

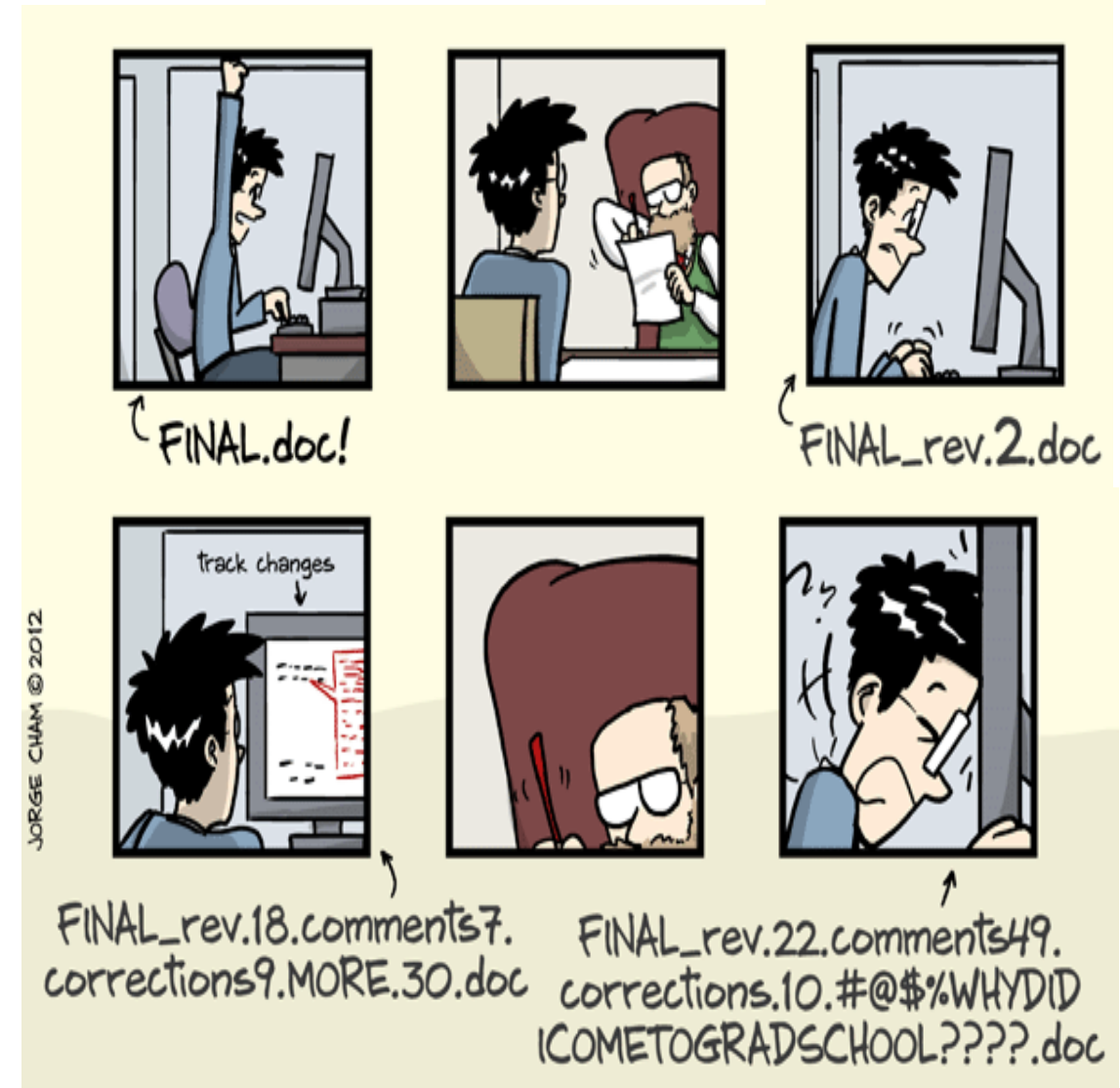


Tools and Programming for Data Science Version Control with Git and GitHub

Study Program Data Science
Prof. Dr. Tillmann Schwörer

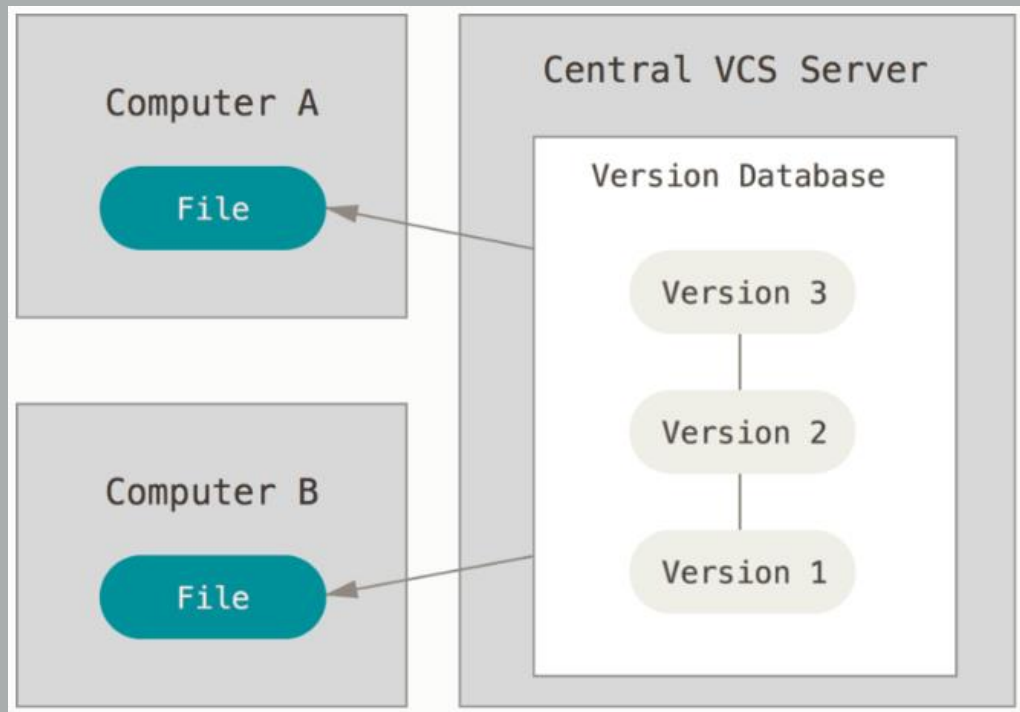
What is a version control system?

- ▶ System to manage changes in code or other documents
 - ◆ Backup: undo changes, restore files, safely experiment
 - ◆ Transparency: who changed? what? when?
 - ◆ Collaboration: work with others on the same project
- ▶ Git: open source, distributed version control system

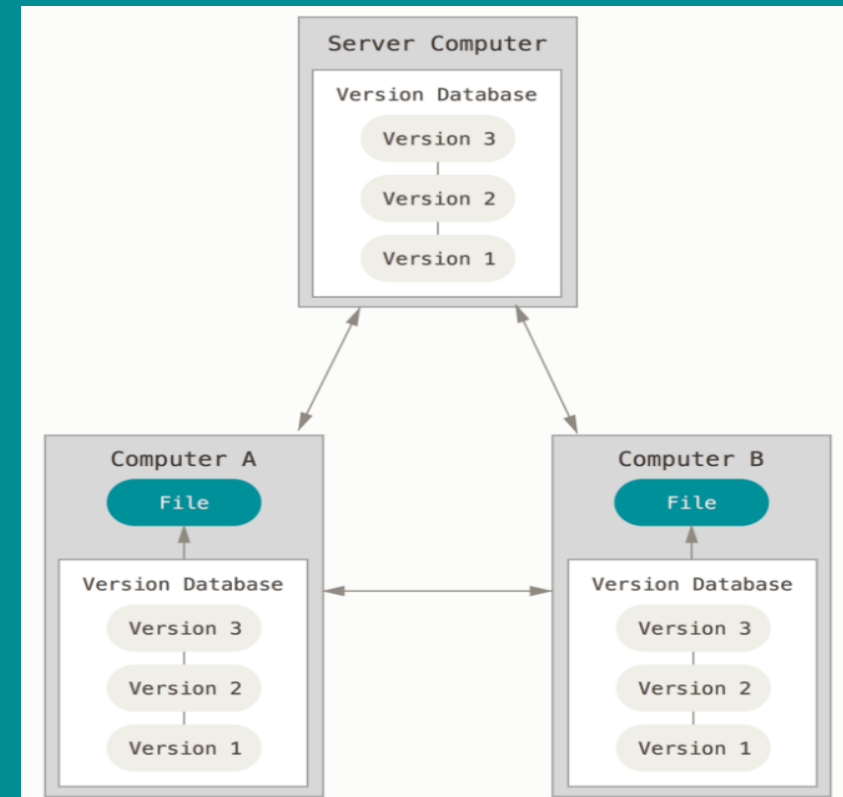


Centralized vs distributed

Centralized: single files are pulled from or pushed to central repository (→ **Subversion**)



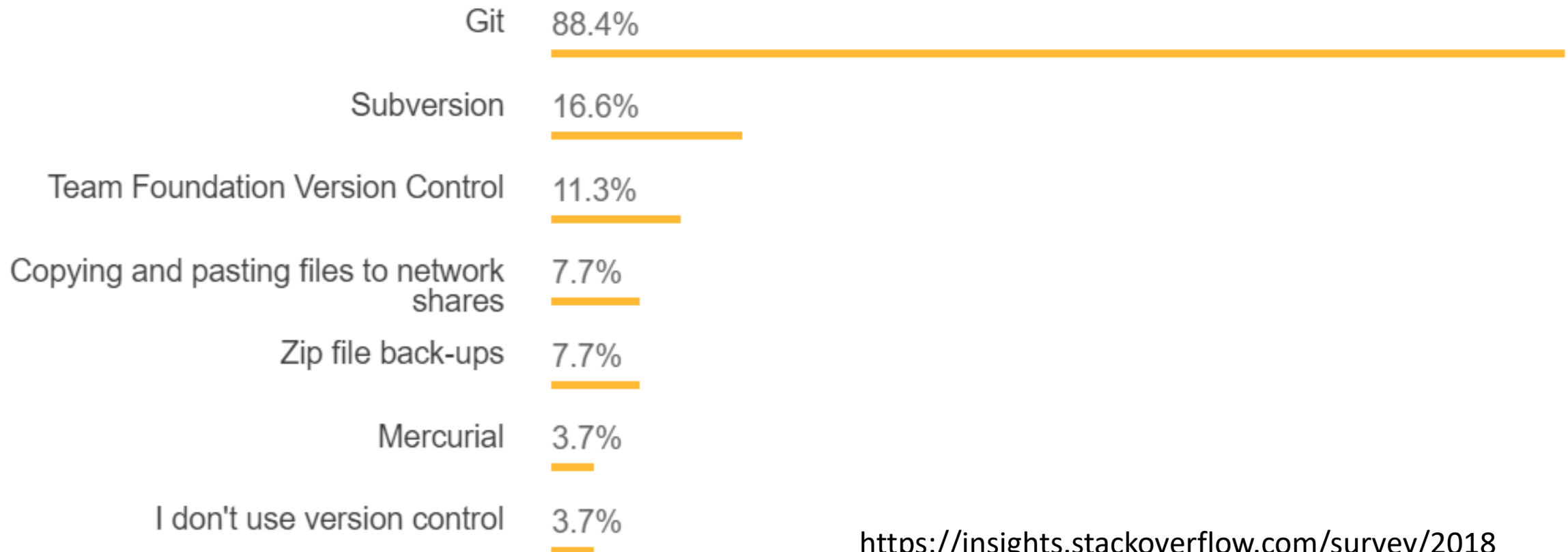
Distributed: Each client has full copy of entire repository (→ **git**)



Most popular version control systems

All Respondents

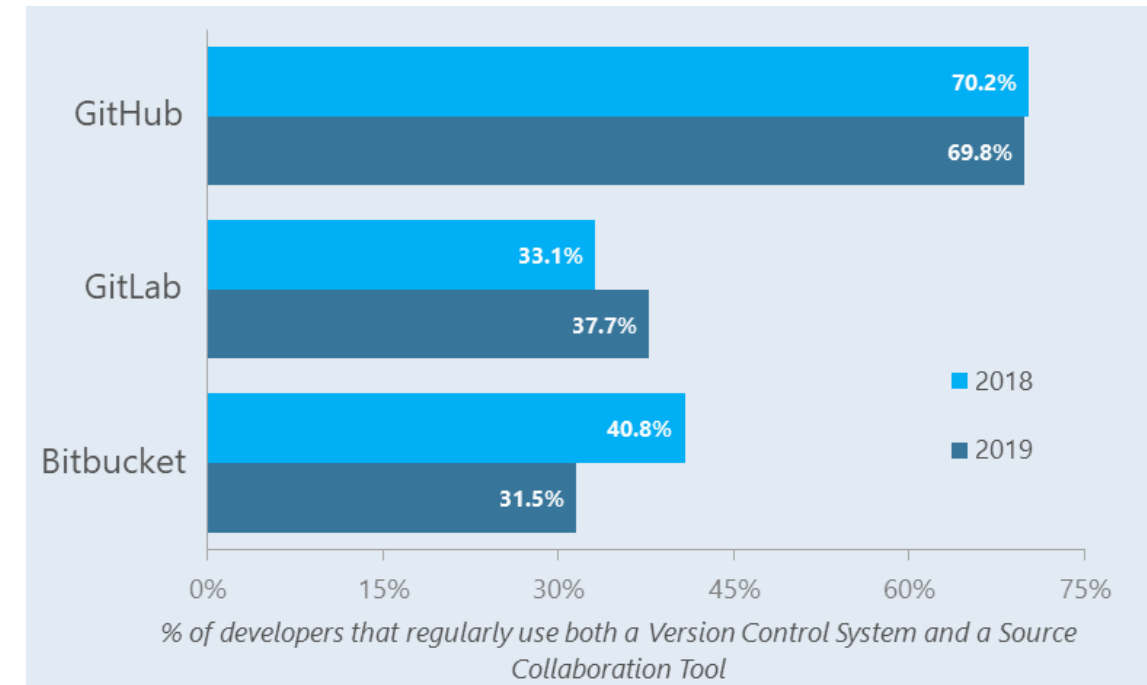
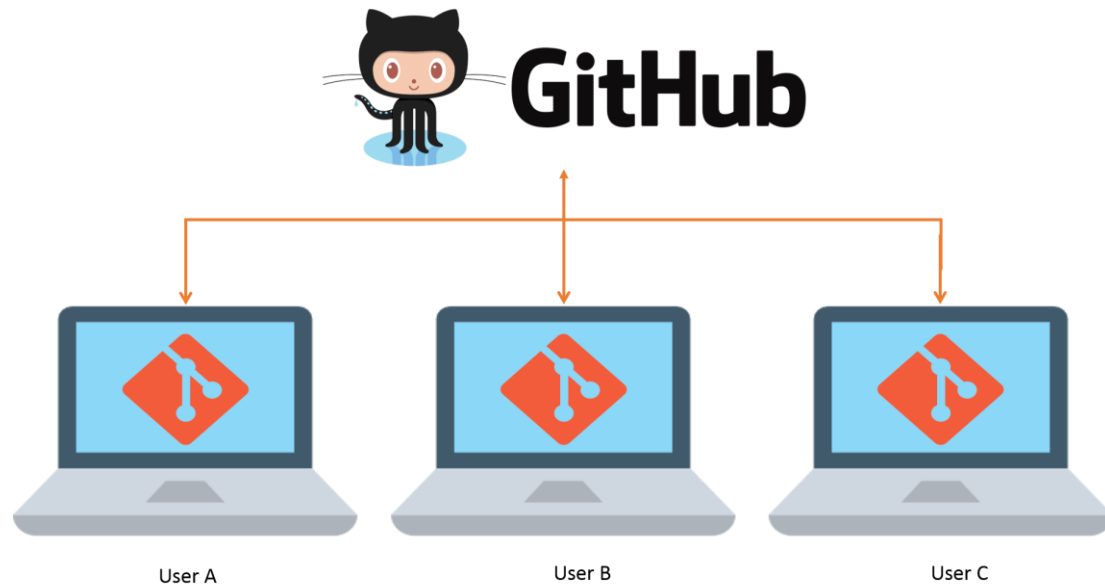
Professional Developers



<https://insights.stackoverflow.com/survey/2018>

How is GitHub different from Git?

- ▶ Hosting service for Git repositories
- ▶ Collaboration platform: bug tracking, feature requests, task management, wikis, etc.



Chaos Computer Club: Entdecke unsere Verfassung



Johannes Rau, Bundespräsident committed on 26 Jul 2002

1 parent dc4dd4d

commit b17a4b2848dac62d37618ebc92c06ccccc5385ff

Showing 1 changed file with 1 addition and 1 deletion.

Unified

Split

2 020a.md

... @@ -1,4 +1,4 @@

1 ## Artikel 20a

2

3 - Der Staat schützt auch in Verantwortung für die künftigen Generationen die natürlichen Lebensgrundlagen im Rahmen der verfassungsmäßigen Ordnung durch die Gesetzgebung und nach Maßgabe von Gesetz und Recht durch die vollziehende Gewalt und die Rechtsprechung.

1 ## Artikel 20a

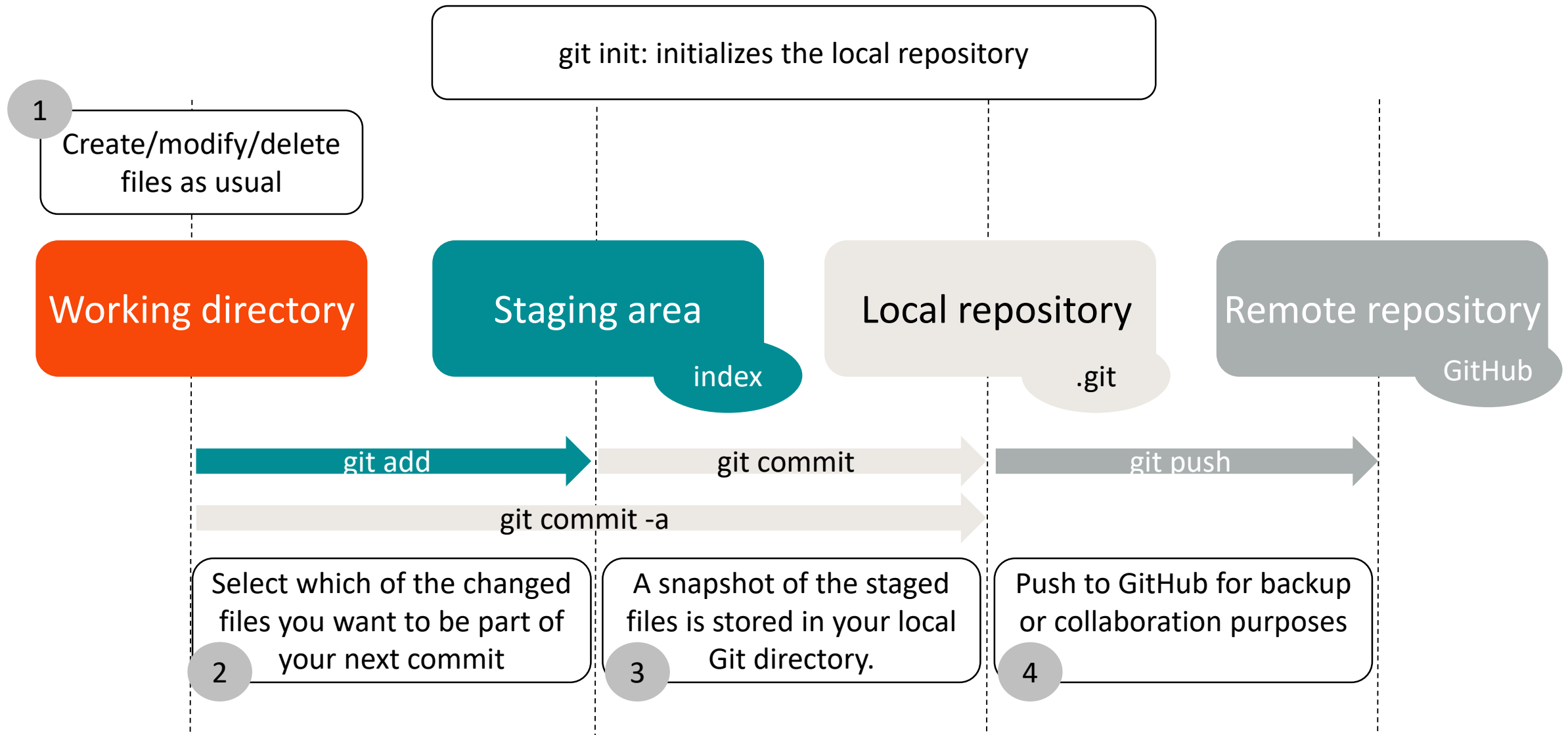
2

3 + Der Staat schützt auch in Verantwortung für die künftigen Generationen die natürlichen Lebensgrundlagen und die Tiere im Rahmen der verfassungsmäßigen Ordnung durch die Gesetzgebung und nach Maßgabe von Gesetz und Recht durch die vollziehende Gewalt und die Rechtsprechung.

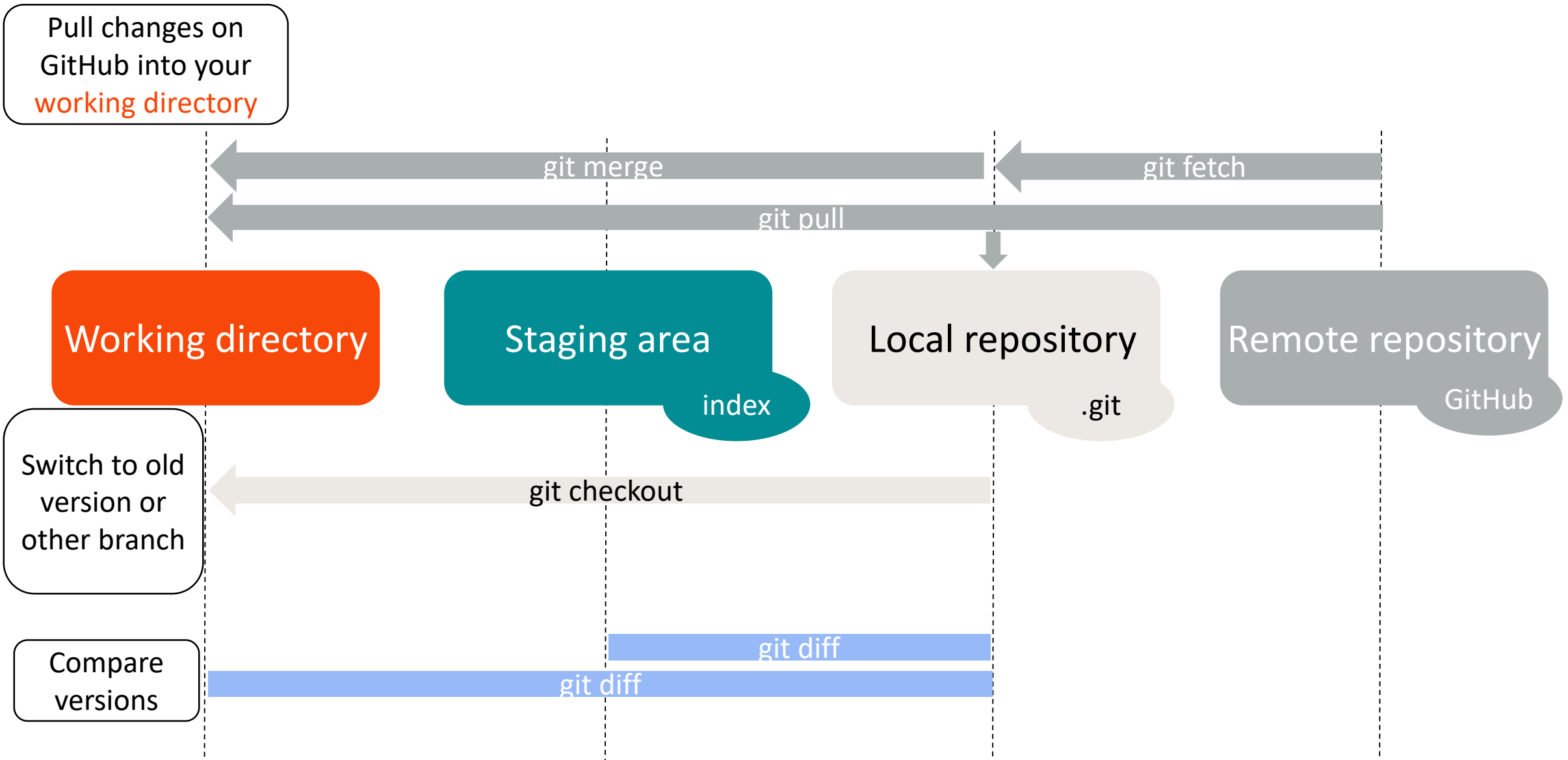
Our Git and Github Agenda



Standard Git Workflow

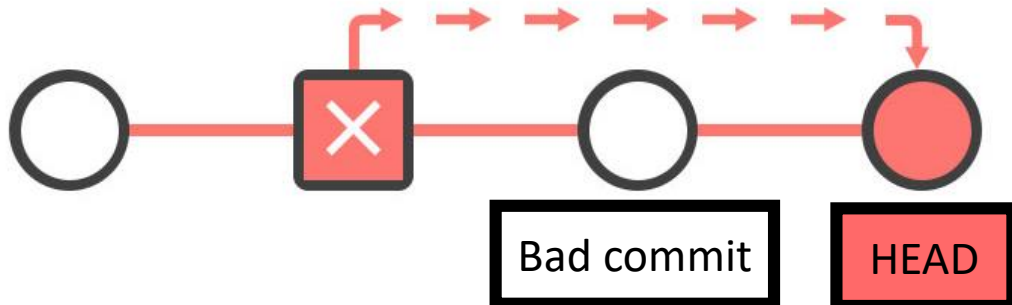


Standard Git Workflow



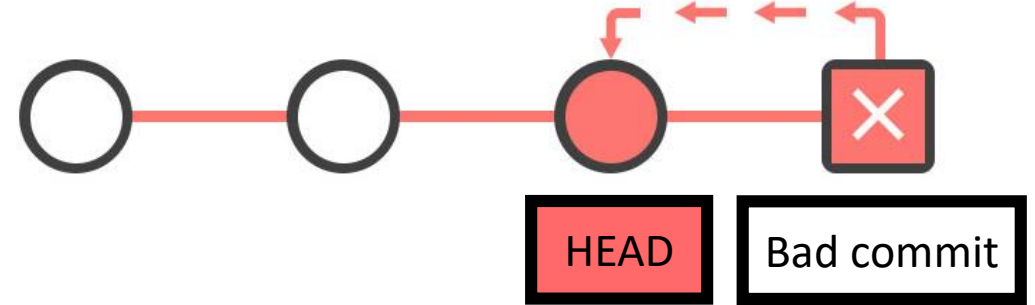
Revert (#*bad-commit*)

- Adds a new commit which reverts the „bad commit“
- Leaves history intact
- Can be safely done, even if you have already pushed to GitHub

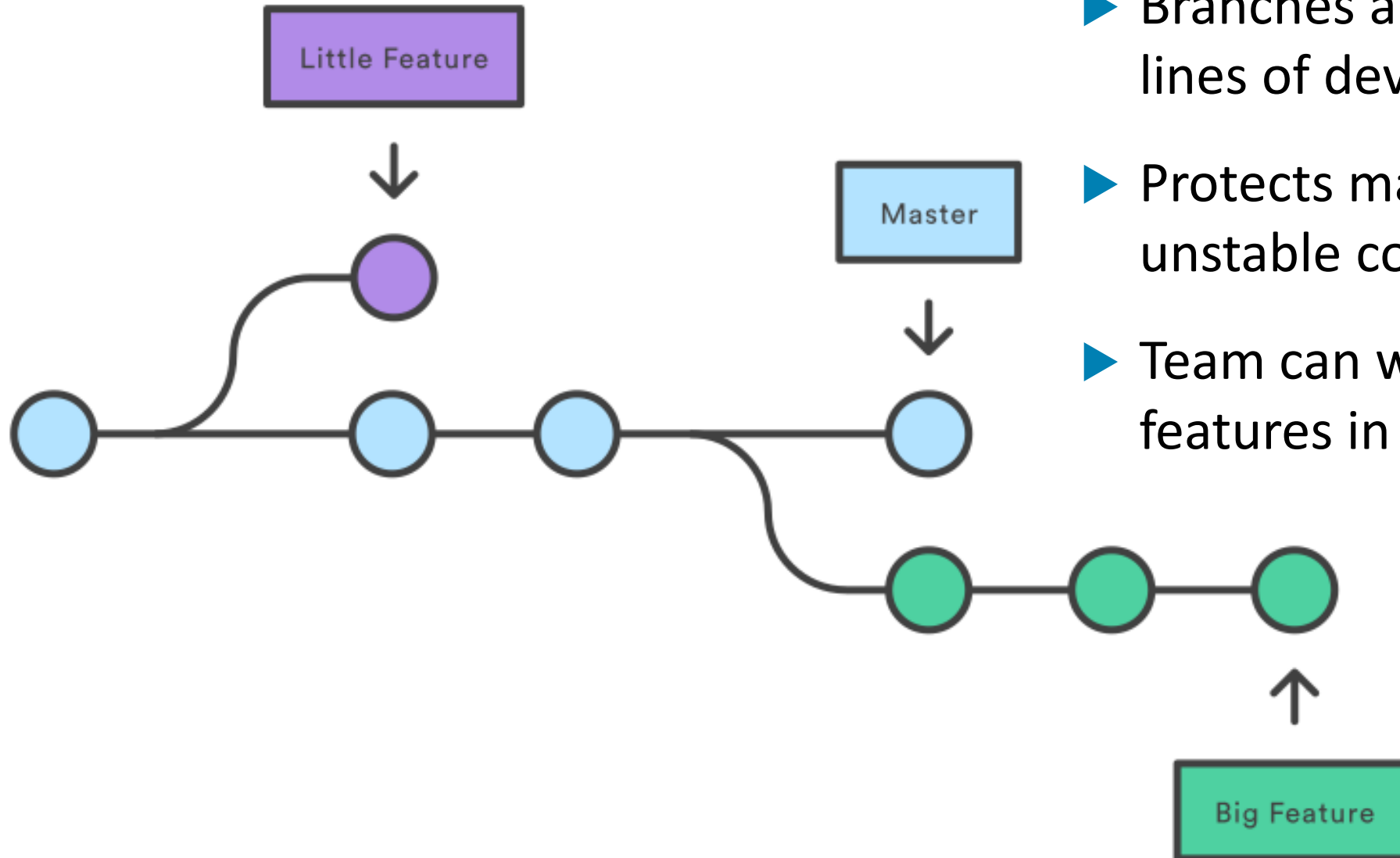


Reset (#*last-good-commit*)

- Undo the „bad commit“
- Rewrites history → use only for local changes!
- Versions:
 - reset hard: resets irreversibly
 - reset mixed (default): rewrites commit history, but keeps changes in working directory

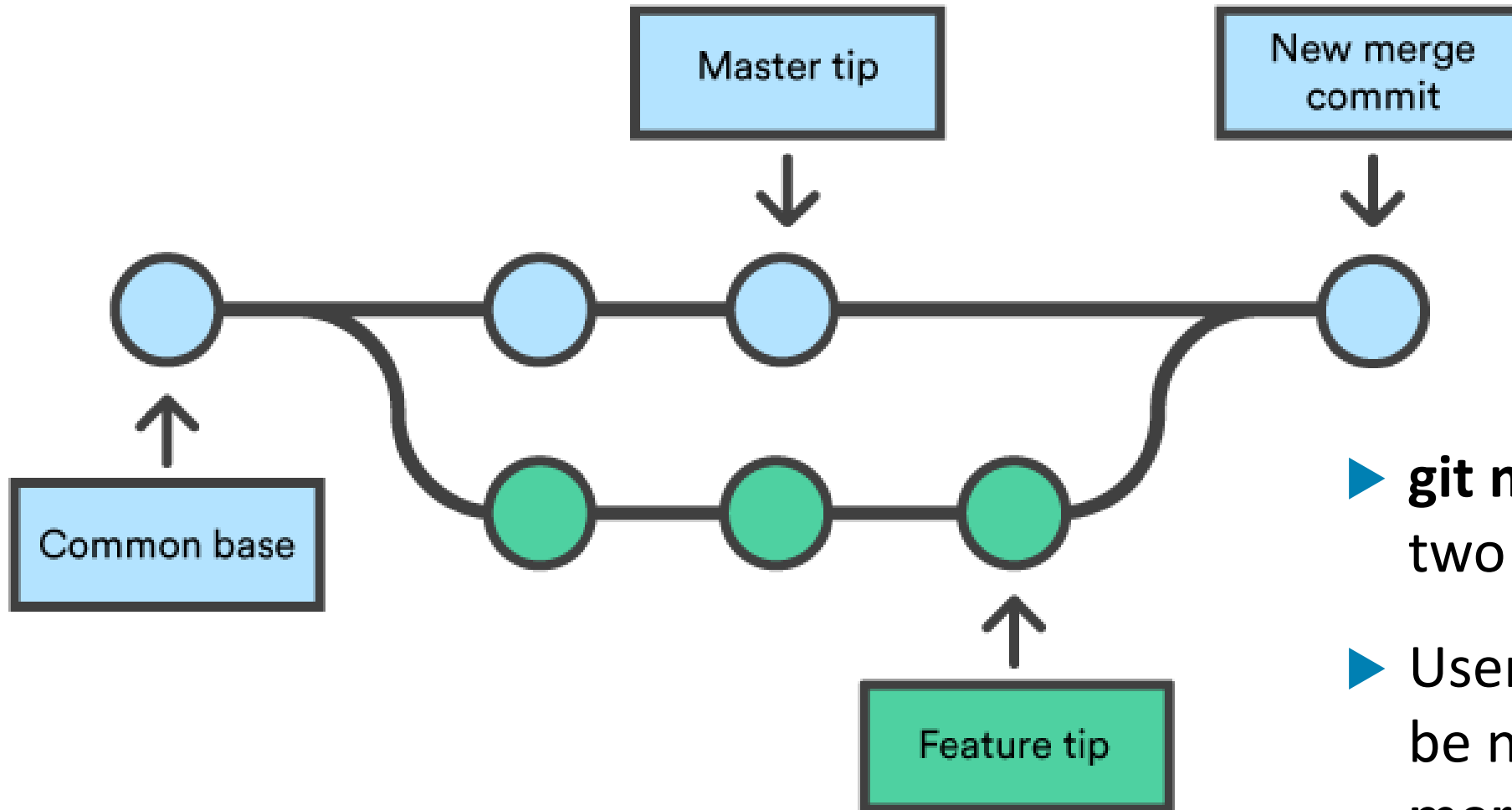


Branching



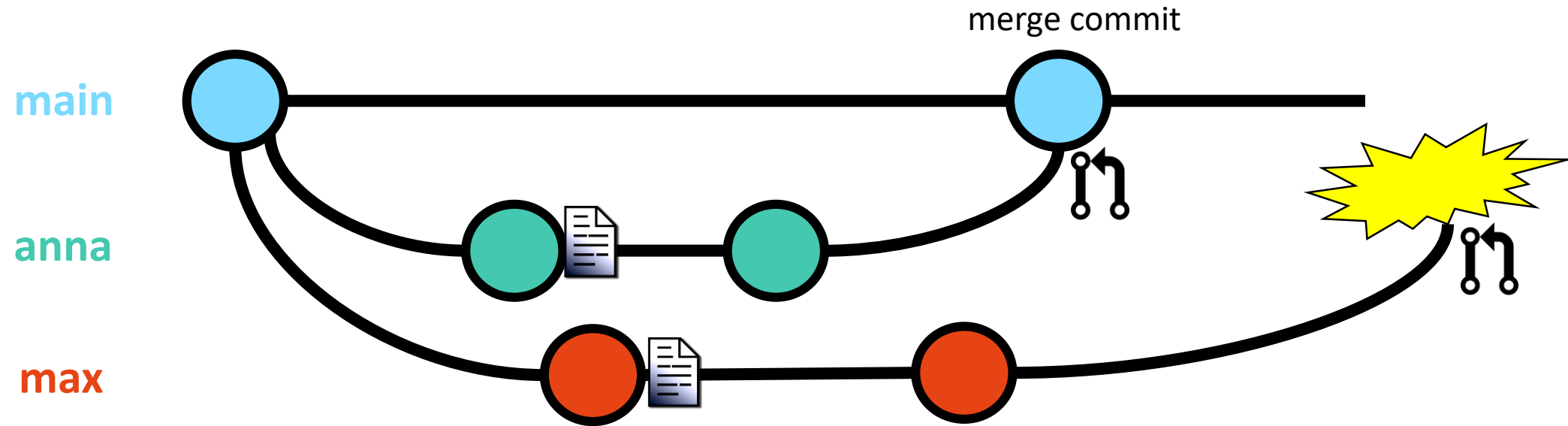
- ▶ Branches are independent lines of development
- ▶ Protects main branch against unstable code
- ▶ Team can work on separate features in parallel

Merging

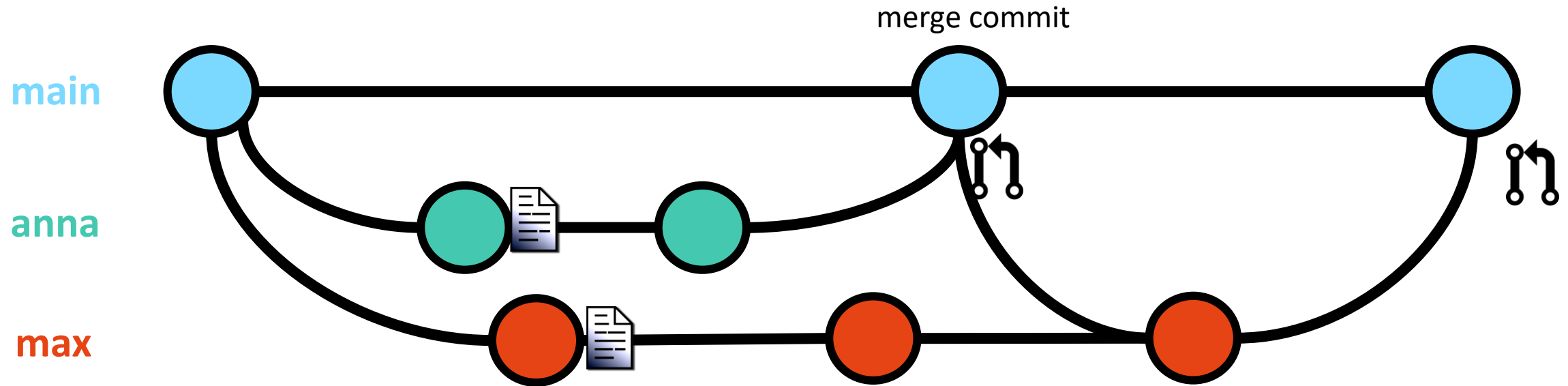


- ▶ **git merge** combines two branches into one.
- ▶ User intervention may be needed to resolve **merge conflicts**

Handling merge conflicts



Handling merge conflicts



Option 1: Before **max** opens a pull request, he checks for changes in **main** and merges them into his feature **branch** (dealing manually with any conflict)

Handling merge conflicts

```
git merge feature-x
```

```
Auto-merging analysis.py
CONFLICT (content): Merge conflict in analysis.py
Automatic merge failed; fix conflicts and then commit the result.
```

Accept Current Change | Accept Incoming Change | Accept Both Changes | Compare Changes

<<<<<< HEAD (Current Change)

df.tail()

=====

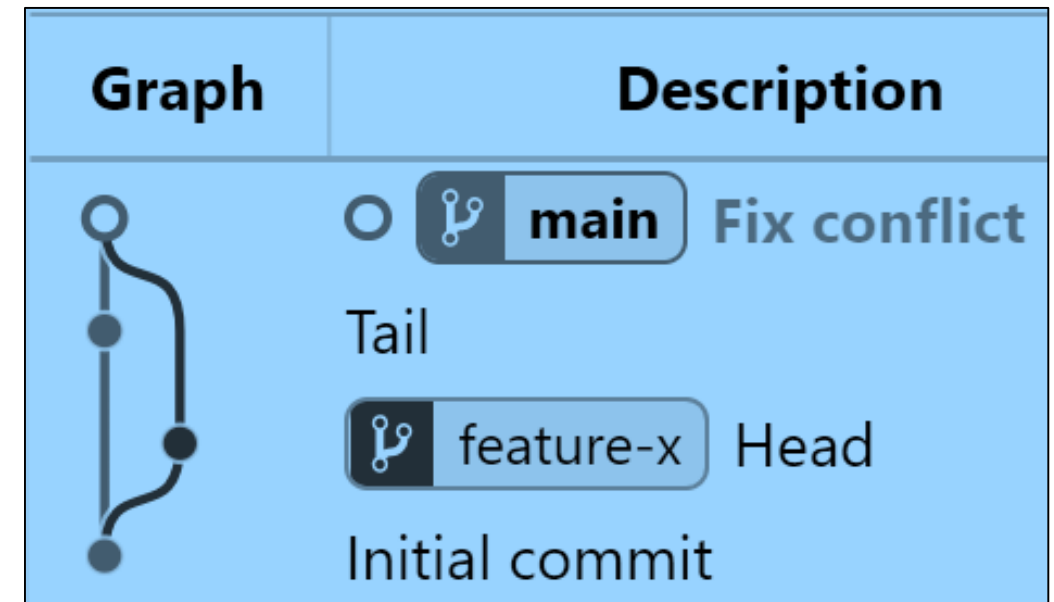
df.head()

>>>>>> feature-x (Incoming Change)

```
df.head()
```

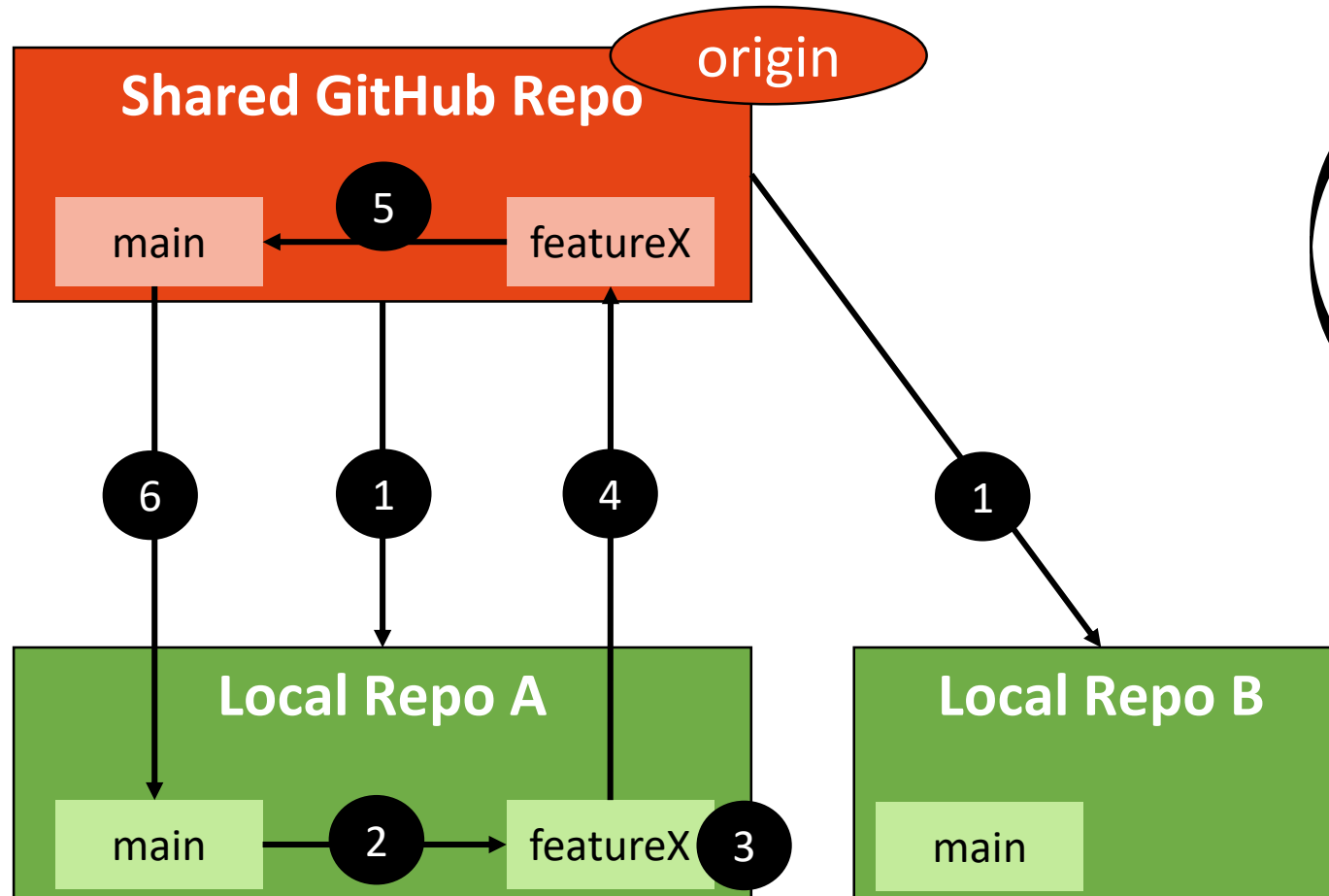
```
df.tail()
```

```
git commit -m "Fix conflict"
```



Feature Branch Workflow

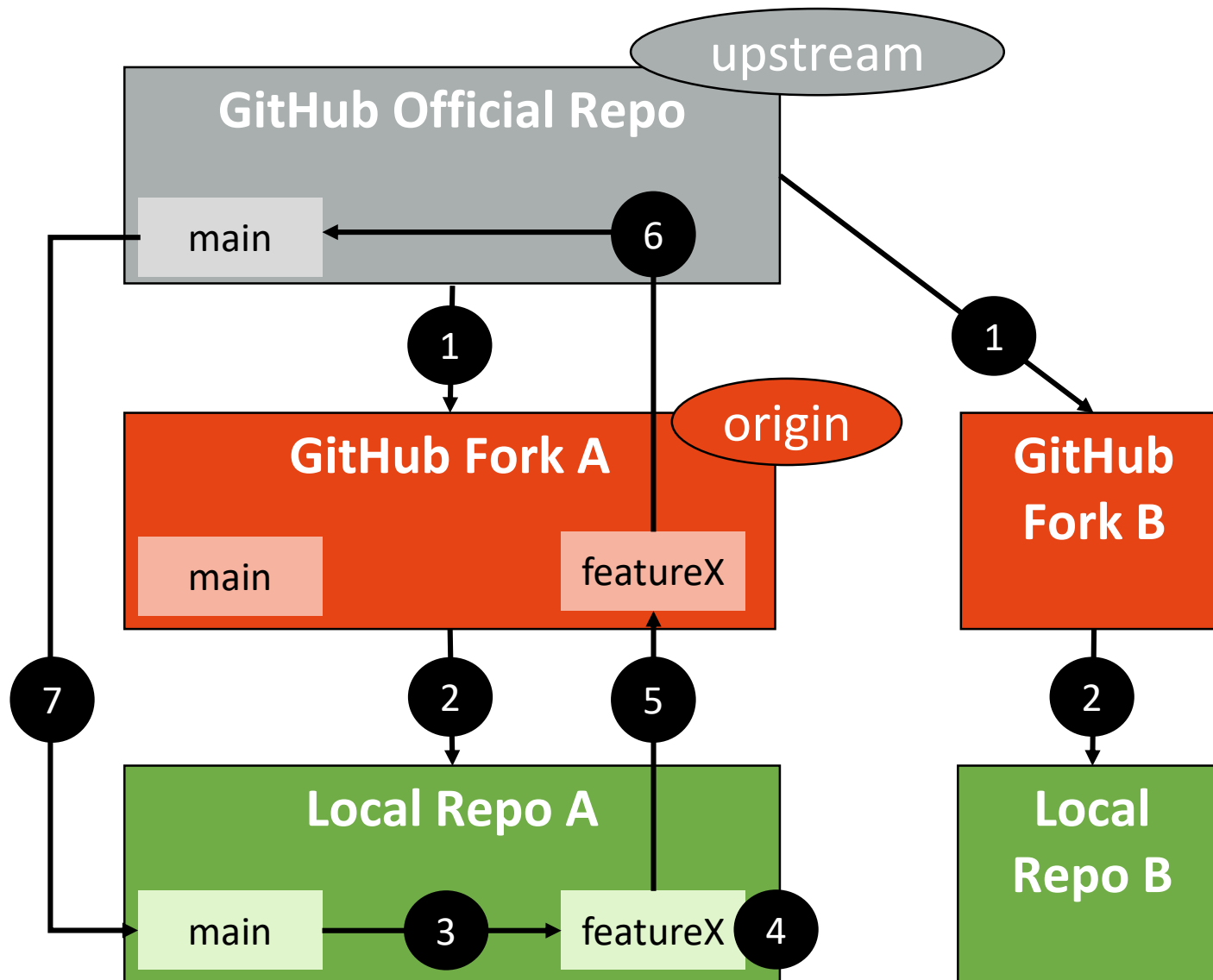
Company Setting



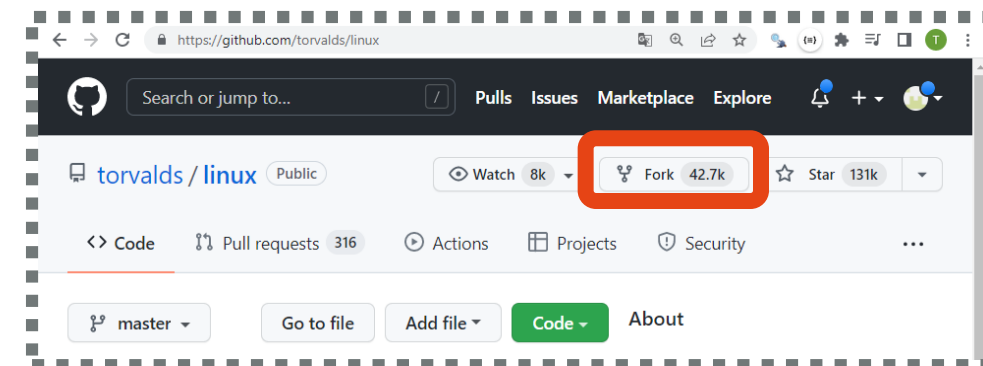
1. `git clone <url>`
2. `git branch featureX`
`git checkout featureX`
3. edit, `git add`, `git commit`
4. `git push`
5. Pull Request → `git merge`
6. `git pull`

Forking Workflow

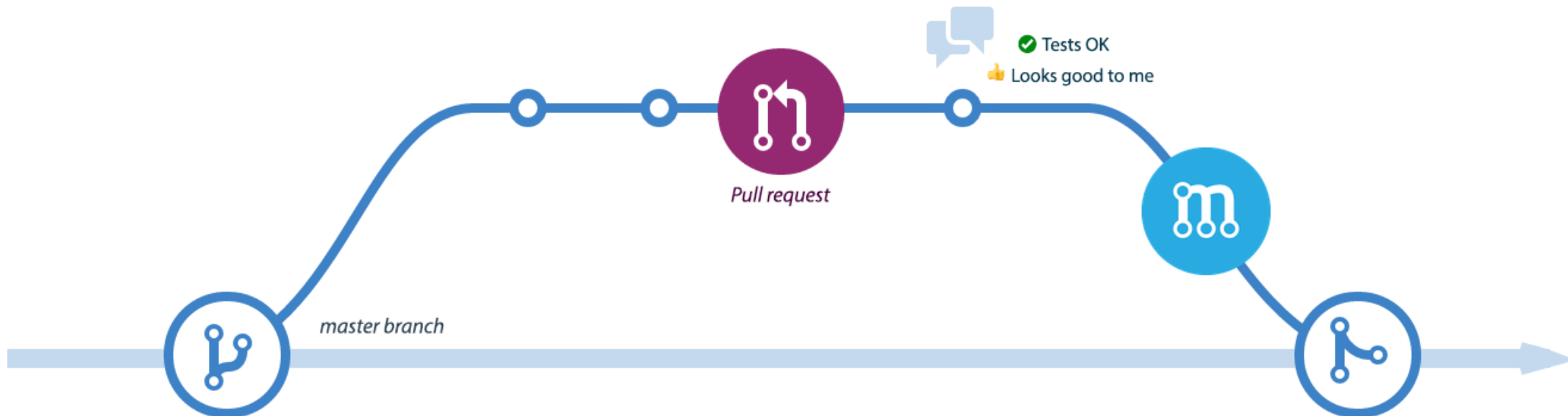
Open Source Setting



1. Fork
2. `git clone <url>`
3. `git branch featureX`
`git checkout featureX`
4. edit, `git add`, `git commit`
5. `git push`
6. Pull Request → `git merge`
7. (`git remote add upstream <url>`)
`git pull upstream`



- ▶ **A Pull request is a GitHub feature**, not a git command: you request to pull changes from your feature branch (“merge request” on Gitlab)
- ▶ May involve multiple iterations of discussions, code reviews, and follow-up commits, before the commit is merged into the main branch



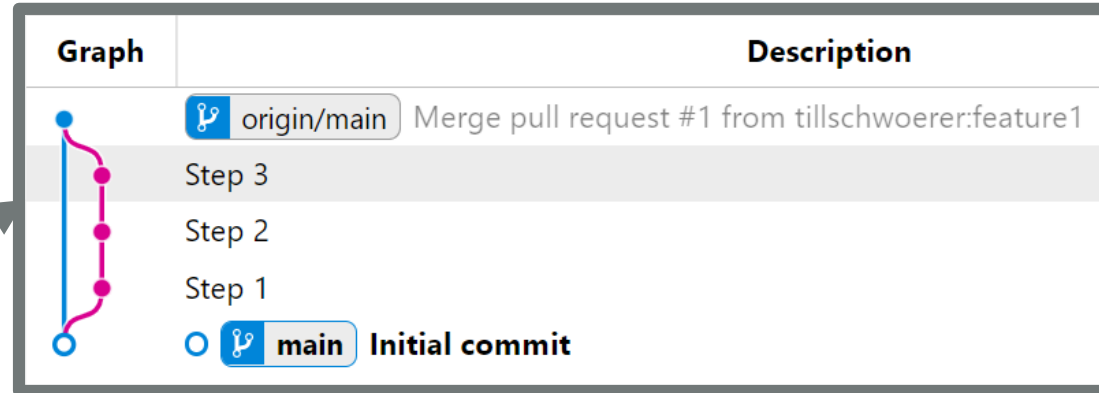
3 Ways of Merging in a Pull Request

Merge pull request You can also [open this](#)

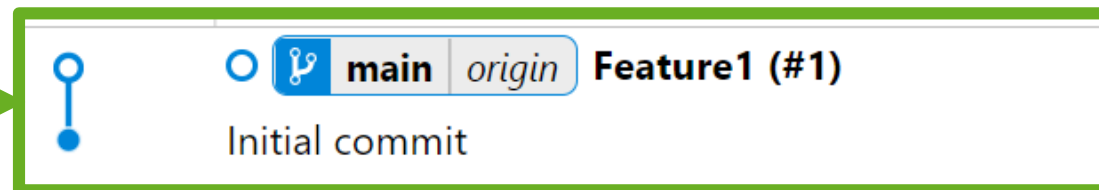
✓ **Create a merge commit**
All commits from this branch will be added to the base branch via a merge commit.

Squash and merge
The 3 commits from this branch will be combined into one commit in the base branch.

Rebase and merge
The 3 commits from this branch will be rebased and added to the base branch.



- + Truthful history of actions
- May pollute project history



- + Short and linear project history
- + Pull requests can still be identified
- Actual single commits are not logged anymore



- + Linear project history
- + Single commits are all part of the history
- Feature branches / pull requests cannot be identified any more

► **.gitignore:**

- ◆ define files that you don't want to track
- ◆ e.g. password files (.env), .ipynb_checkpoints, very large binary files

► **Large files**

- ◆ GitHub file size limit: 100 MB
- ◆ [Git Large File Storage](#): tracking files up to 2 GB
- ◆ Git LFS stores references to the file in the Github repository. The actual file is stored separately.

Further Resources

- ▶ [Tutorial](#) on Git (Workflows) (available in **German**)
- ▶ [Git reference book](#) (available in **German**)
- ▶ [Git commands: cheat sheet](#)
- ▶ [10-minute reads](#) mostly on GitHub topics (by GitHub)
- ▶ [Glossary of Git and GitHub terms](#) (by Github)