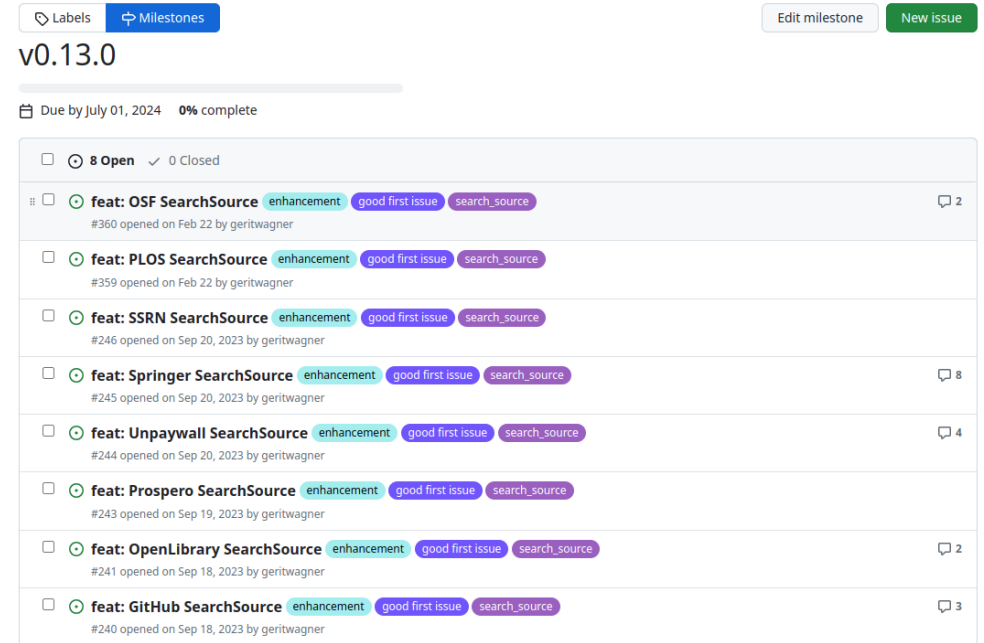




Introduction to Git

Check-in: Group formation

- Milestone
- Anyone not yet part of an issue discussion?
- Challenges related to the setup?



The screenshot shows a GitHub milestone page for version v0.13.0. At the top, there are tabs for 'Labels' and 'Milestones', with 'Milestones' selected. To the right are buttons for 'Edit milestone' and 'New issue'. Below the milestone title, a progress bar indicates '0% complete' and a due date of 'July 01, 2024'. A summary bar shows '8 Open' and '0 Closed' issues. The main content is a list of 8 issues, each with a title, labels, and a comment count. The issues are:

Issue Title	Labels	Comments
feat: OSF SearchSource	enhancement, good first issue, search_source	2
feat: PLOS SearchSource	enhancement, good first issue, search_source	0
feat: SSRN SearchSource	enhancement, good first issue, search_source	0
feat: Springer SearchSource	enhancement, good first issue, search_source	8
feat: Unpaywall SearchSource	enhancement, good first issue, search_source	4
feat: Prospero SearchSource	enhancement, good first issue, search_source	0
feat: OpenLibrary SearchSource	enhancement, good first issue, search_source	2
feat: GitHub SearchSource	enhancement, good first issue, search_source	3

Tip! You can use `shift + j` or `shift + k` to move items with your keyboard.

Git: A distributed version control system

Advantages:

- Every repository has a full version history
- Most operations run locally
- Reliable data handling, ensuring integrity and availability
- Efficient data management for versions and branches
- Scalable collaboration mechanisms for large teams and complex projects

Caveats:

- Need to learn and understand the underlying model
- Not built for binary files or large media files



Learning objectives

Understand and use git to develop software in teams.

Part 1: Branching

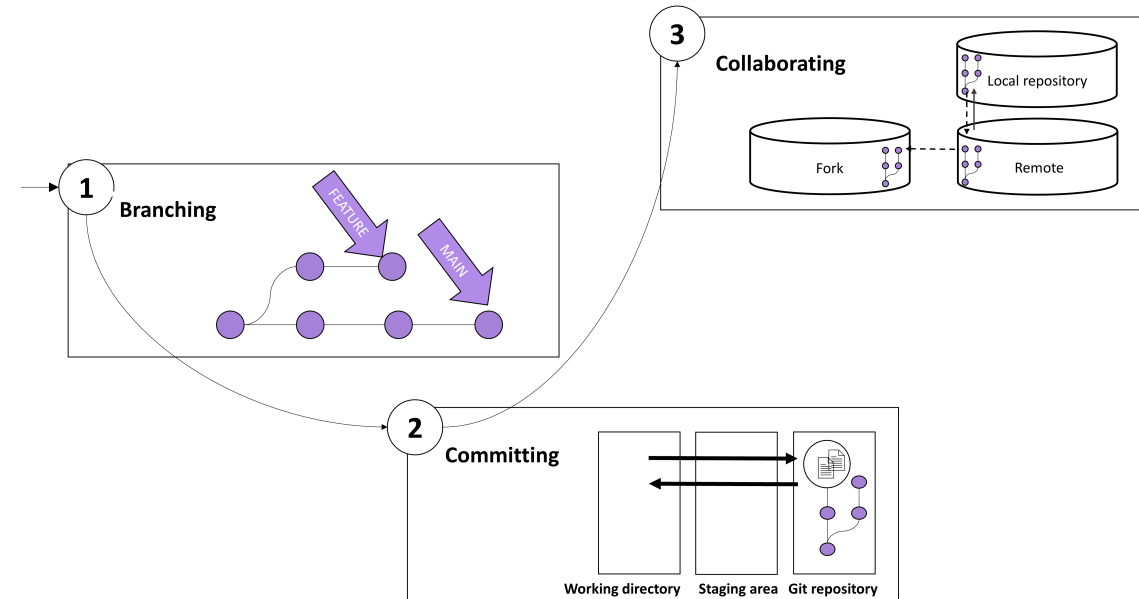
Part 2: Committing

Part 3: Collaborating

Each part starts with the **concepts** before the **practice** session.

In the practice sessions:

- Form groups of two to three students
- Work through the exercises
- Create a *cheat sheet* summarizing the key commands



* Note: This session is based on our [unique and peer-reviewed approach](#).

Start the Codespace

Open the notebook for practicing Git branching:



Open in GitHub Codespaces

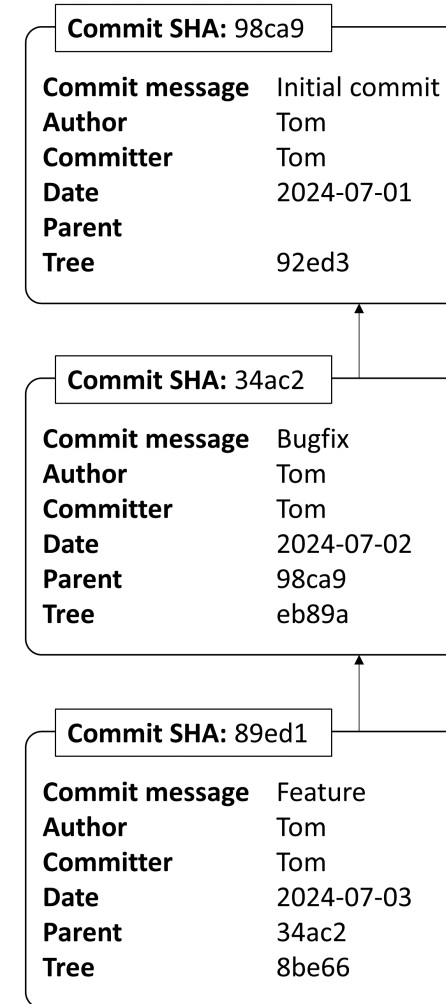
The setup can run in the background, while we focus on the concepts.



Part 1: Branching

Commits

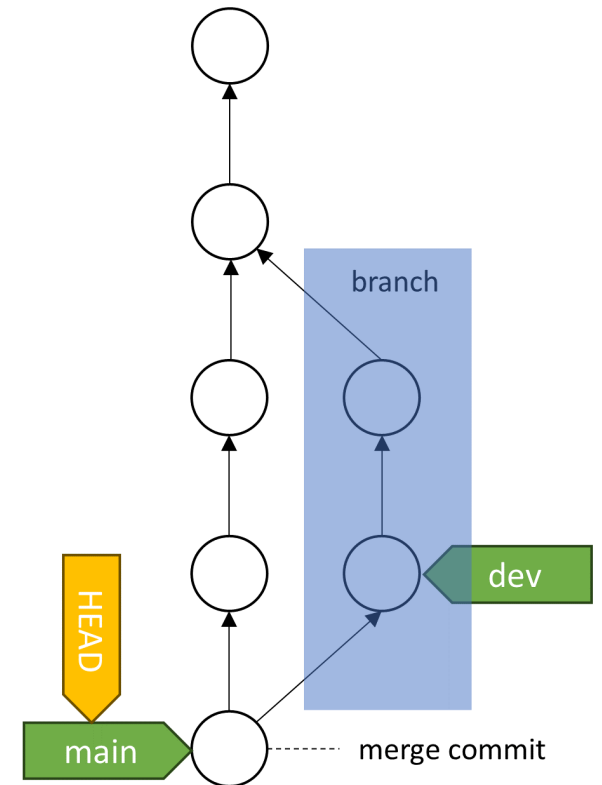
- A **commit** refers to a snapshot (version) of the whole project directory, including the metadata and files
- Commits are identified by the **SHA** fingerprint of their metadata and content*, e.g., 98ca9
- Commits are created in a sequence, with every commit pointing to its **parent** commit(s)
- The **tree** object contains all files (and non-empty directories); it is identified by a SHA hash
- Commits are created by the **git commit** command



* If any of the metadata or content changes, the SHA will be completely different.

The DAG, branches, and HEAD

- Commits form a **directed acyclic Graph (DAG)**, i.e., all commits can have one or more children, and one or more parents (except for the first commit, which has no parent). Closed directed cycles are not allowed.
- With the **git branch <branch-name>** command, a separate line of commits can be started, i.e., one where different lines of commits are developed from the same parent. The branch pointer typically points at the latest commit in the line.
- With the **git switch <branch-name>** command, we can select the branch on which we want to work. Switch effectively moves the HEAD pointer, which points to a particular branch and indicates where new commits are added.
- With the **git merge <other-branch>** command, separate lines of commits can be brought together, i.e., creating a commit with two parents. The *merge commit* integrates the contents from the <other-branch> into the branch that is currently selected. The <other-branch> is not changed.
- Per default, Git sets up a branch named "main".



Note: Arrows point from children to parent commits.



Practice: Branching

Open the notebook for practicing Git branching:



Open in GitHub Codespaces



Part 2: Committing

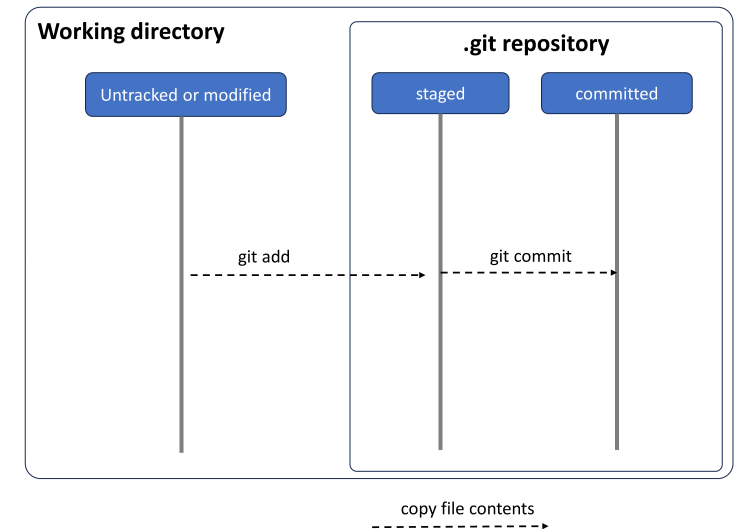
The working directory and .git repository

All working file contents reside in the working directory; staged and committed file contents are stored in the `.git` directory (a subfolder of the working directory).

Git allows us to stage (select) specific file contents for the next commit.

- With **git add <file-name>**, contents of an *untracked or modified* file are copied to the `.git` repository and added to the staging area, i.e., explicitly marked for inclusion in the next commit.
- With **git commit**, *staged* files contents are included in a *commit*.

The **git init** command creates the `.git` directory.

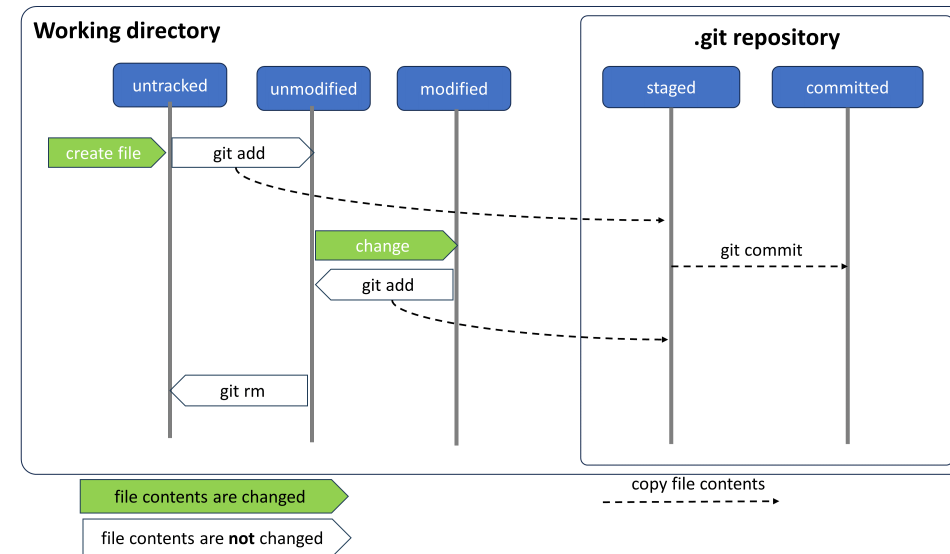


The three states of a file

Files in the working directory can reside in three states:

- New files are initially **untracked**, i.e., Git does not include new files in commits without explicit instruction.
- With *git add*, file contents are staged, and the file is tracked. Given that the file in the working directory is identical with the staged file contents, the file is **unmodified**.
- When users change a file, it becomes **modified**, i.e., the file in the working directory differs from the file contents in the staging area.
- With *git add*, the file contents are staged again, and the file becomes **unmodified**.
- With *git rm*, files are no longer tracked.

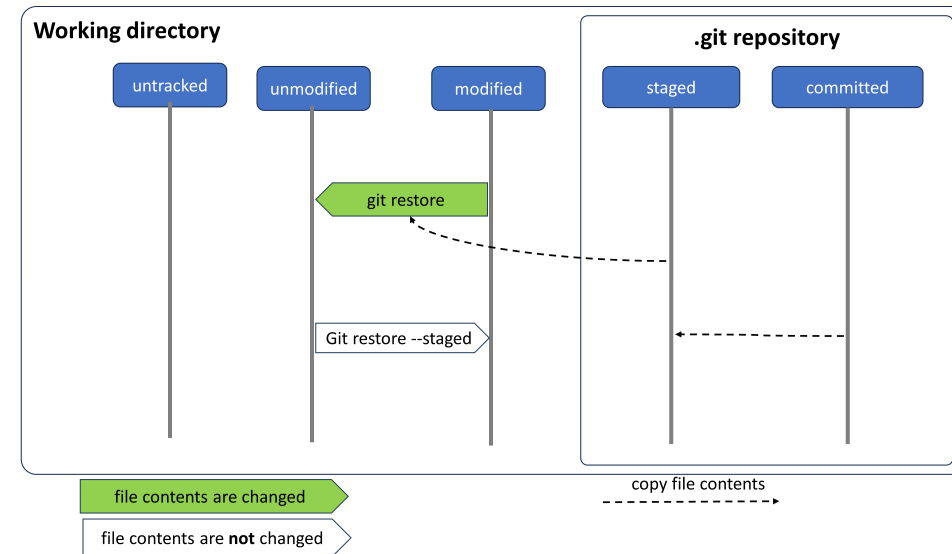
Note: *git add* and *git rm* do not change the contents of the file in the working directory.



Resetting changes

To undo changes that are not yet committed, it is important to understand whether they are staged or unstaged:

- If changes are not yet staged, the file is currently *modified*. A **git restore <file-name>** replaces the file in the working directory with the staged version. As a result, the file is *unmodified* because it corresponds to the staged file.
- If the file is currently *unmodified*, a **git restore --staged <file-name>**, Git discards the staged changes by using the last committed version. The file contents in the working directory do not change, but the file becomes *modified* because it differs from the staged version.



Practice: Committing

Open the notebook for practicing Git committing:



Open in GitHub Codespaces

Transfer challenges I

1. Consider how the **git switch** (or the revert/pull/checkout) command affects the git areas. How does it affect the working directory?
2. Git provides the option to edit prior commits using an interactive rebase, such as the **git rebase -i**. How would that affect the following commits?

Transfer challenge: Git merge conflicts

Open the notebook for practicing the resolution of Git merge conflicts (related to branching and committing):



Open in GitHub Codespaces



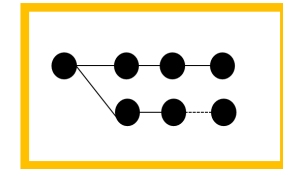
Part 3: Collaborating

Collaborating

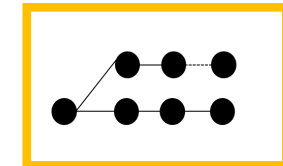
- The distributed model of Git means that every repository has a full version history, (almost) all operations can be completed locally, and every repository can be developed autonomously.
- To collaborate, a *remote* repository is needed, initially named "origin"
- If the remote repository exists, the **git clone** command retrieves a local copy
- To create a remote repository (named "origin"), and push a specific branch:

```
git remote add origin REMOTE-URL  
git push origin main
```

remote „origin“ repository



clone ↓ ↑ push



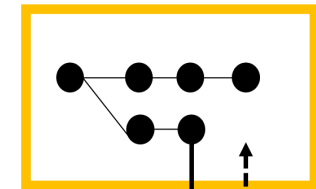
local repository

* If the remote repository does not exist, you have to add the remote origin and push the repository

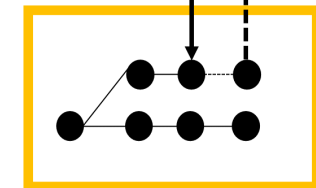
Collaborating on branches

- To retrieve changes, use the **git pull** command
- To share changes, use the **git push** command
- Most remote operations, including pull, push, pull requests refer to branches
- In some cases, **branches must be selected explicitly**, and in other cases, git automatically selects branches, i.e., it remembers the typical branch to pull or push

remote „origin“ repository



pull

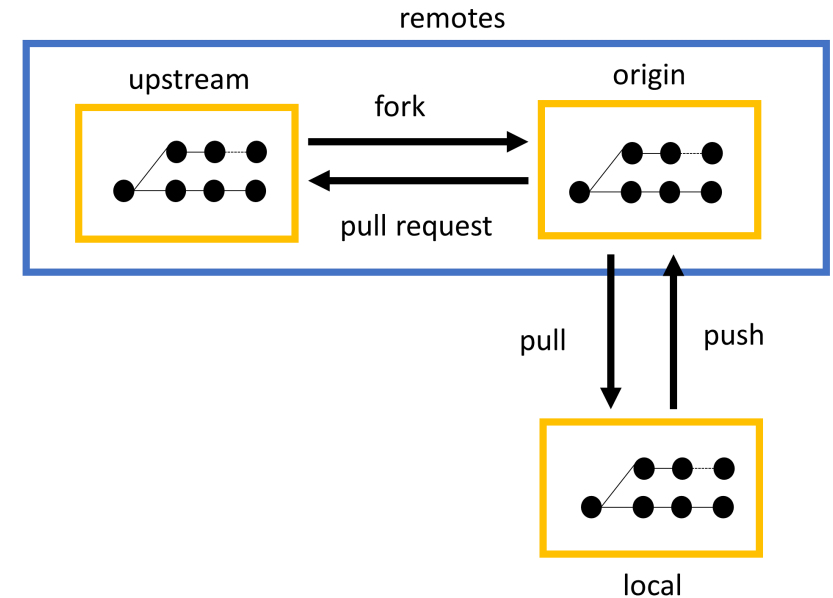


local repository

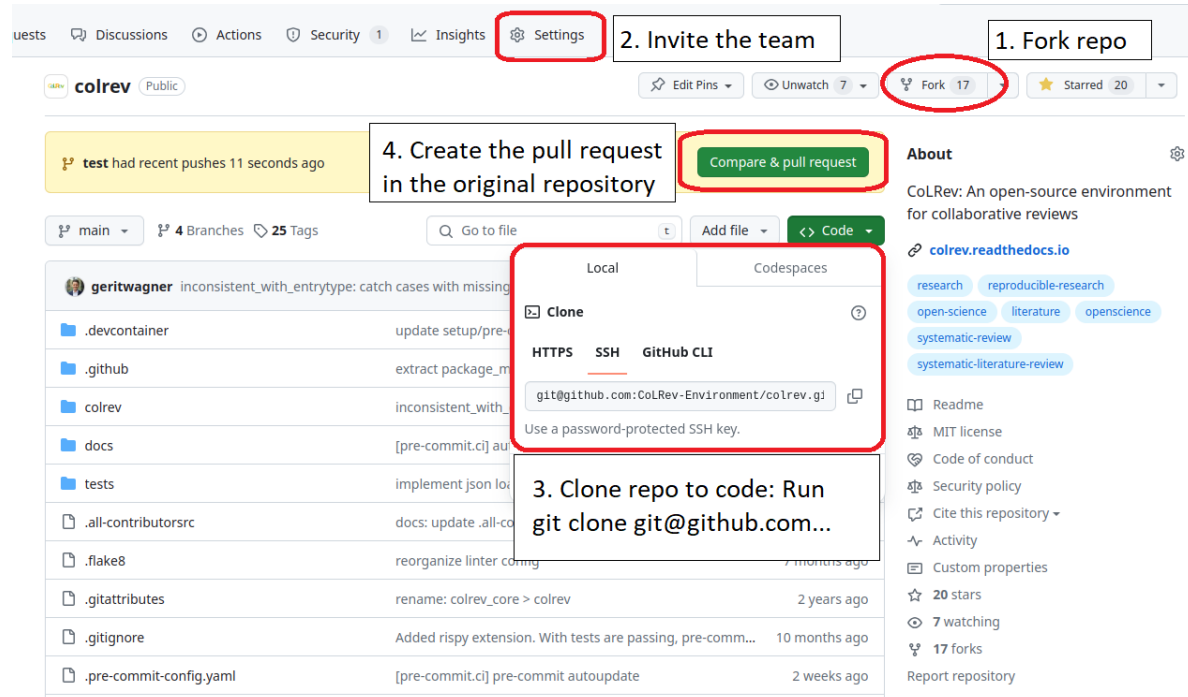
Collaborating with forks

This model works if you are a maintainer of the remote/origin, i.e., if you have write access.

- In Open-Source projects, write-access is restricted to a few maintainers
- At the same time, it should be possible to integrate contributions from the community
- **Forks** are remote copies of the upstream repository
- Contributors can create forks at any time and push changes
- Contributors can open a **pull request** to signal to maintainers that code from the fork can be merged
- Pull requests are used for code review, and improvements before code is accepted or rejected



Fork, invite, clone, and pull-request on GitHub



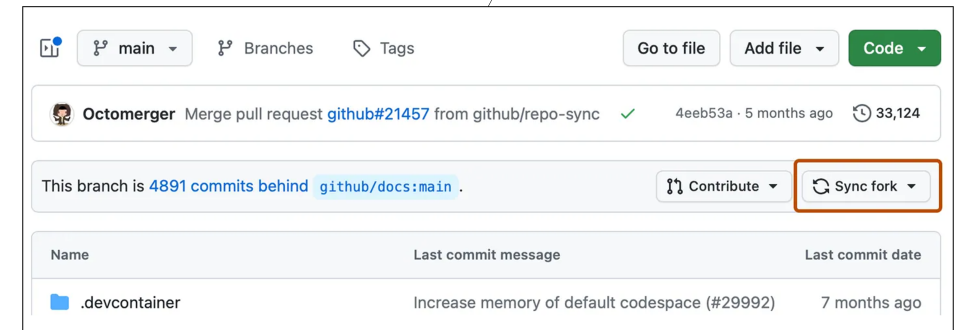
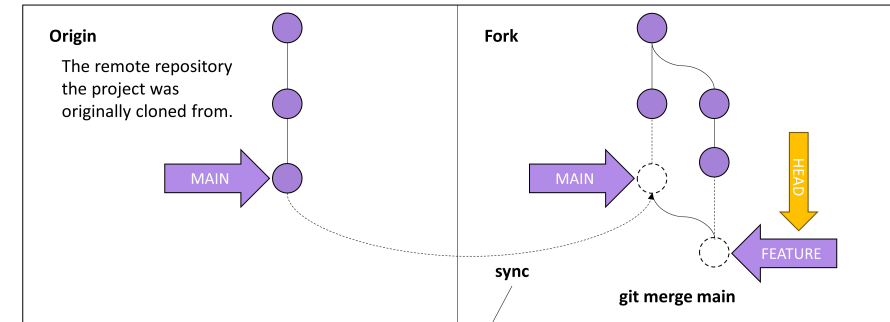
The screenshot shows the GitHub repository page for 'colrev'. Annotations include:

- 1. Fork repo**: A red circle highlights the 'Fork 17' button in the top right.
- 2. Invite the team**: A red box highlights the 'Settings' tab in the top navigation bar.
- 3. Clone repo to code: Run git clone git@github.com...**: A red box highlights the 'Clone' dialog box, which shows the SSH URL: `git@github.com:CoLRev-Environment/colrev.git`.
- 4. Create the pull request in the original repository**: A red box highlights the 'Compare & pull request' button.

The repository page also shows a file tree with folders like `.devcontainer`, `.github`, `colrev`, `docs`, and `tests`, and a list of recent commits.

Work in a forked repository

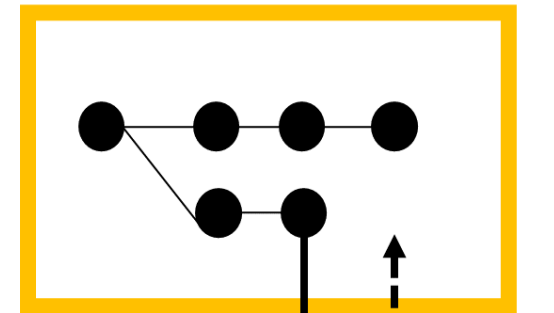
- In the fork, it is recommended to create working branches instead of committing to the `main` branch.
- It is good practice to regularly **sync** the `main` branches (on GitHub), and merge the changes into your working branches (locally or on GitHub).
- Syncing changes may be necessary to get bug fixes from the original repository, and to prevent diverging histories (potential merge conflicts in the pull request).



Remotes and branches

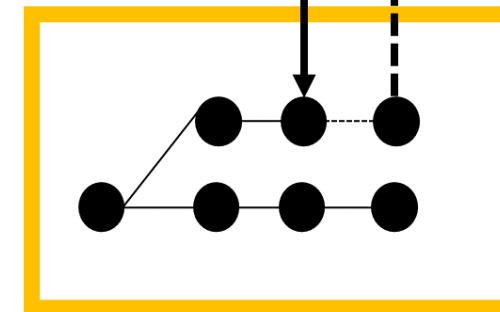
- Most remote operations, including pull, push, pull requests refer to branches
- In some cases, **branches must be selected explicitly**: pull requests, or pulling new branches
- In other cases, git automatically selects branches, i.e., it remembers the typical branch to pull or push

remote „origin“ repository



pull

push



local repository

Transfer challenges II

- Once a pull request has been opened, how can new changes (commits) be added?
- Assume that you discovered a typo in a very old commit. One option would be to run an interactive rebase and fix the typo. Why could such cases of "rewriting history" be problematic in collaborative settings?
- When pulling changes, there are two strategies to handle diverging branches: `--merge` or `--rebase`. How do the results differ between these strategies?
- GitHub offers the possibility to edit files directly. Are all three git areas available in this setting?

Which branching / merging strategy should we select?

Recommended branch setup in your fork:

1. Work on a shared **feature branch**, such as `unpaywall_search` . This is where your latest, working version is developed
2. Do not commit directly to `remotes/fork/main` . This branch should be kept in-sync with `remotes/origin/main`
3. Regularly merge `remotes/origin/main` into `remotes/fork/main` and `remotes/fork/main` into your feature branch using merge commits (i.e., `sync`, which will fast-forward, `git fetch` , `git switch feature_branch` and `git merge main`)



Survey

Please share your feedback to help us improve!

Project organization

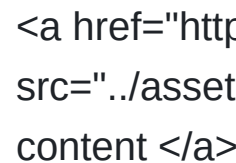
- Select a team leader who creates the fork and invites group members
- Plan how tasks could be completed in separate branches
- Avoid working on the `main` branch and synchronize it regularly with the original repository
- Regularly check whether branches should be synchronized (merged)

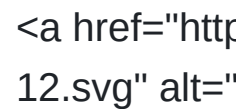
Task: complete one merge between branches.

Note: we will distribute a survey asking for the current state of your project after the merge. Your input will help us prepare the best practice session.






We value your feedback and suggestions

We encourage you to share your feedback and suggestions on this slide deck:

[Suggest specific changes by directly modifying the content](https://github.com/digital-work-lab/open-source-project/edit/main/slides/02-git.md) 

[Provide feedback by submitting an issue](https://github.com/digital-work-lab/open-source-project/issues/new) 

Your feedback plays a crucial role in helping us align with our core goals of **impact in research, teaching, and practice**. By contributing your suggestions, you help us further our commitment to **rigor, openness** and **participation**. Together, we can continuously enhance our work by contributing to **continuous learning** and collaboration across our community.

Visit this [page](https://digital-work-lab.github.io/handbook/docs/10-lab/10_processes/10.01.goals.html) to learn more about our goals:      .