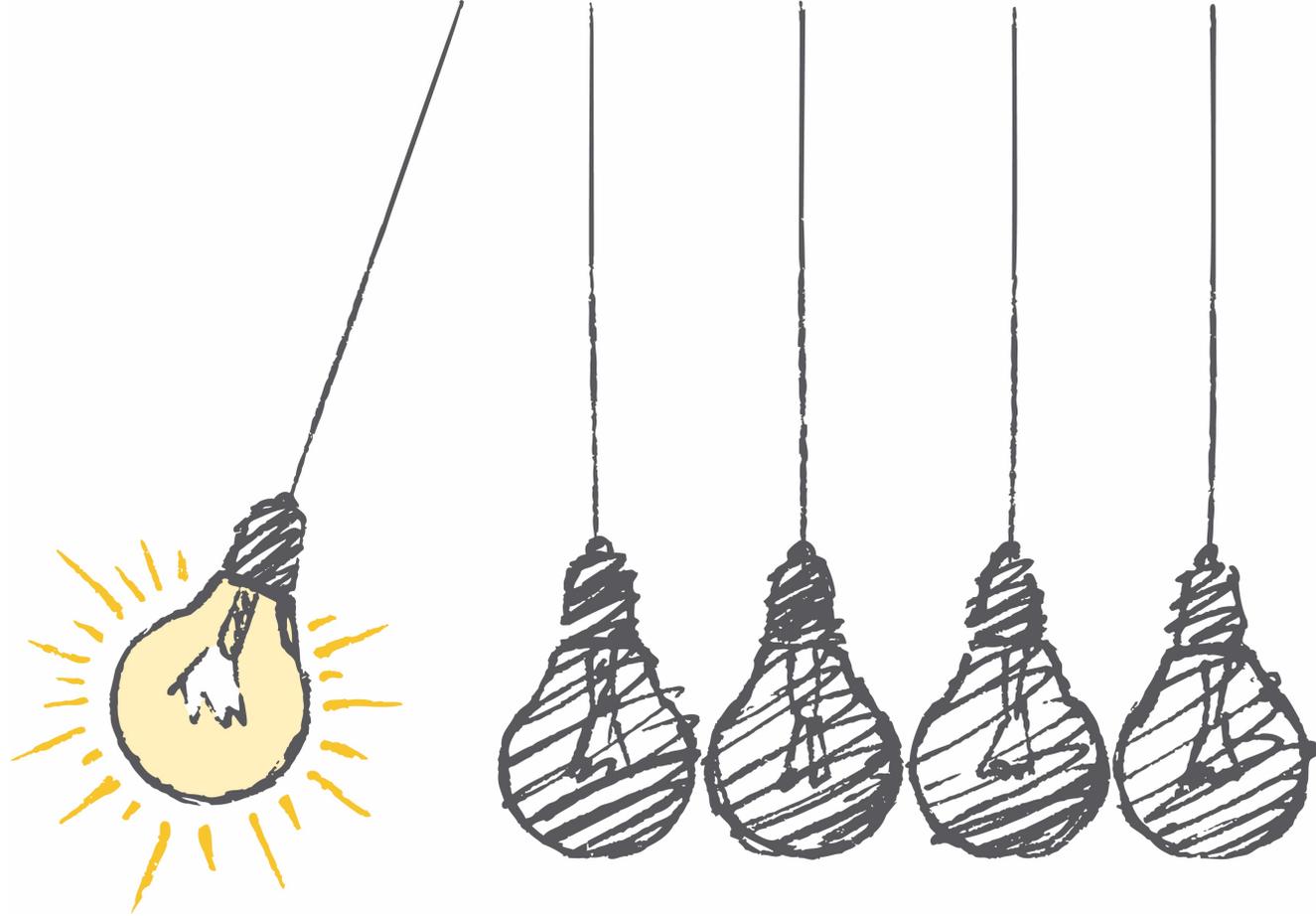
It's hidden because you never need to touch these. NEVER.

Just pretend this doesn't exist – everything important happens entirely inside your Git client.



So how does it *work*?

Part 1: Review commits and branches



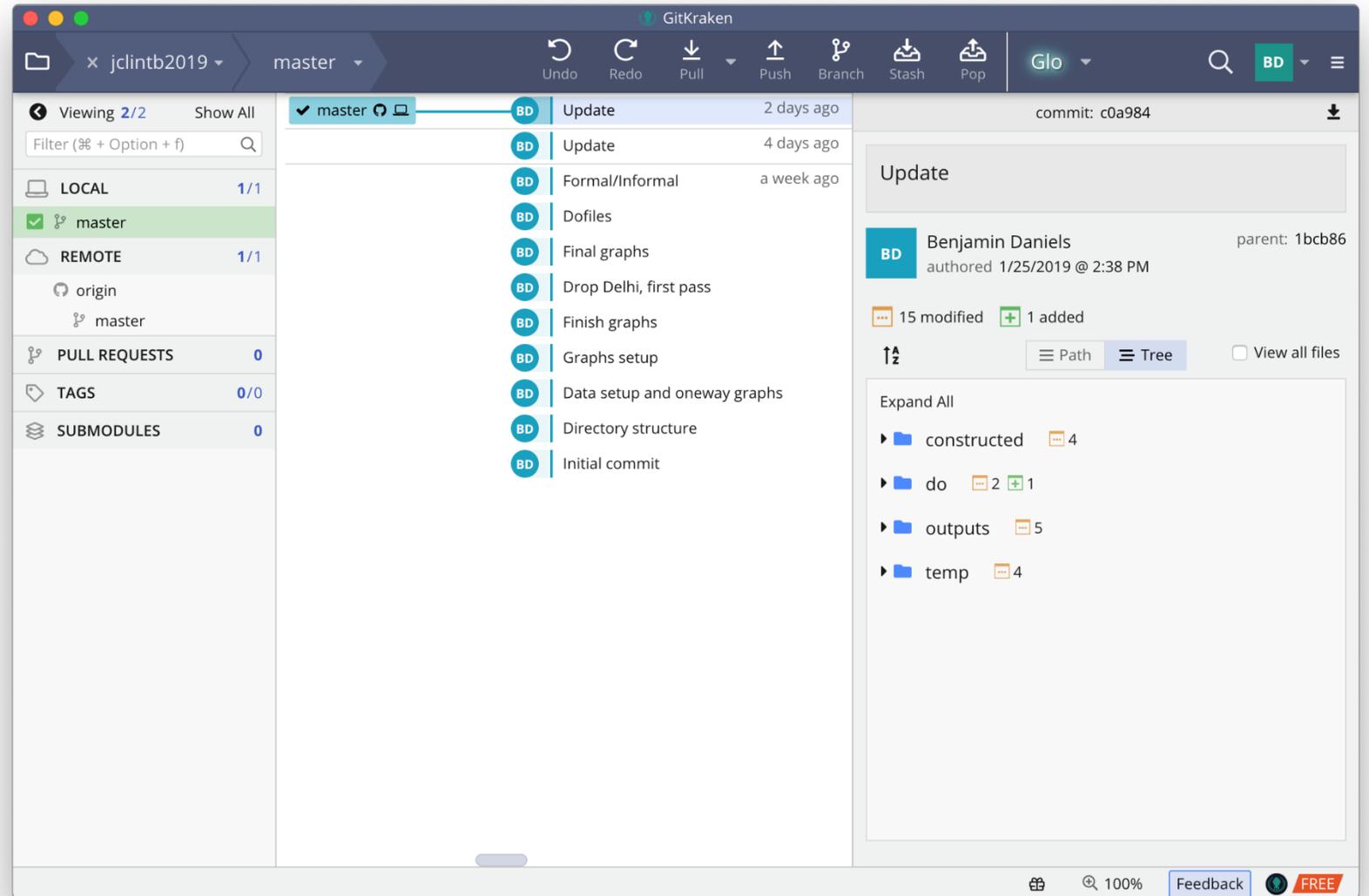
Norad

Git creates commits

The most important thing Git does is create “commits”.

A commit is a snapshot of the entire working directory. Each commit has a *name* and a *timestamp*.

You have to tell Git when to create new commits. This is probably the biggest difference from whatever you do now.



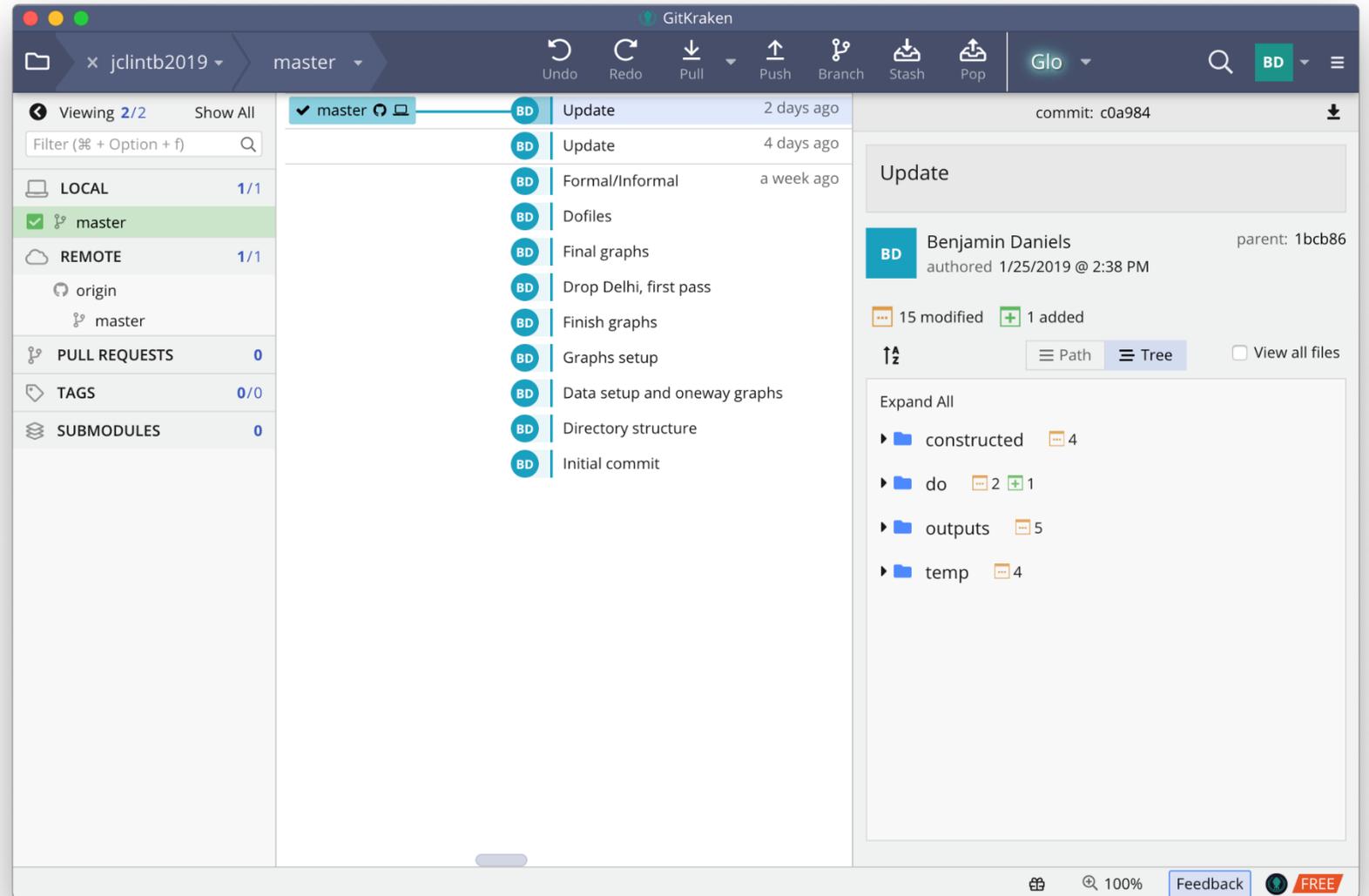
Commits name past versions

Git *watches* your saved changes

But it only *stores* them when you tell it to: this is a “commit”

You can name these so that they make sense and are easy to find when you need them in the future (you *always* need them when you least expect, right?)

Required Reading: Git vs Dropbox
<https://michaelstepner.com/blog/git-vs-dropbox>

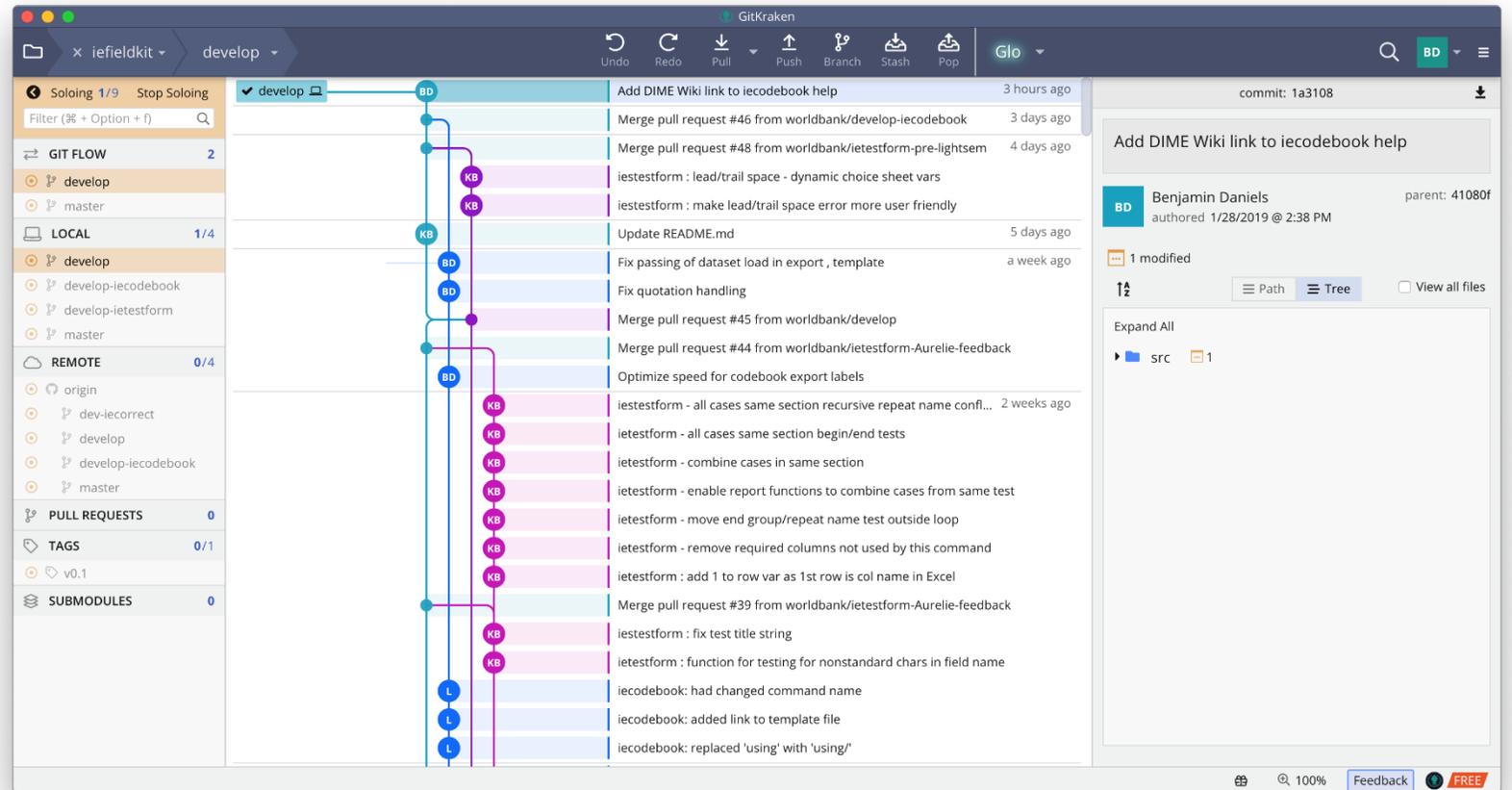


Git creates branches

Branches are the “killer feature” of Git. Nearly every advanced Git usage you learn will be about how to manage branches.

Branches enable different people to work on the same thing at the same time, and they enable you to view different versions of your files.

Branching allows you to move forwards or backwards in time; and to move “horizontally” in time through various concurrent versions.



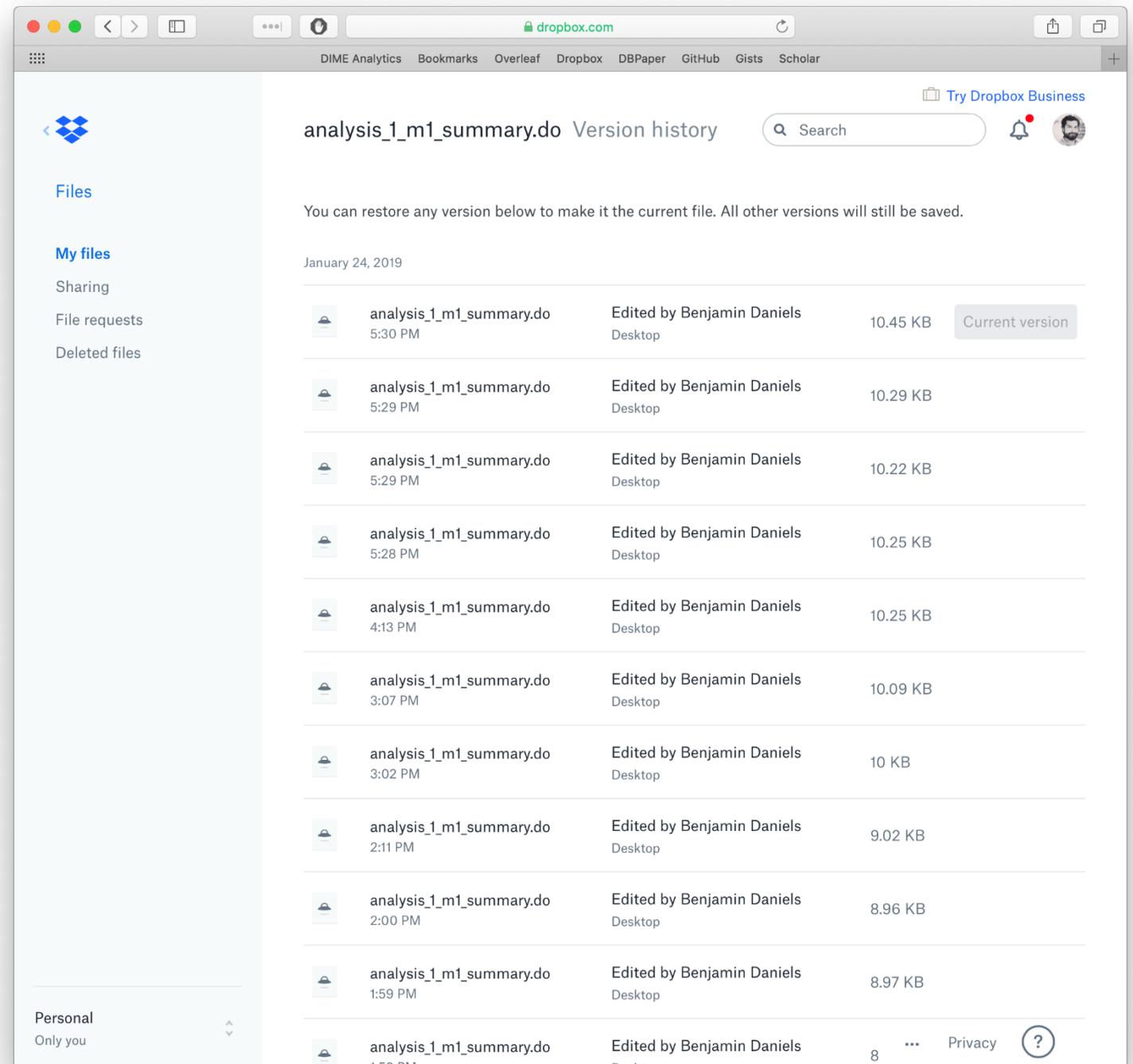
How is that better than Dropbox?

Dropbox stores a new snapshot *every time you save* each file.

You save your work *often* (right?)

So Dropbox stores *many* similar images of your files

But you can't tell which is which!



The screenshot shows a web browser window displaying the Dropbox version history for a file named 'analysis_1_m1_summary.do'. The browser's address bar shows 'dropbox.com'. The page title is 'analysis_1_m1_summary.do Version history'. A search bar and user profile icons are visible in the top right. The main content area shows a list of file versions, all edited by Benjamin Daniels on a Desktop. The most recent version, from January 24, 2019, at 5:30 PM, is marked as the 'Current version' and has a size of 10.45 KB. Other versions range from 10.22 KB to 10.09 KB, with timestamps from 1:58 PM to 5:29 PM. A 'Privacy' link and a help icon are visible at the bottom right of the list.

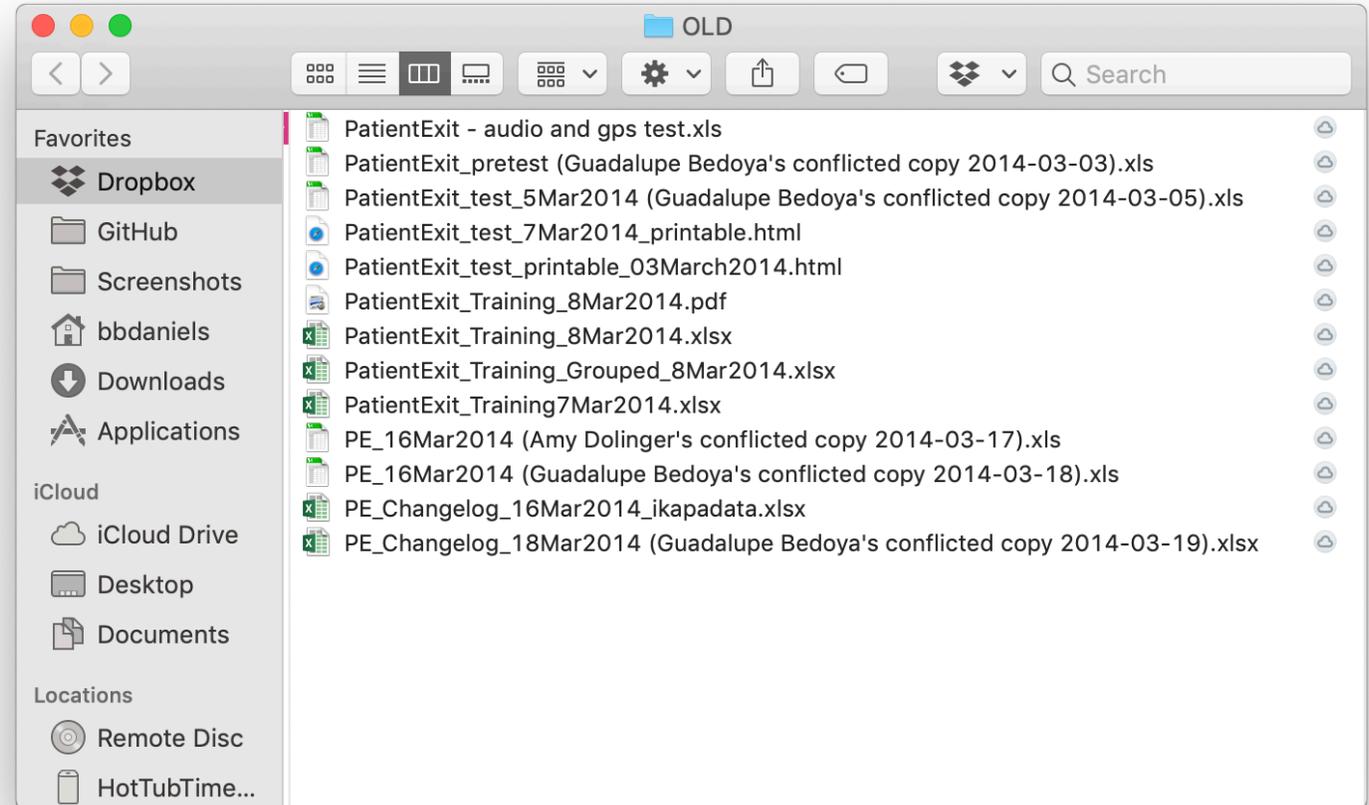
Version	Edited by	Location	Size	Notes
analysis_1_m1_summary.do 5:30 PM	Benjamin Daniels	Desktop	10.45 KB	Current version
analysis_1_m1_summary.do 5:29 PM	Benjamin Daniels	Desktop	10.29 KB	
analysis_1_m1_summary.do 5:29 PM	Benjamin Daniels	Desktop	10.22 KB	
analysis_1_m1_summary.do 5:28 PM	Benjamin Daniels	Desktop	10.25 KB	
analysis_1_m1_summary.do 4:13 PM	Benjamin Daniels	Desktop	10.25 KB	
analysis_1_m1_summary.do 3:07 PM	Benjamin Daniels	Desktop	10.09 KB	
analysis_1_m1_summary.do 3:02 PM	Benjamin Daniels	Desktop	10 KB	
analysis_1_m1_summary.do 2:11 PM	Benjamin Daniels	Desktop	9.02 KB	
analysis_1_m1_summary.do 2:00 PM	Benjamin Daniels	Desktop	8.96 KB	
analysis_1_m1_summary.do 1:59 PM	Benjamin Daniels	Desktop	8.97 KB	
analysis_1_m1_summary.do 1:58 PM	Benjamin Daniels	Desktop	8	Privacy ?

How is that better than Dropbox?

In Dropbox, there can be only *one* living version of a file at any time – otherwise a “conflicted copy” is created. This means nobody can edit the same file at the same time: Dropbox has no concurrency.

Worse, if you “roll back” a file to a previous version to see what it looked like, this affects everyone’s *current* version, even if you don’t want it to.

Finally, Dropbox versions are costly (computationally) – so they delete them often without telling you.



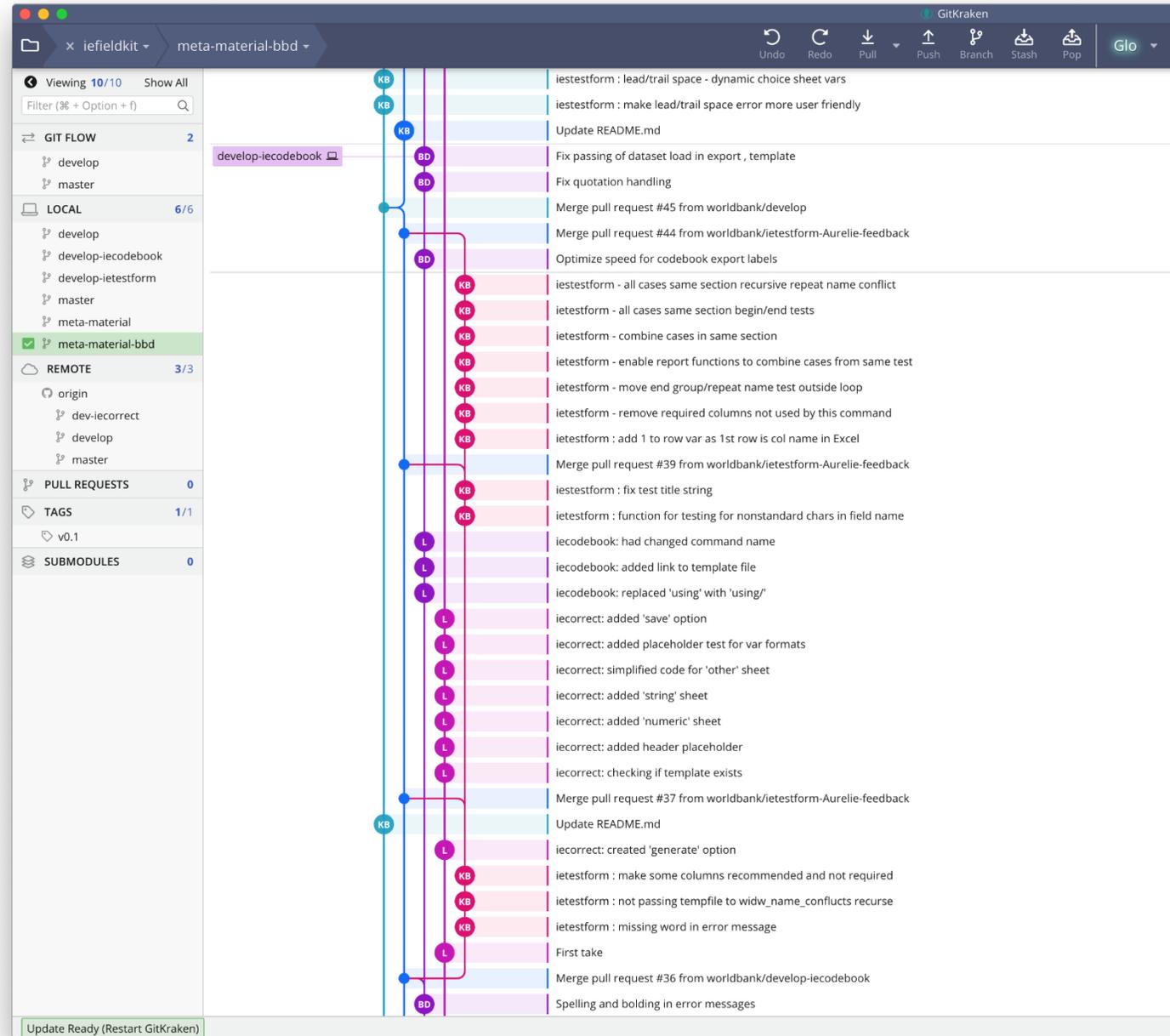
This can get complicated.

To make sure we know what we are doing, we need to *organize* our commits and branches.

This is non-trivial: it was more than 5 years between the release of Git and the development of a “successful” branching model.

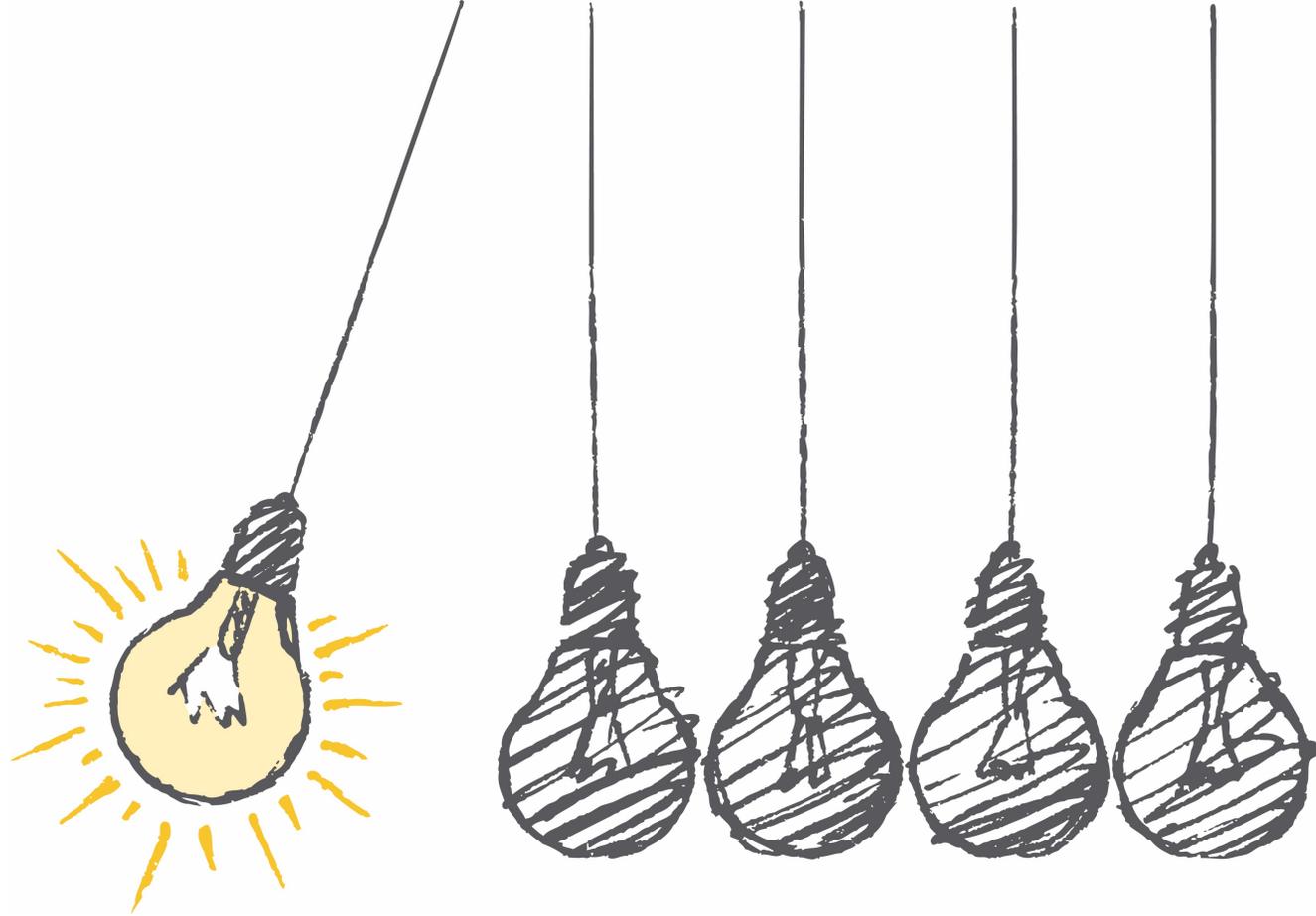
You have to think in ways you are not used to thinking when you start using Git, so these workflows will not feel intuitive at all.

But they are worth learning, since Git will not be as useful otherwise.



So how does it *work*?

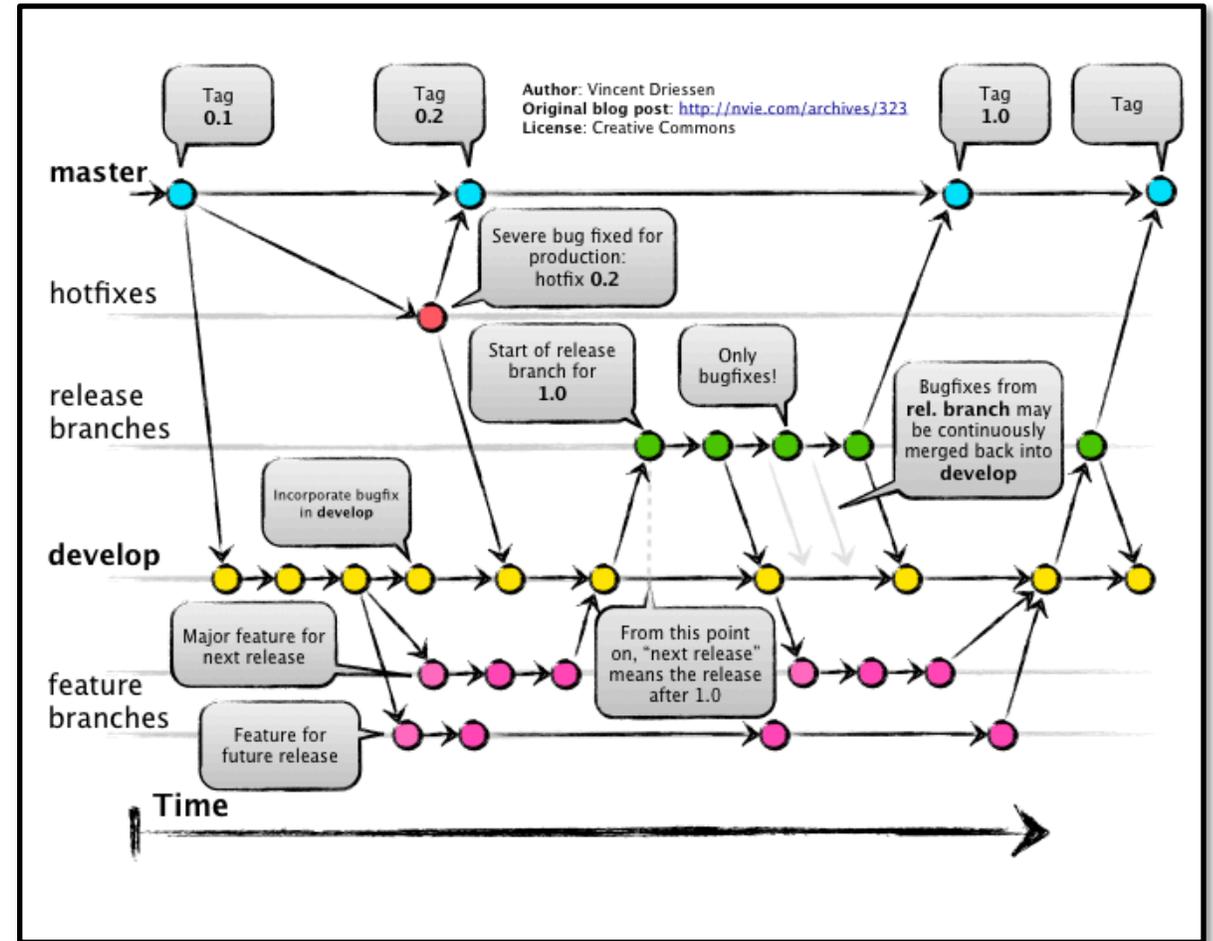
Part 2: Organize your commits and branches



Norad

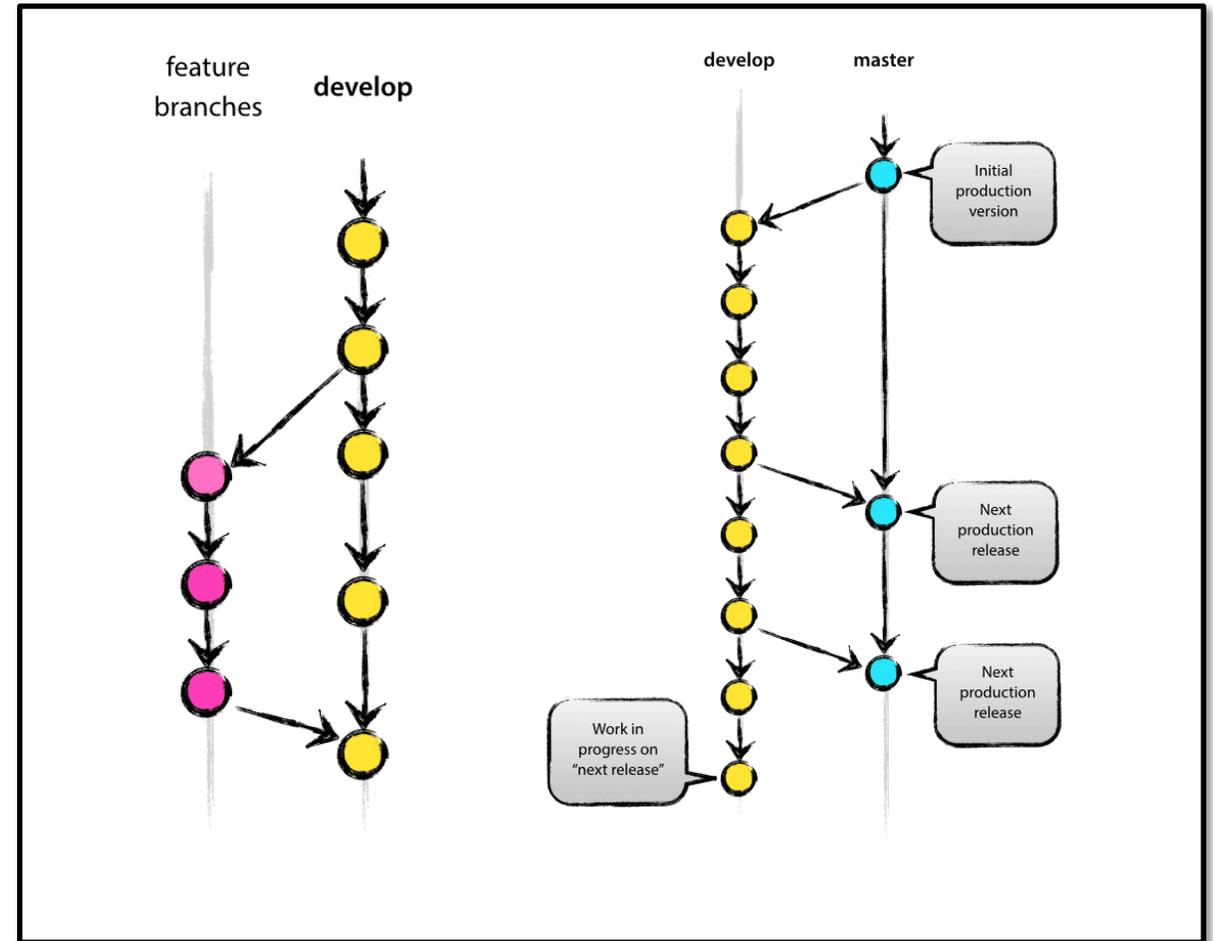
Git Flow organizes your commits and branches

- Git Flow is a similar idea to a standardized folder structure
- When you move to another project using Git Flow, you will know what everything is for and how you should be working
- Required Reading: Git Flow <https://nvie.com/posts/a-successful-git-branching-model/>



Core structure of Git Flow workflow

1. [master] branch is protected, always ready to run, and has sparse commits from [develop]
2. [develop] branch is protected, “nearly” ready to run, and has frequent commits from [feature] branches
3. [feature] branches are frequent, specific, personal, and hold all new work



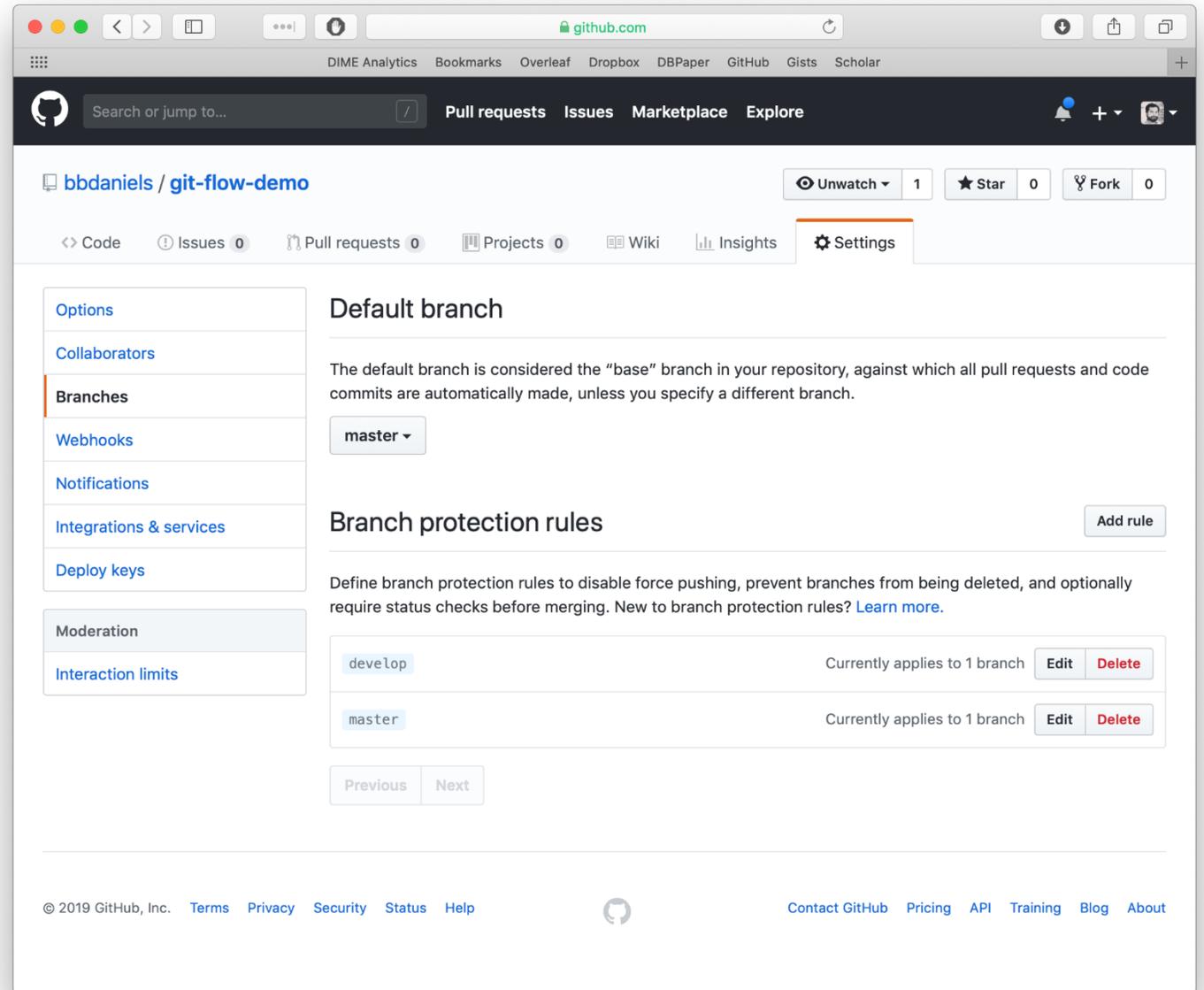
Two “protected” branches

These have rules set on GitHub which prevent:

- “Force pushing” (this can be used to overwrite history and is not allowed)
- Merging without approval

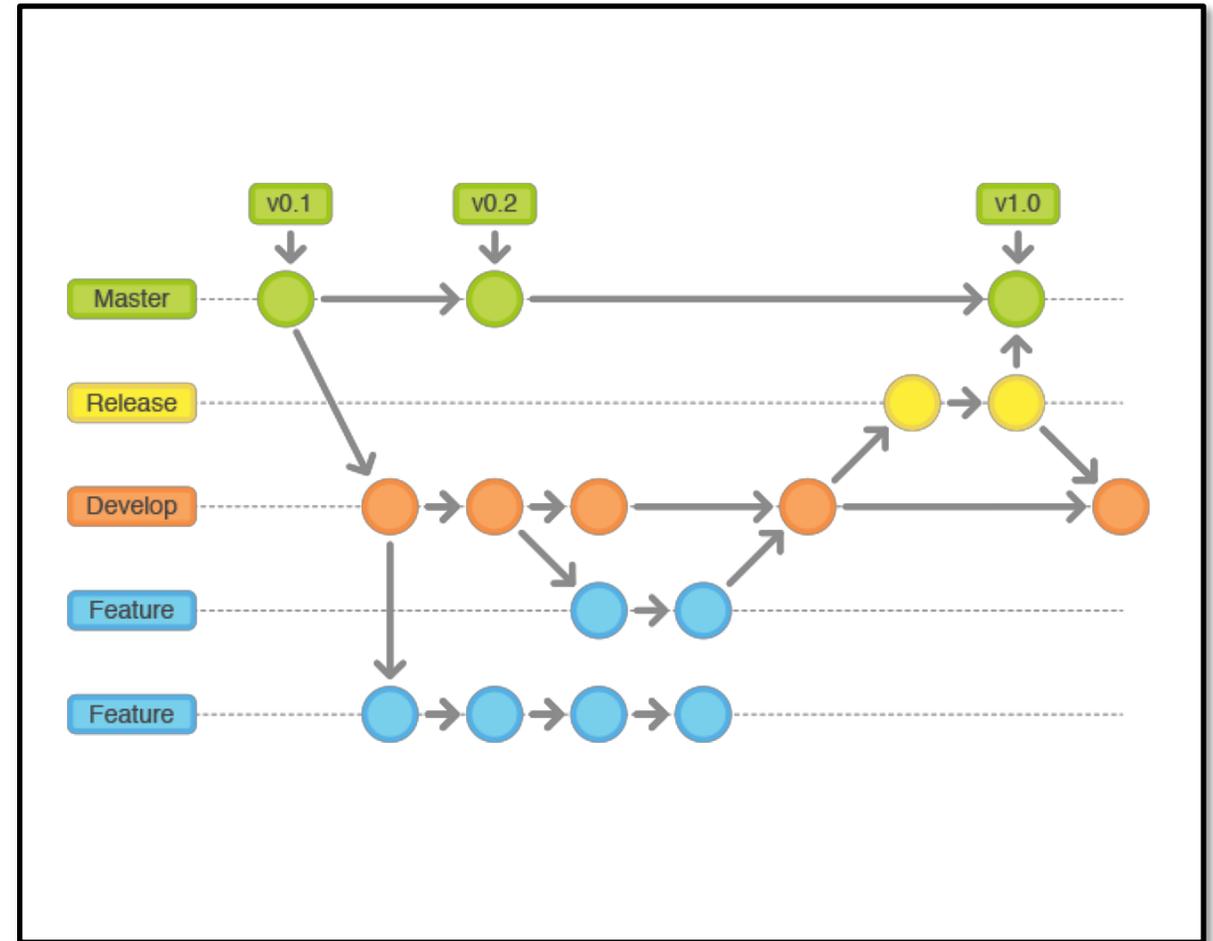
With these protections in place, the [master] branch is a preserved common status record that nobody edits directly.

[develop] is the common workspace where ready changes are merged after being built in a [feature] branch.



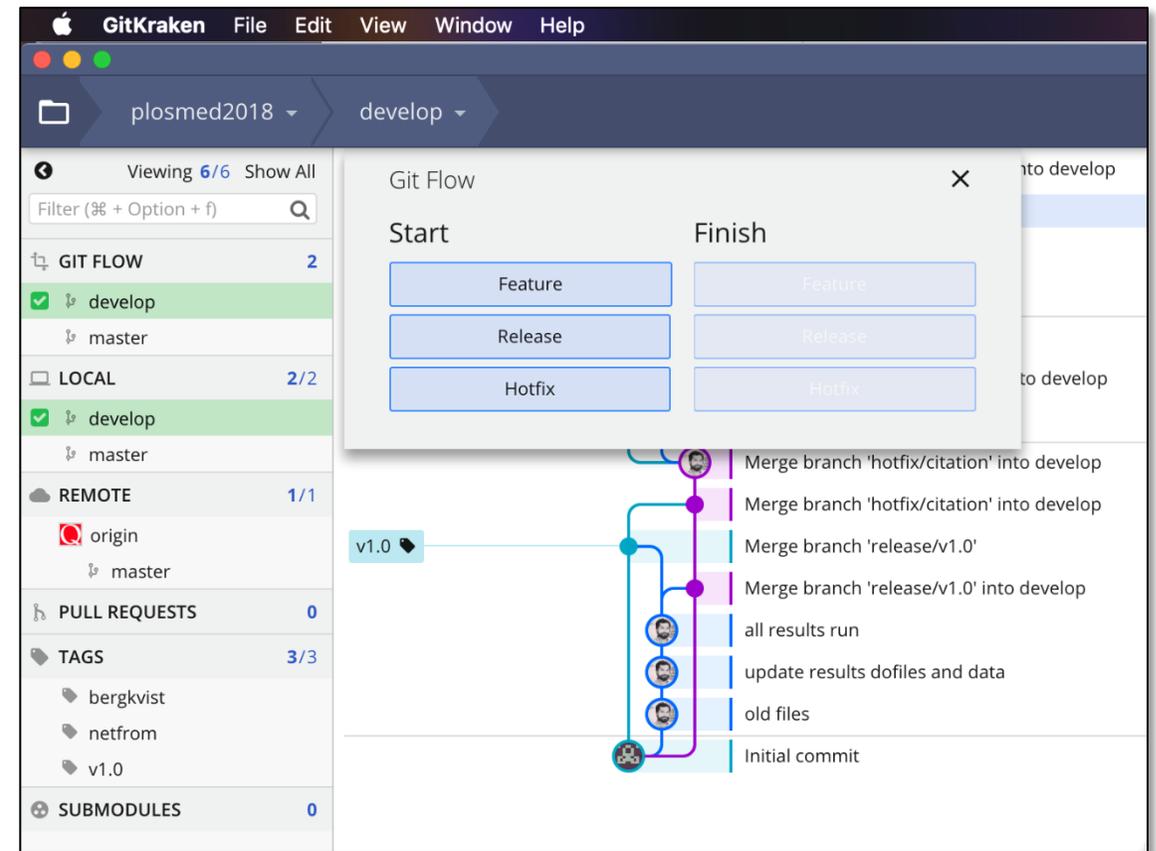
Benefits of Git Flow

- Anyone (like the PI) can *always* open [master] branch and see “where we are”.
- Anyone (like the PI) can create a [feature] branch and ask “what happens if I do X?”
- You can always look back and see what happened when someone tried something else (abandoned [feature] branches)



Using Git Flow with GitKraken

- Git Flow is so successful it is the default in many organizations, and GitKraken has powerful automatic features for supporting Git Flow.
- Send me your GitHub username so I can give you edit access on the git-flow-demo repo



Clone the [git-flow-demo](#) repo locally

Repository Management

Open

Clone

Init

- Clone with URL
- GitHub.com
- GitHub Enterprise
- GitLab.com
- GitLab (self-hosted)
- Bitbucket.org
- Visual Studio Team Services

Clone a Repo

Where to clone to

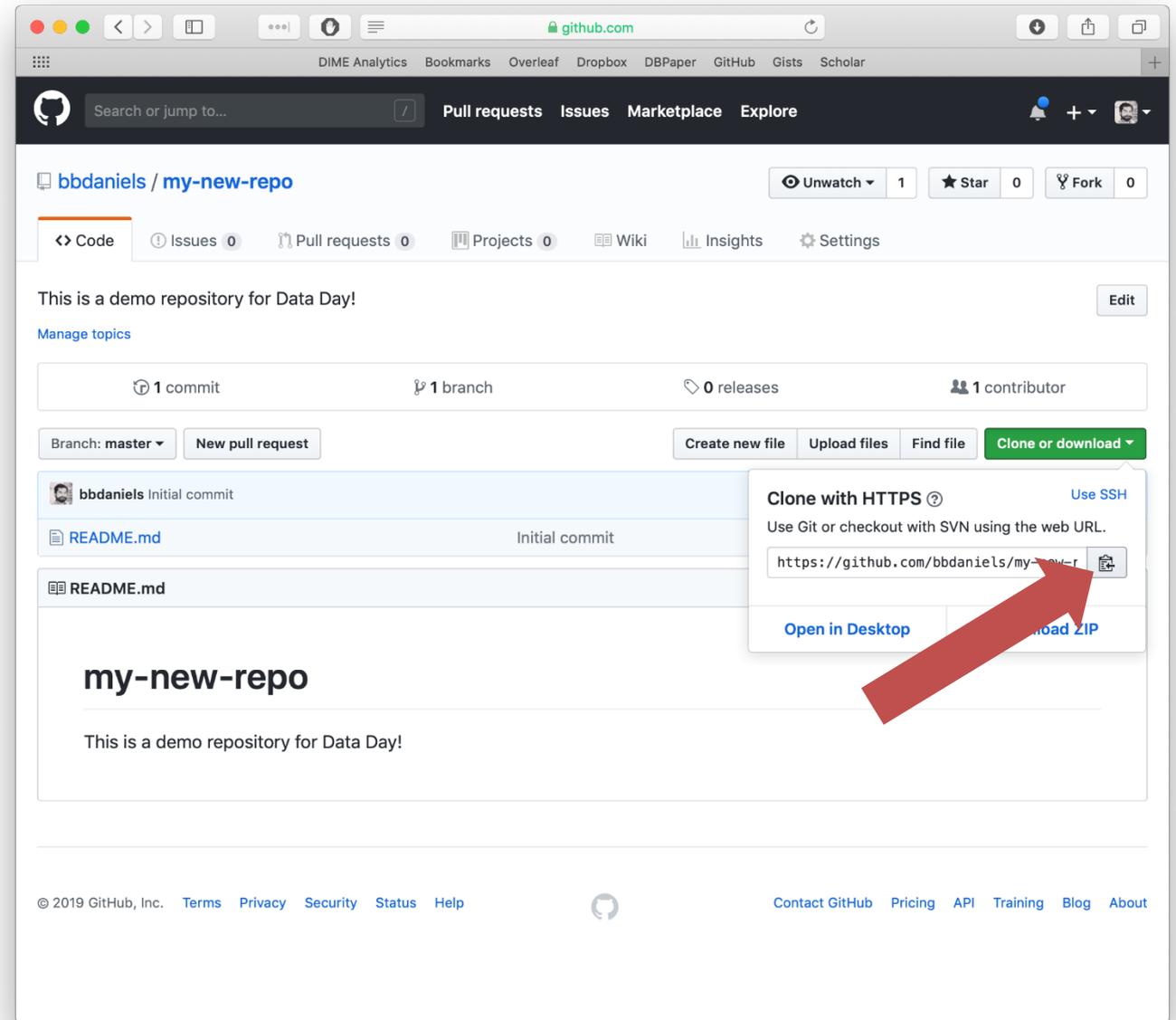
URL

Full path

Clone the repo with your desktop client

Cloning the repository makes a complete copy of it at the new location, including the entire history (remember, the repository is the history).

Cloning from GitHub to your local computer is a good way to start a repository.

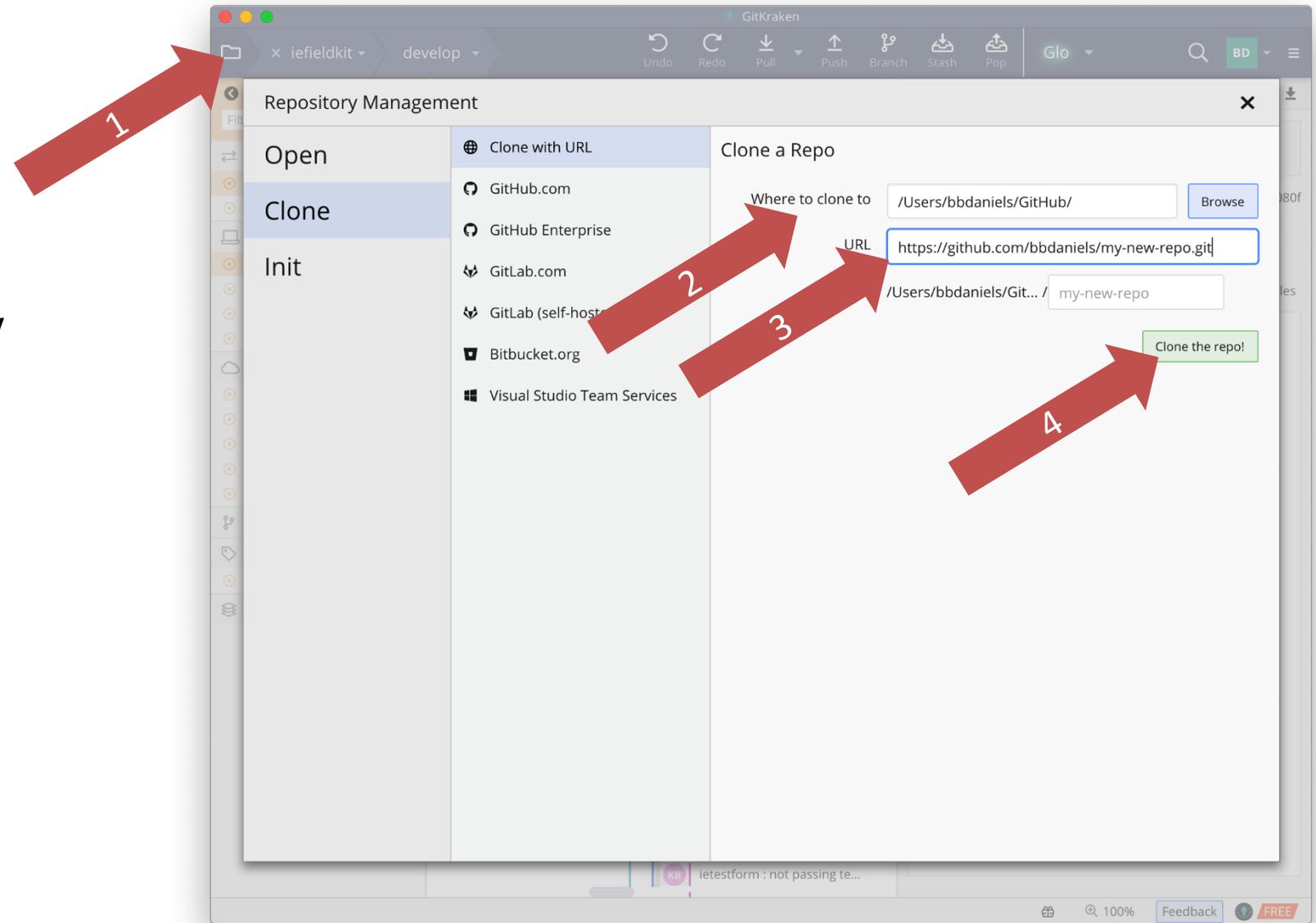


Clone the repo with your desktop client

Cloning the repository makes a complete copy of it at the new location, including the entire history (remember, the repository is the history).

Cloning from GitHub to your local computer is a good way to start a repository.

(I know it seems like a lot of steps now, but this becomes second nature. Bear with me, and use this presentation as a guide!)



Initialize Git Flow in “Preferences”

Exit Preferences

Default Profile
Benjamin Daniels
bbdaniels@gmail.com

General

Profiles

Authentication

UI Preferences

Repo-Specific Preferences
git-flow-demo
/Users/bbdaniels/GitHub/git-flow-demo/

Git Flow

Git Flow

Branches

Master

Develop

Prefixes

Feature

Release

Hotfix

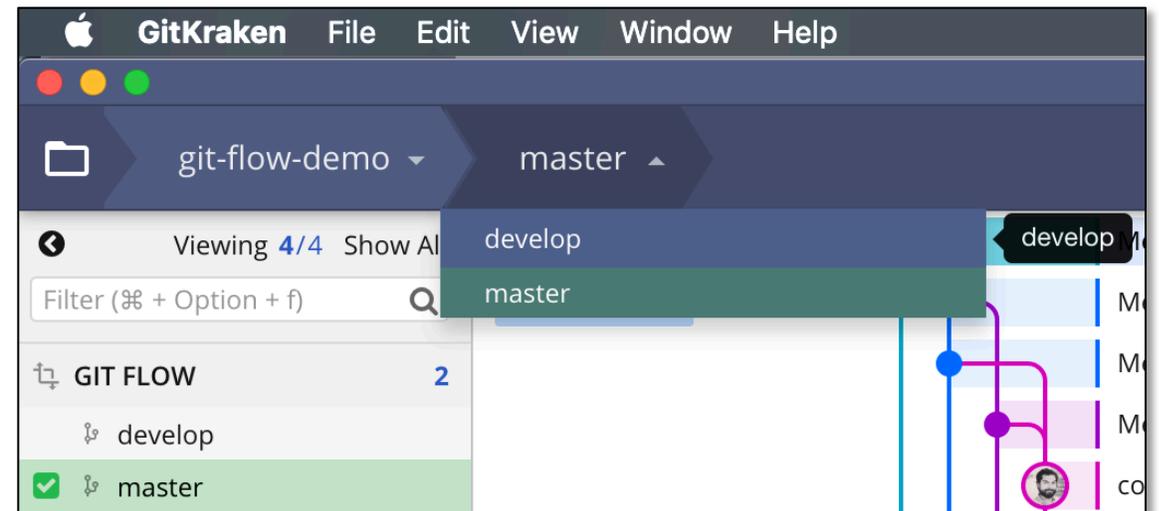
Version Tag

[Initialize Git Flow](#)

Checkout “develop”

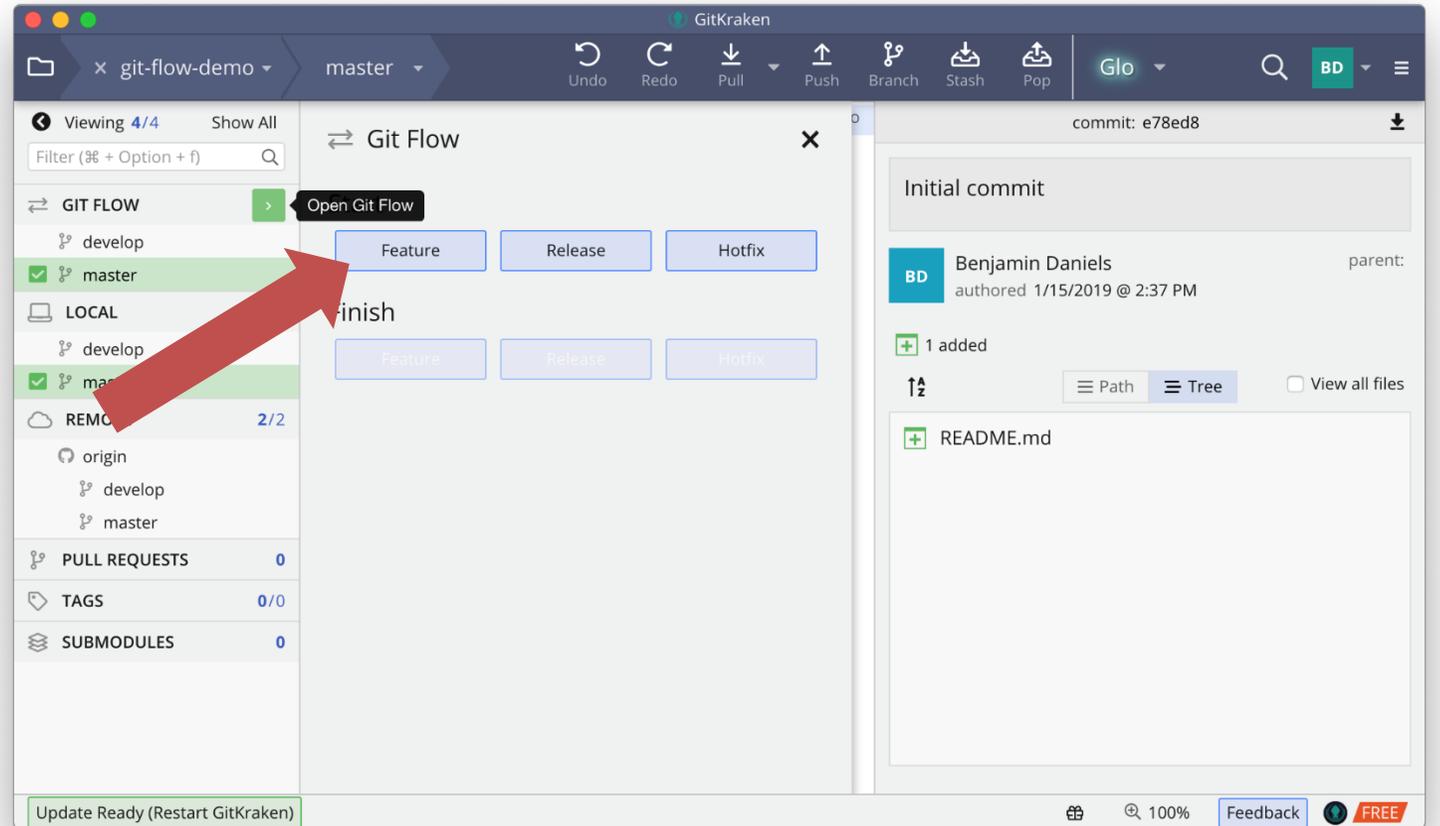
When you “check out” a branch, everything on your computer is instantly updated to what that branch looks like.

Always be conscious to remember “where” you are – this is the core functionality of Git. You’ll get used to it!



Start a new [feature] branch

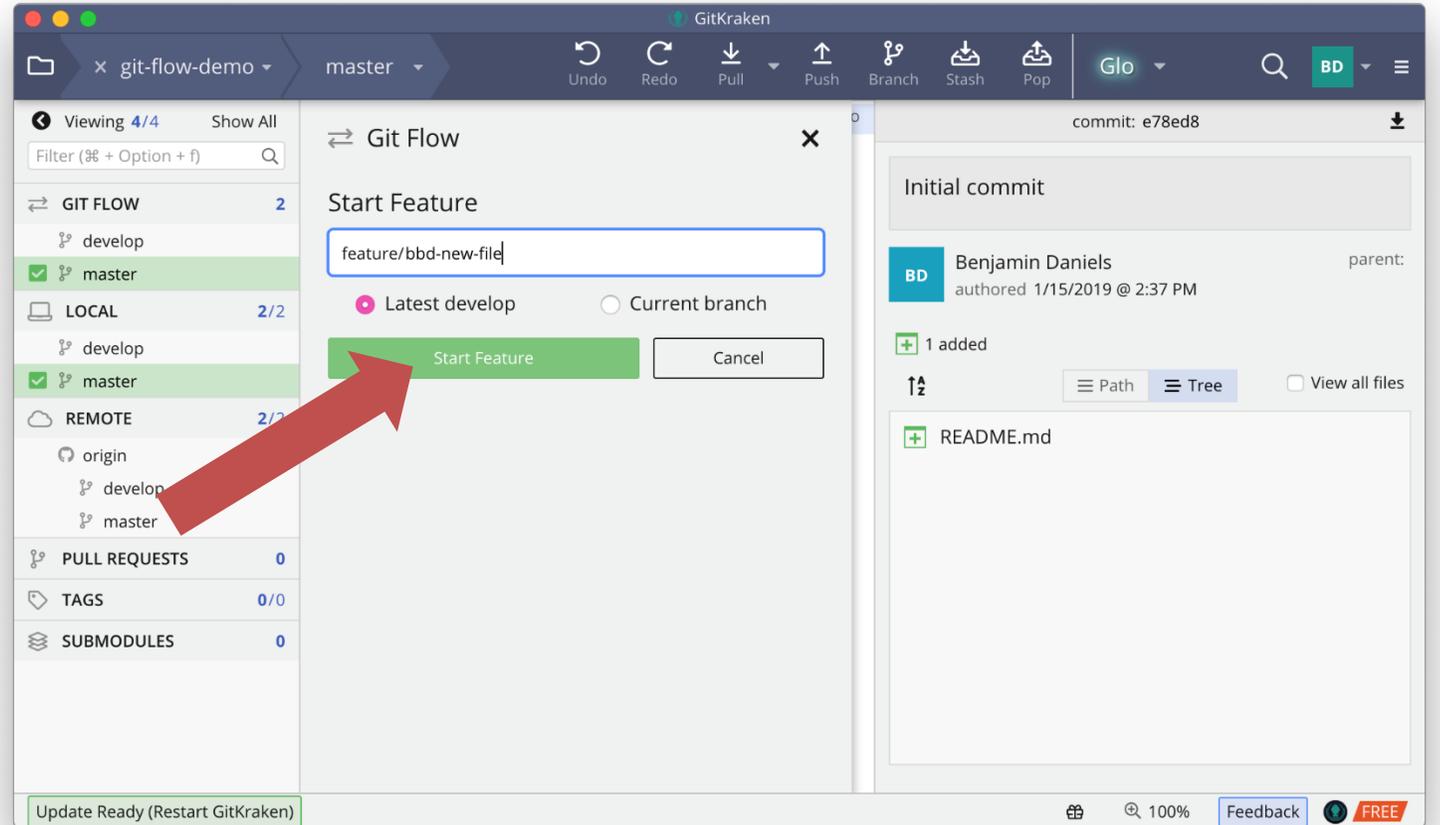
On your local machine, use the Git Flow popout to start a new feature. Call it [xx-new-file] where xx is your initials.



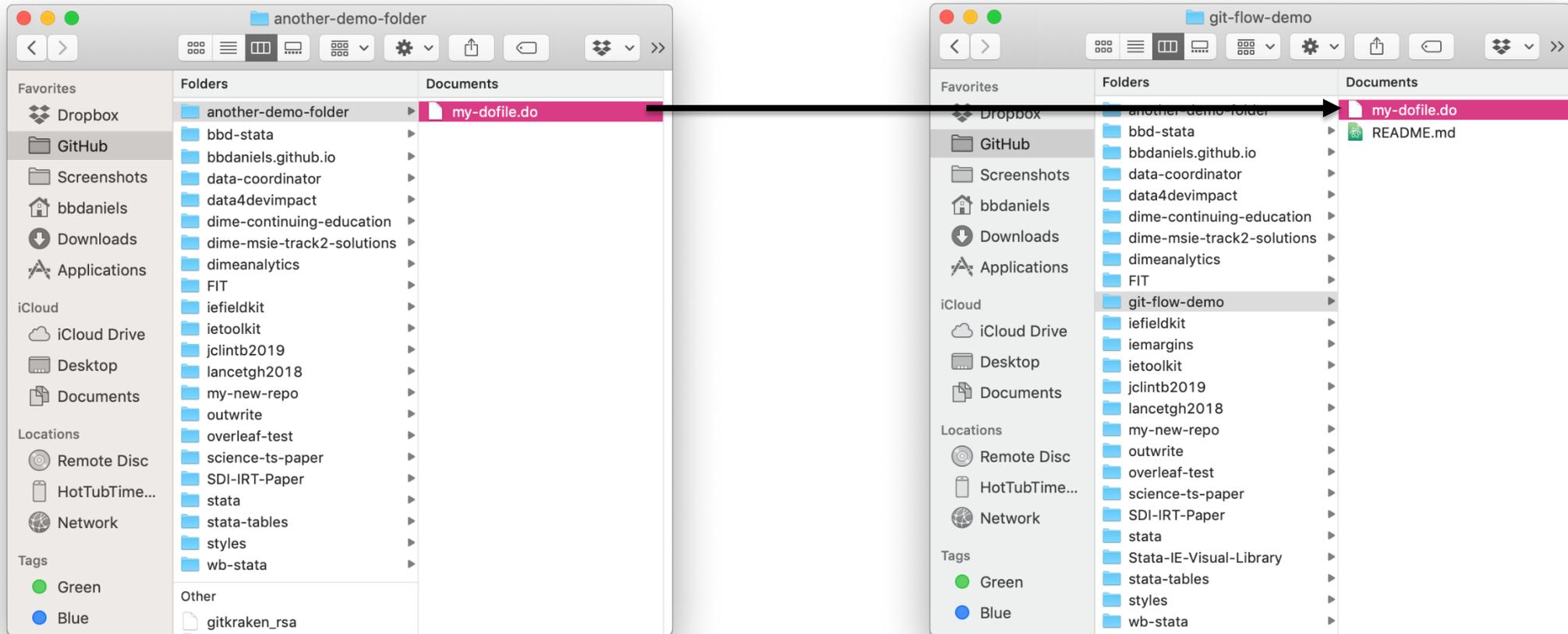
Start a new [feature] branch

On your local machine, use the Git Flow popout to start a new feature. Call it [xx-new-file] where xx is your initials.

[feature] branch names should be single-person and single-purpose: the name is a great place to be clear about those.

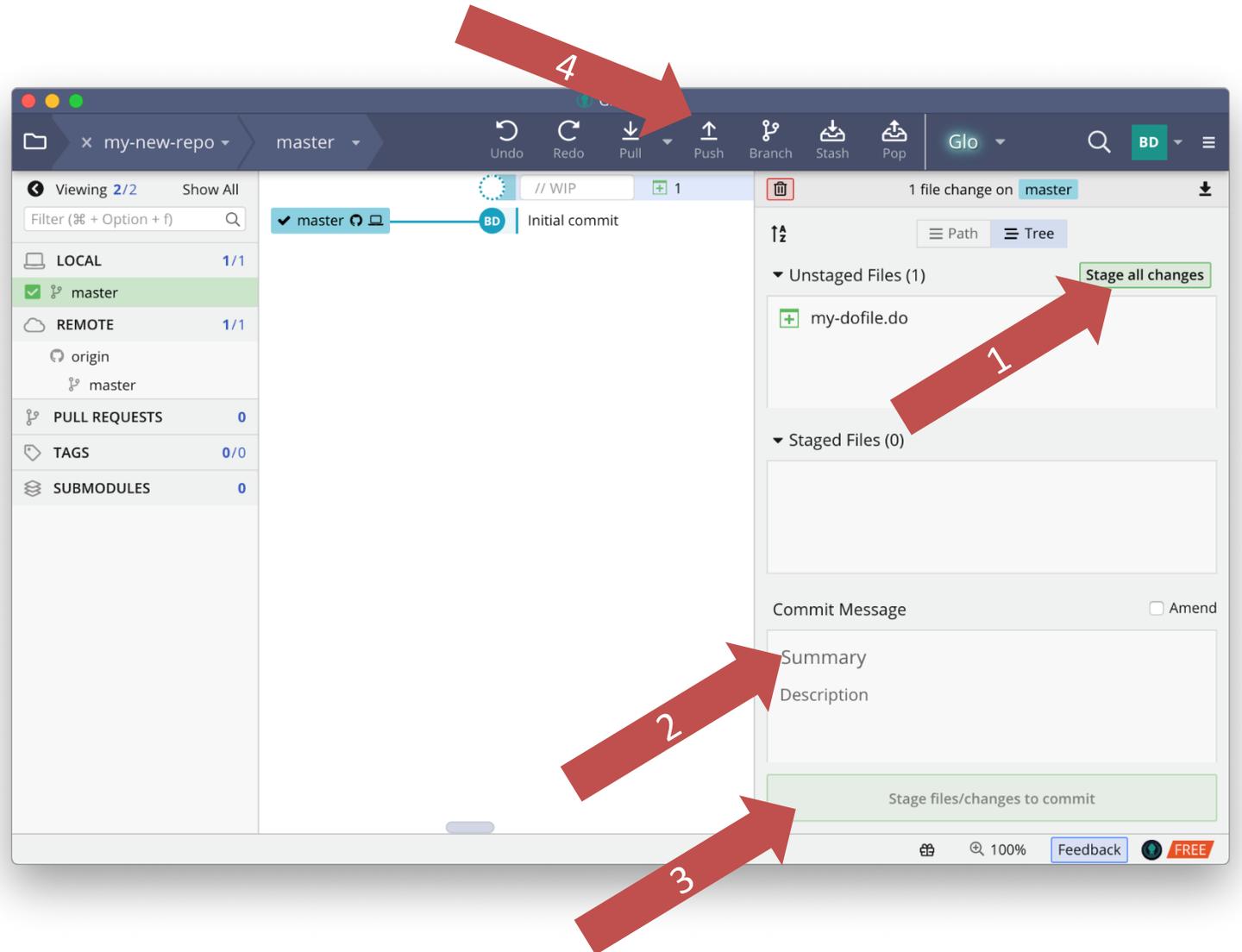


Add some kind of file locally (xx-file initials)



Git notes changes as “work in progress”

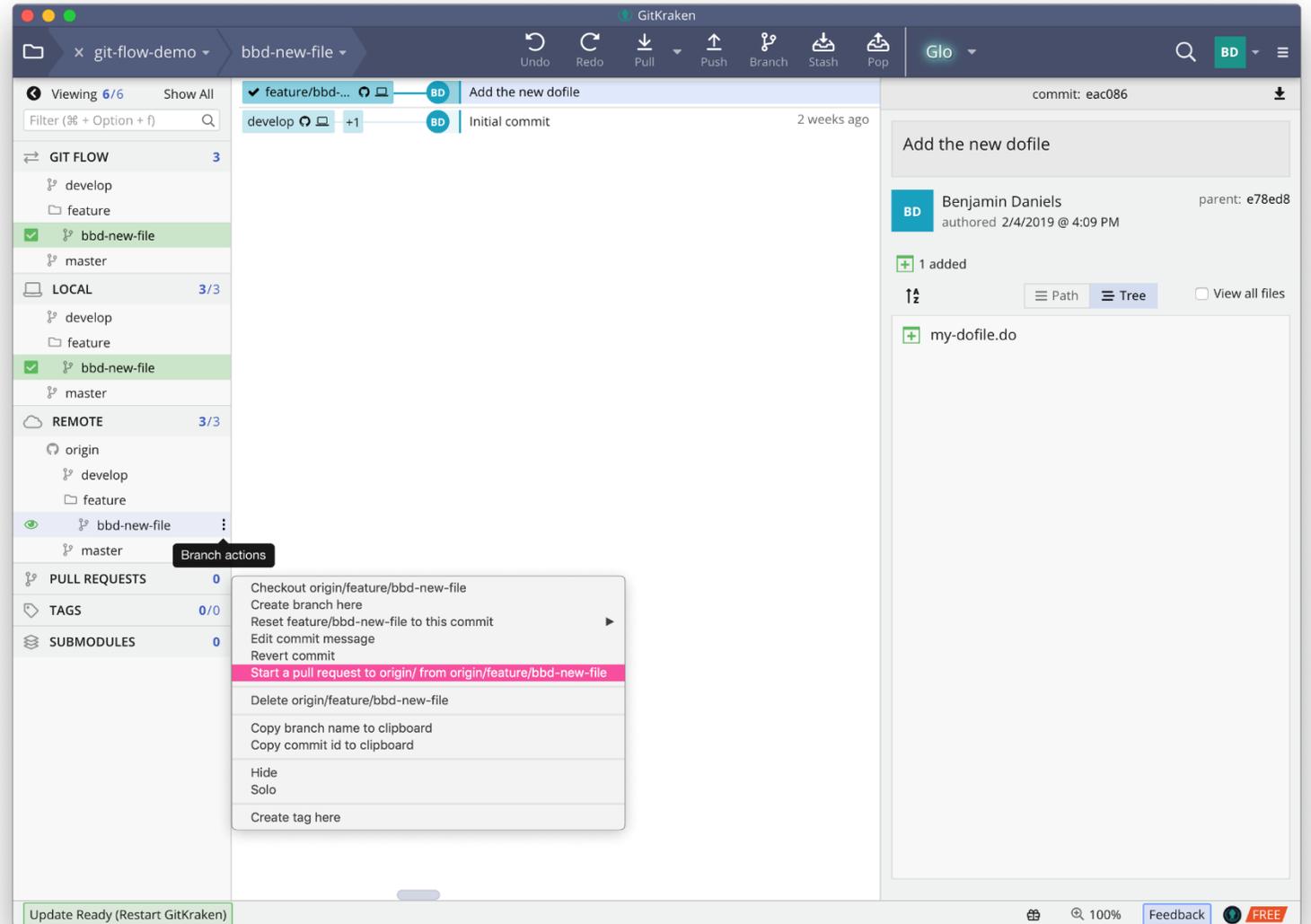
1. “Stage” your changes to add them to the queue to be committed.
2. “Name” your changes informatively – a short sentence will do nicely.
3. “Commit” your changes to add them to the version history.
4. “Push” your changes to sync the remote (GitHub) repository with your local copy.



Start a pull request

This can be done by using the context menu from the feature branch under the “remote/origin” header – because this is a GitHub operation, not a Git operation

You can also drag-and-drop the feature branch onto “origin/develop”.



Make sure you target *develop*

Next, as the administrator, I'll resolve these on GitHub.

Once the branch is pushed and PRed, you can delete the branch locally.

The screenshot shows the GitKraken application interface for creating a Pull Request. The main window is titled "Create Pull Request" and displays the "From Repo" and "To Repo" sections. The "From Repo" is set to "bbdaniels/git-flow-demo" and the "Branch" is "feature/bbd-new-file". The "To Repo" is also "bbdaniels/git-flow-demo", and the "Branch" dropdown menu is open, showing "develop" selected. A red arrow points to the "develop" option in the dropdown. The "Title" field contains "Add the new dofile" and the "Description" field contains "Pull request description". The "Reviewers" field is empty, and the "Assignees" field is empty. The "Labels" field is empty. The "Create Pull Request" button is highlighted in green. The left sidebar shows the "GIT FLOW" section with "develop" and "feature" branches, and the "LOCAL" section with "develop", "feature", and "bbd-new-file" branches. The "REMOTE" section shows "origin" with "develop", "feature", "bbd-new-file", and "master" branches. The "PULL REQUESTS", "TAGS", and "SUBMODULES" sections are empty. The bottom status bar shows "Update Ready (Restart GitKraken)".

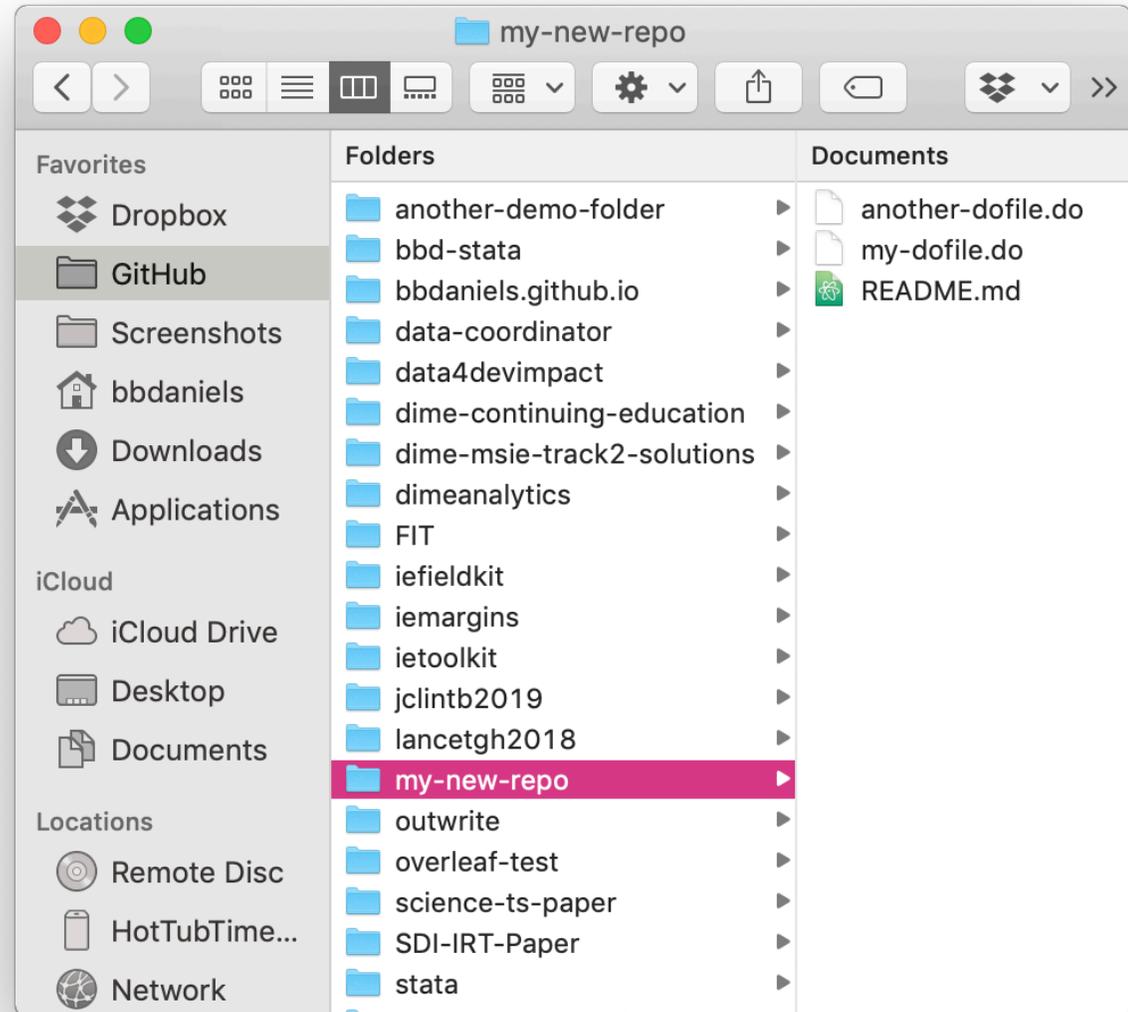
Merging consolidates changes

In the local working directory, you will now see that *both* files now exist in the same commit, “Merge pull request #1 from bbdaniels/feature/...”, on the *develop* branch.

Then I will merge these to *master* so that all those changes are only reflected in a single *master* commit: a new version!

If people made *conflicting* changes to the same files, this process also gives you a chance to resolve those, and this and other workflows will be covered in a later session.

That’s all for now!



Thank you!

