Git and version control in data science

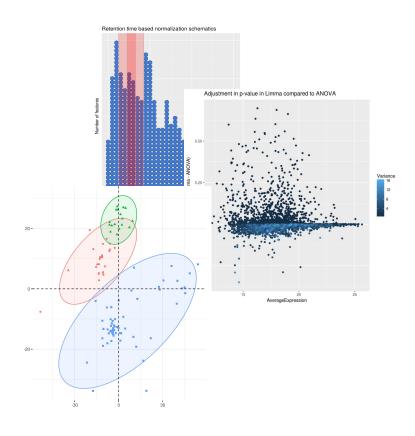
Jakob Willforss, Department of Immunotechnology, LTH 24th April 2020

I use Git when I...

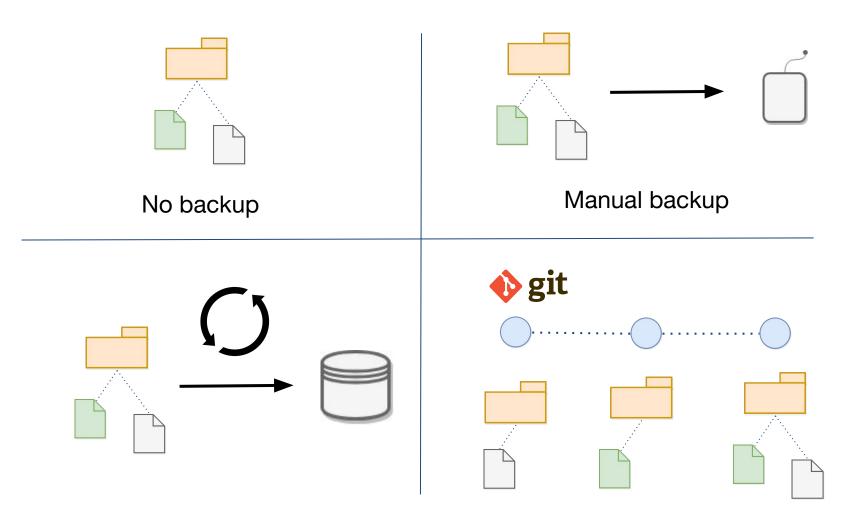
...develop software



...analyze data



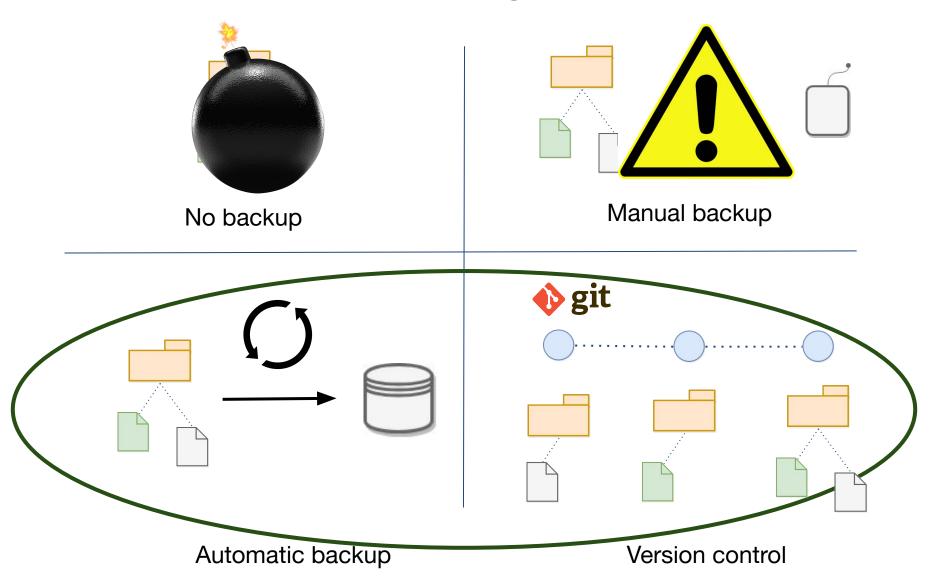
How we manage our files



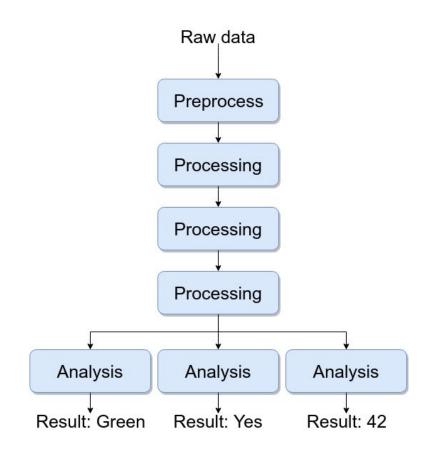
Automatic backup

Version control

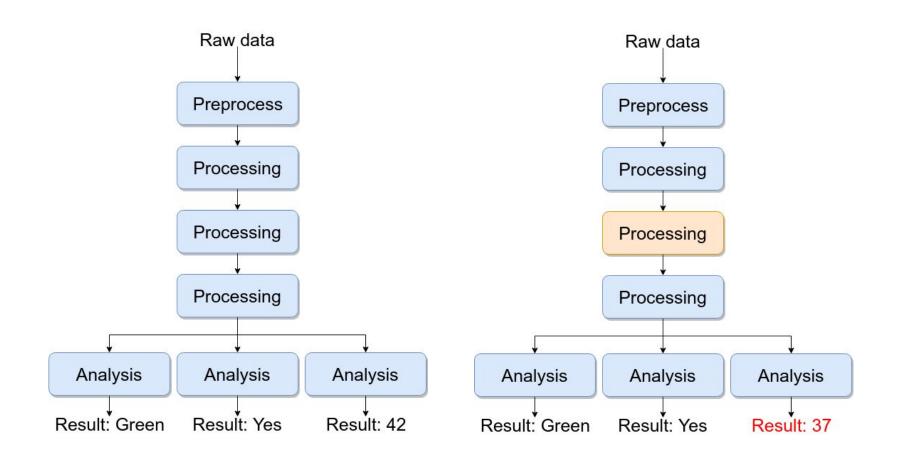
How we manage our files



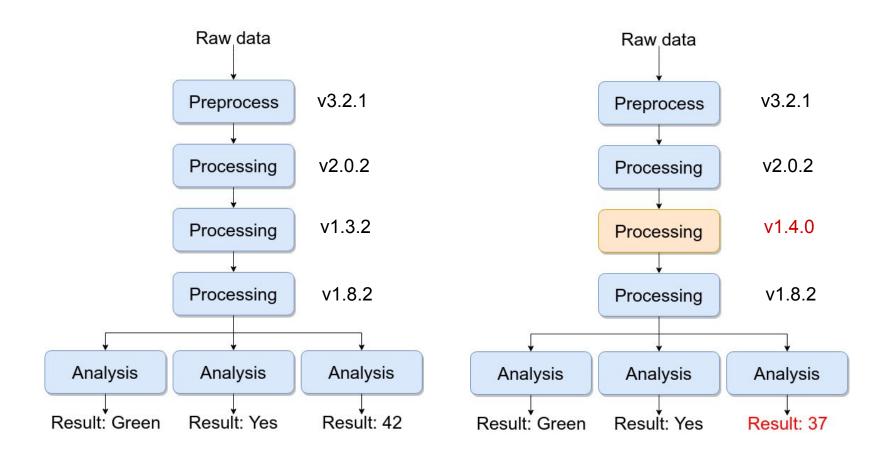
Reproducible research



Reproducible research

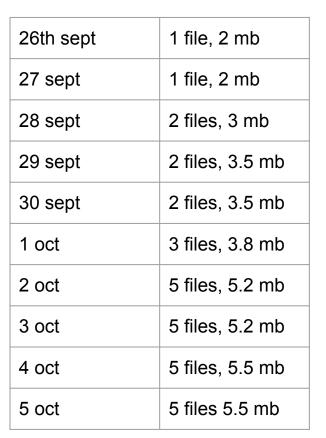


Reproducible research



Version control vs. backup







26th sept	Add first file	Tags
26 sept	Make important edits	
29 sept	Include second analysis	v2.0 Analysis 2
30 sept	Extend second analysis	
30 sept	Add visualizations for first analysis	
2 oct	Found errors in ANOVA, corrected	
2 oct	Preparing third analysis	
4 oct	Start third analysis	v3.0 Analysis 3

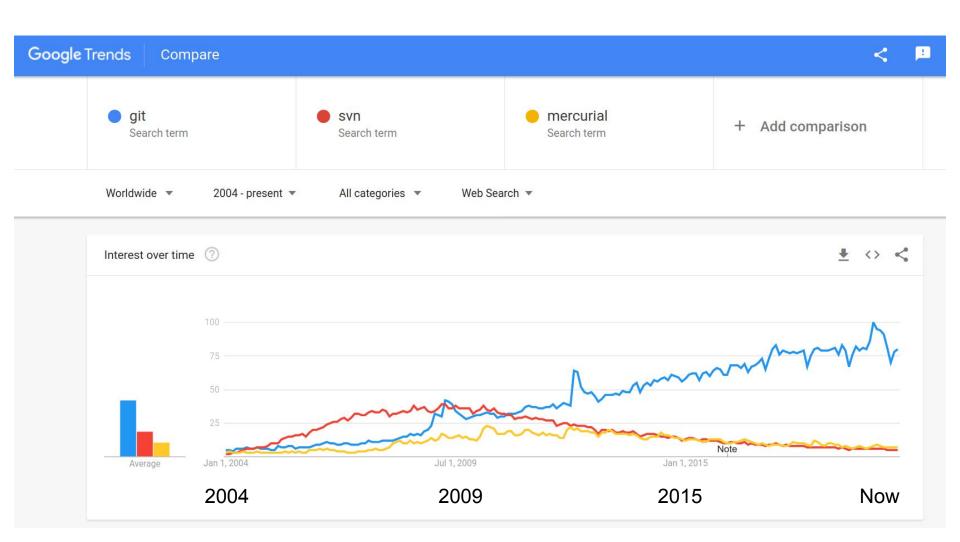
Git

Developed 2005 by Linus Torvald to maintain Linux source code



Linus Torvald

From wikimedia commons



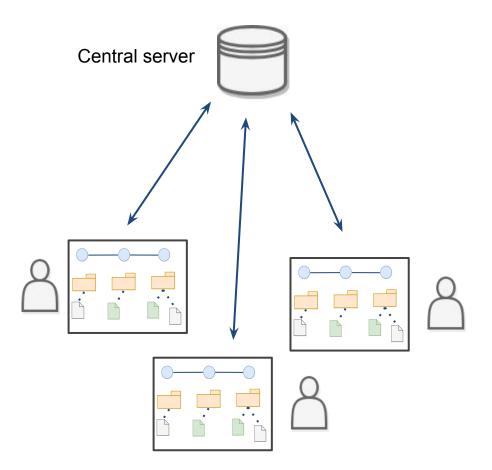
Centralized version control

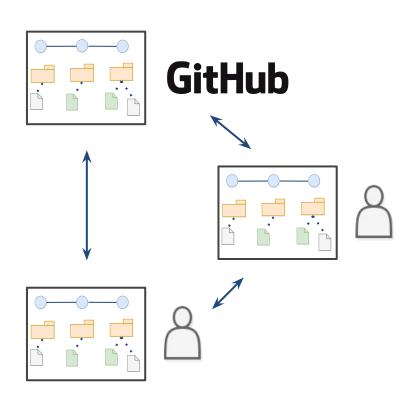
Distributed version control











What does version control do?

Determine what changed, who changed it and why

Organized way of collaborating in code

Way of presenting code

Allow navigating in history of files

Why don't everyone use it?

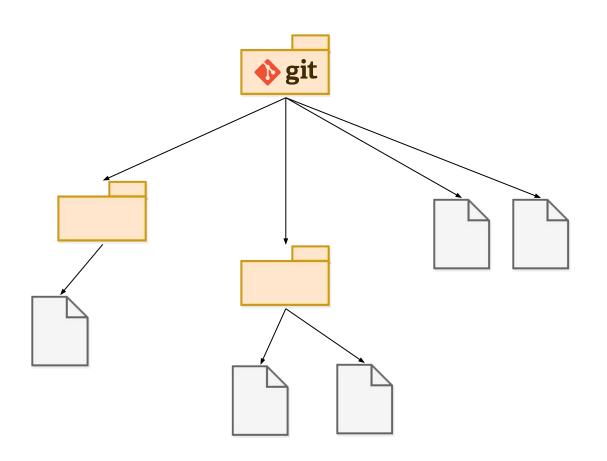
Learning threshold

Tricky to manage data

Let's dive in

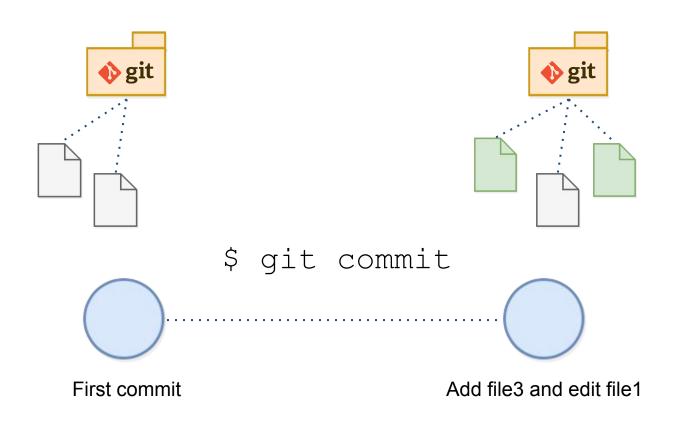


The 'repository' and the 'file tree'



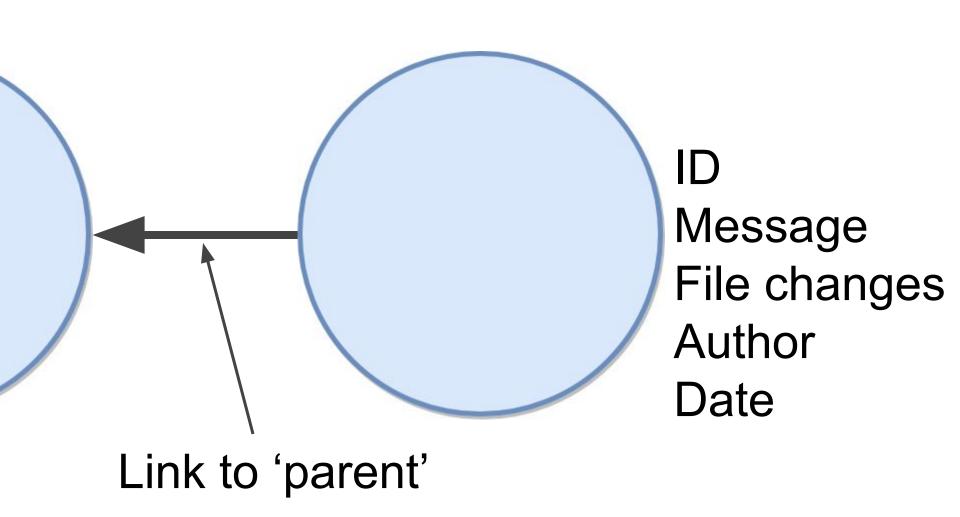
The 'repository' keeps track of changes in the 'file tree'

What is a commit?

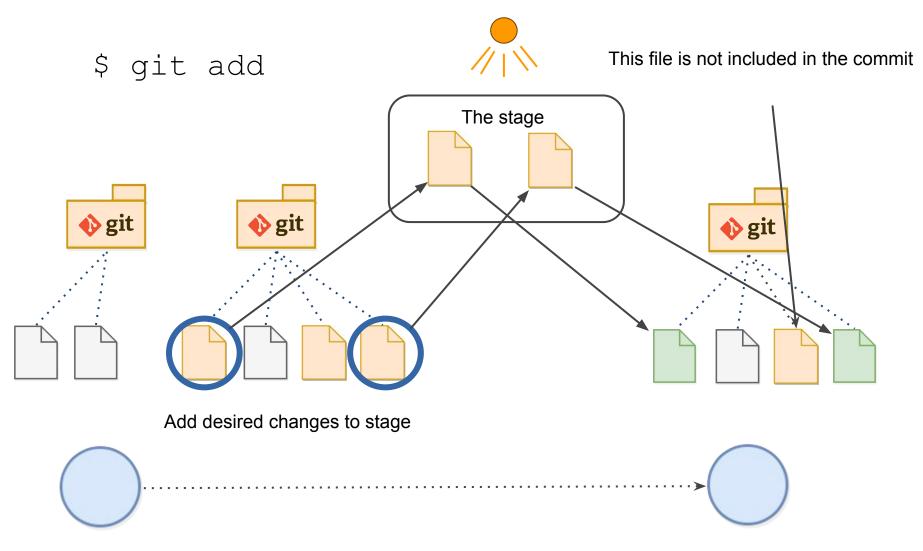


A snapshot of particular state in the file tree

What is a commit?



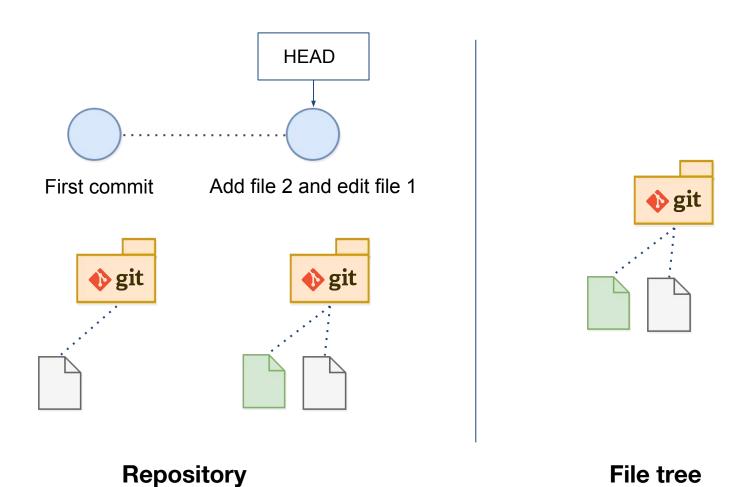
What is the 'stage'?



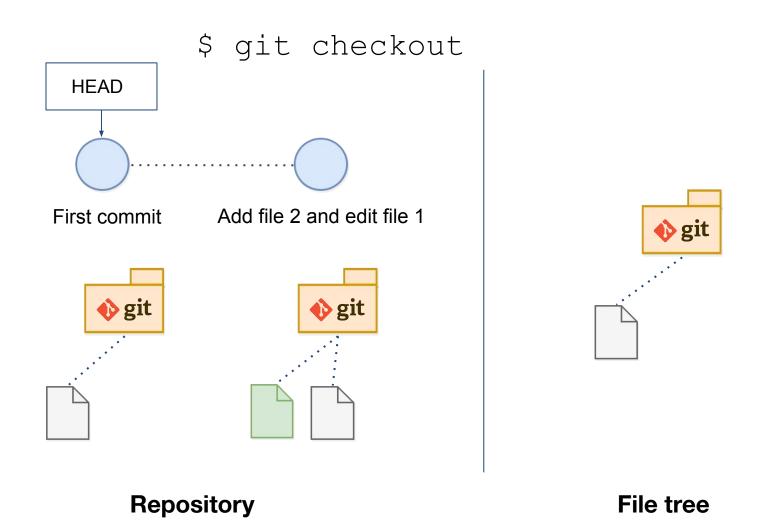
First commit

Edit file 1 and add file 4

Navigating the history and the HEAD

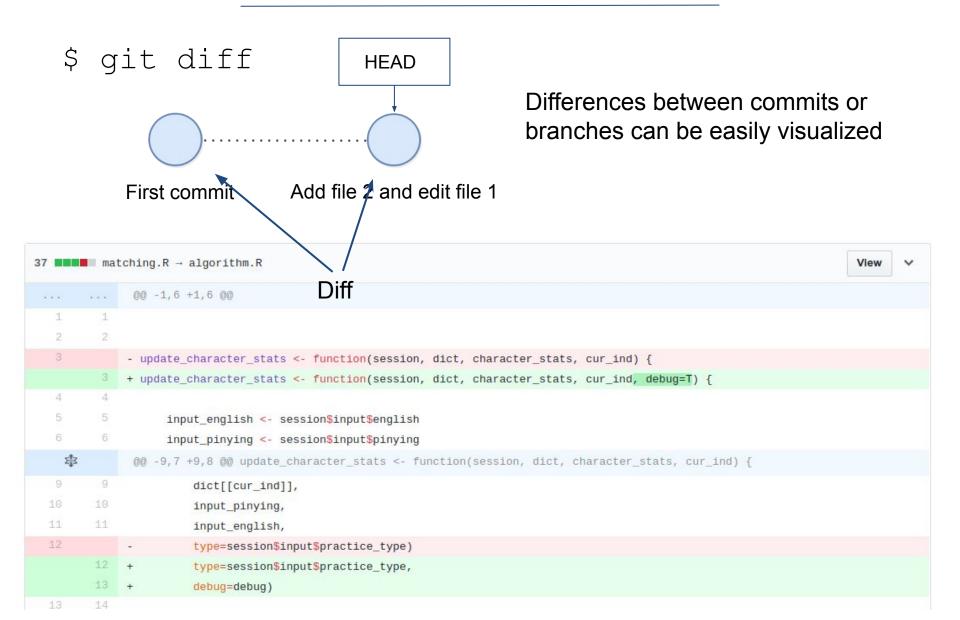


Navigating the history and the HEAD



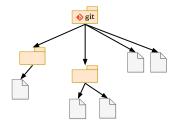
By moving the HEAD, the file tree changes

Comparing changes



Recap

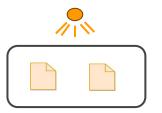
Repository and file tree

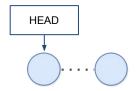




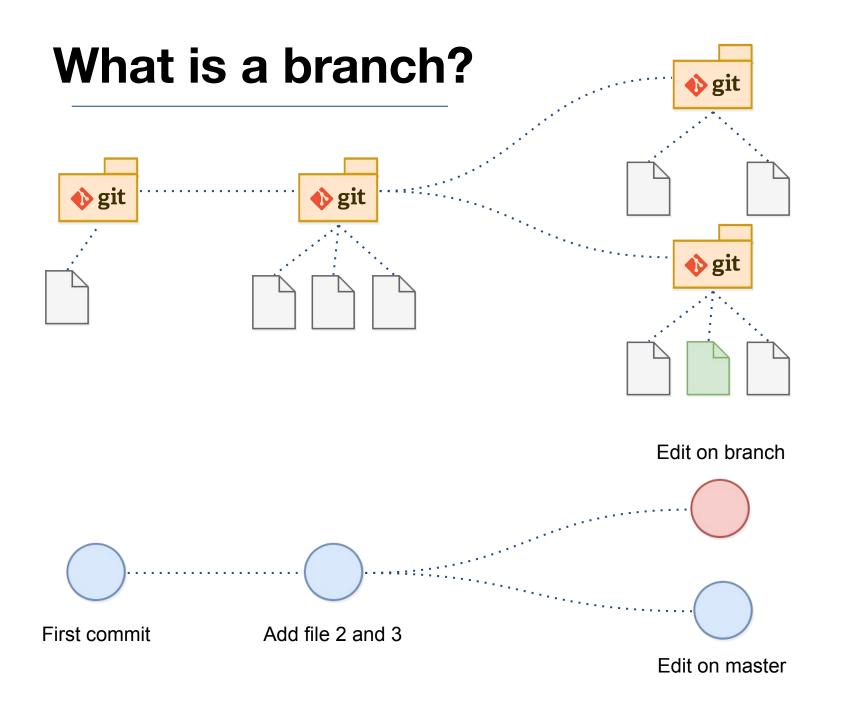
Commit

Stage

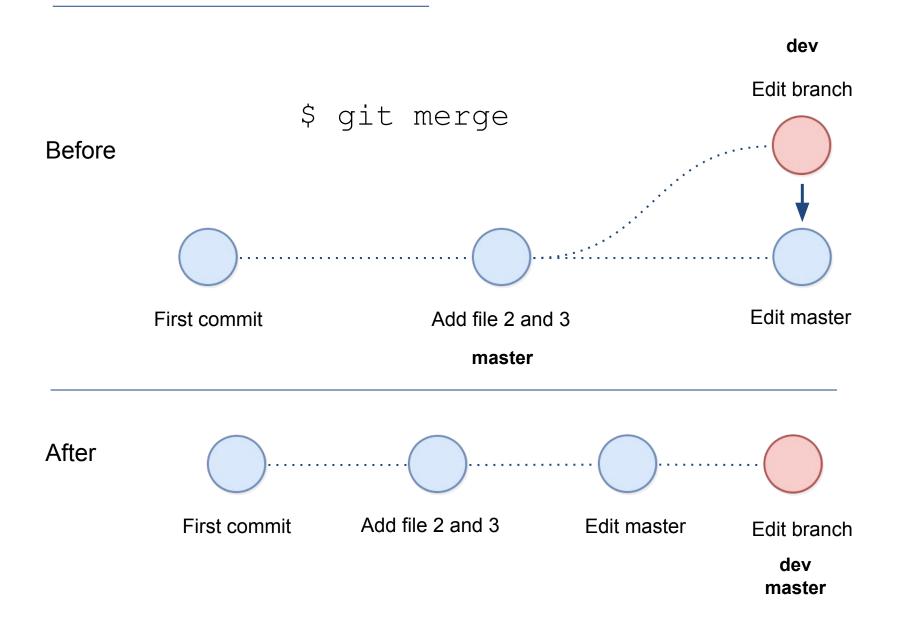




The HEAD



Merge branches

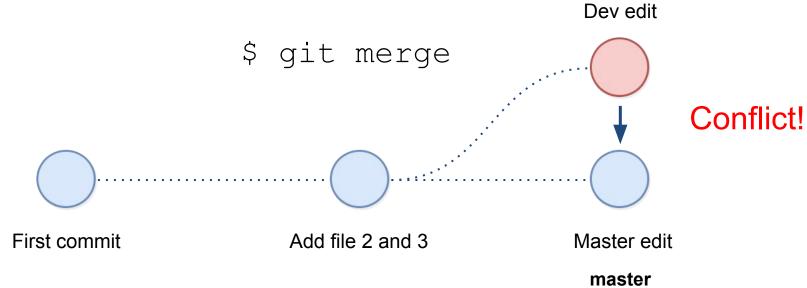


Merge conflicts

If merging conflicting changes,
you get a "merge conflict"

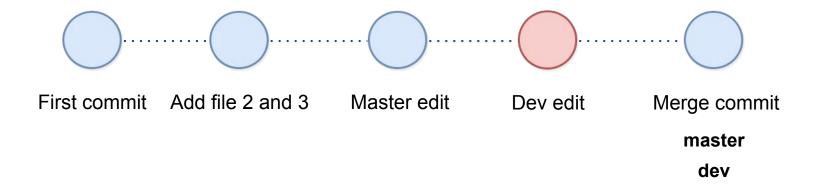
dev

Dev ex



Git stores both changes - You decide what to keep!

Resolved merge conflict

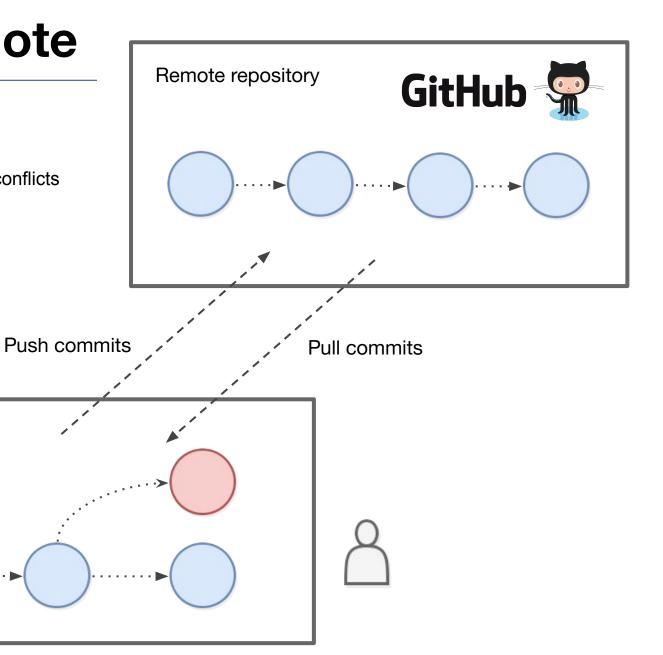


After resolving all commits are there, followed by a merge commit

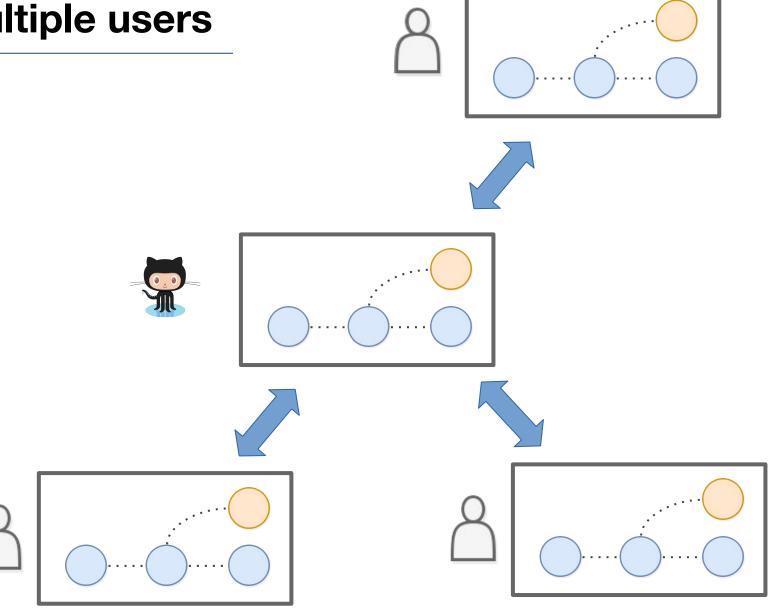
The remote

We can get merge conflicts here too!

Local repository

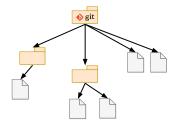


Multiple users



Recap

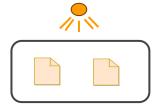
Repository and file tree

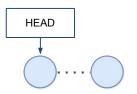




Commit

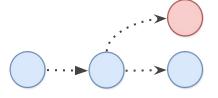
Stage

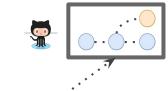




The HEAD

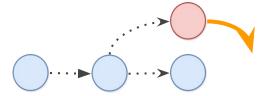
Branch





Remote

Merging



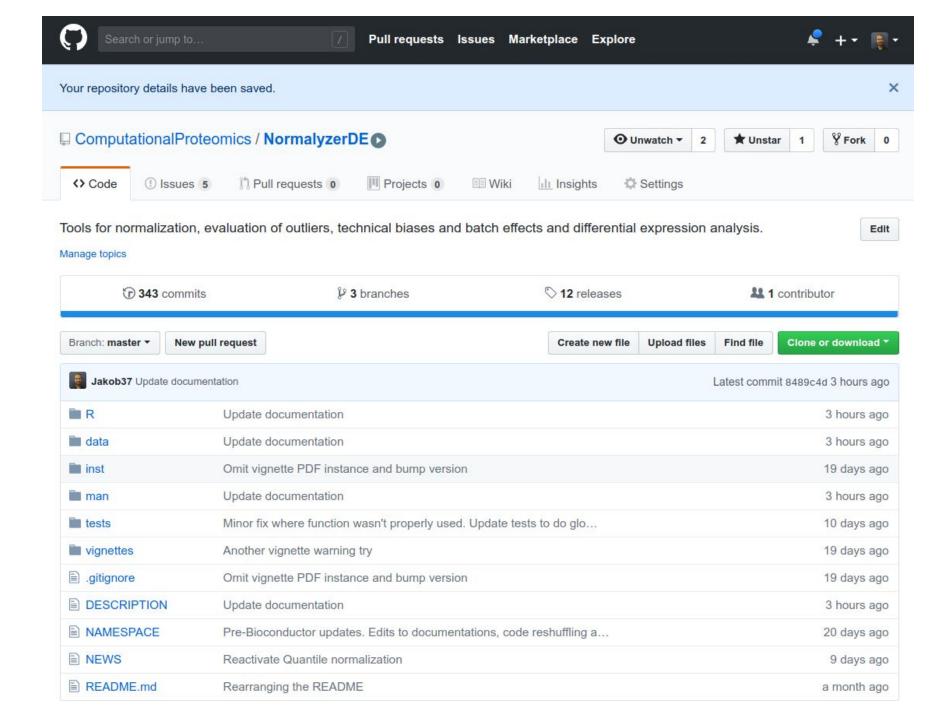
GitHub 🐃

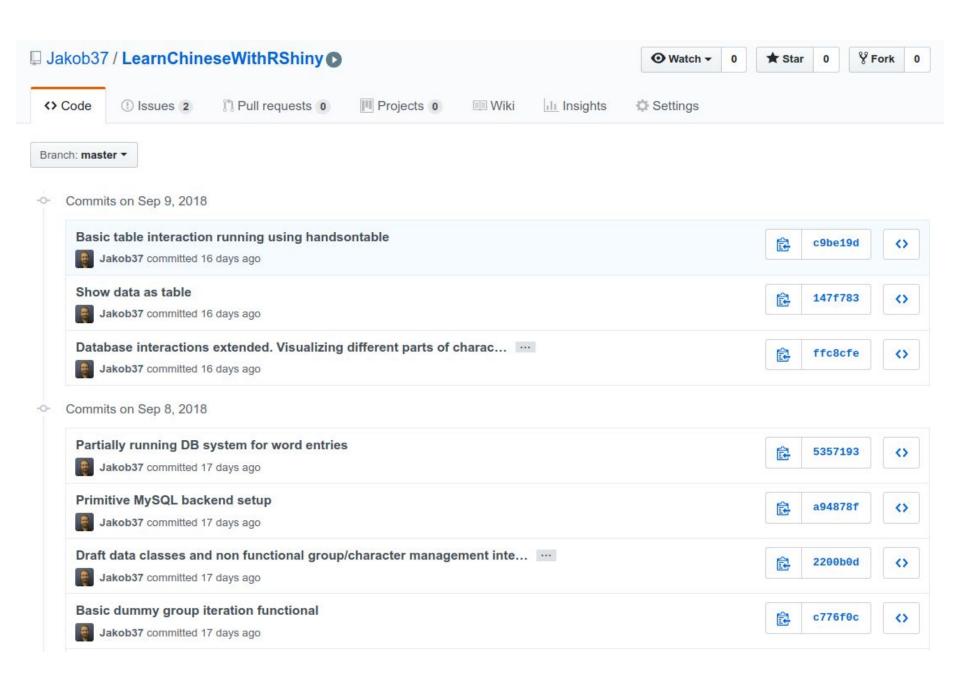


Social platform for code

Common way to make code public

Allows interacting with other peoples code





Updated working algorithm **Browse files** p master Jakob37 committed on Aug 2 1 parent 85aa4de commit c9a8d66b39c8a247ffd9dc87b7ad512724761637 Showing 4 changed files with 57 additions and 32 deletions. Unified Split 37 ■■■■ matching.R → algorithm.R View 00 -1,6 +1,6 00 2 - update_character_stats <- function(session, dict, character_stats, cur_ind) {</pre> + update_character_stats <- function(session, dict, character_stats, cur_ind, debug=T) { 4 4 5 5 input_english <- session\$input\$english 6 input_pinying <- session\$input\$pinying</pre> 牵 00 -9,7 +9,8 00 update character stats <- function(session, dict, character stats, cur_ind) { 9 9 dict[[cur_ind]], input_pinying, input_english,

character_stats[[cur_ind]]\$right <- character_stats[[cur_ind]]\$right + 1

character_stats[[cur_ind]]\$wrong <- character_stats[[cur_ind]]\$wrong + 1

00 -18,10 +19,15 00 update_character_stats <- function(session, dict, character_stats, cur_ind) {

type=session\$input\$practice_type)

type=session\$input\$practice_type,

debug=debug)

if (correct) {

12 +

13 +

串

Getting started

How to use it

Use it from terminal: Or graphical interface:

https://git-scm.com https://desktop.github.com

How to learn it

Visualizing Git - See the commands in the browser https://git-school.github.io/visualizing-git

Codecademy - Interactive online tutorial (free): https://www.codecademy.com/learn/learn-git

DataCamp - Comprehensive online tutorials (not free): https://www.datacamp.com/courses/introduction-to-git-for-data-science

For deeper understanding: https://www.sbf5.com/~cduan/technical/git

My materials: http://ponderomatics.com/git.html

Handling large data

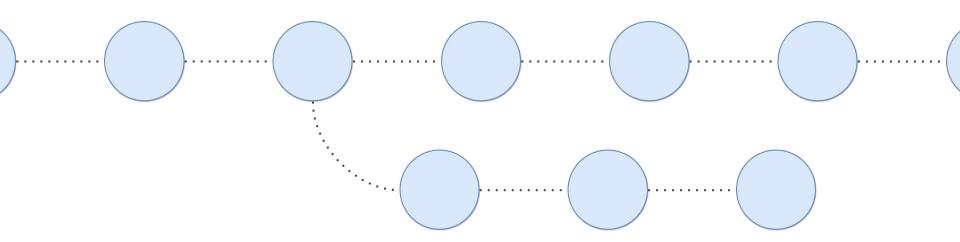
Git Large File Storage

https://git-lfs.github.com/

Quilt

https://quiltdata.com/

Thank you for listening!



Jakob Willforss, jakob.willforss@immun.lth.se