



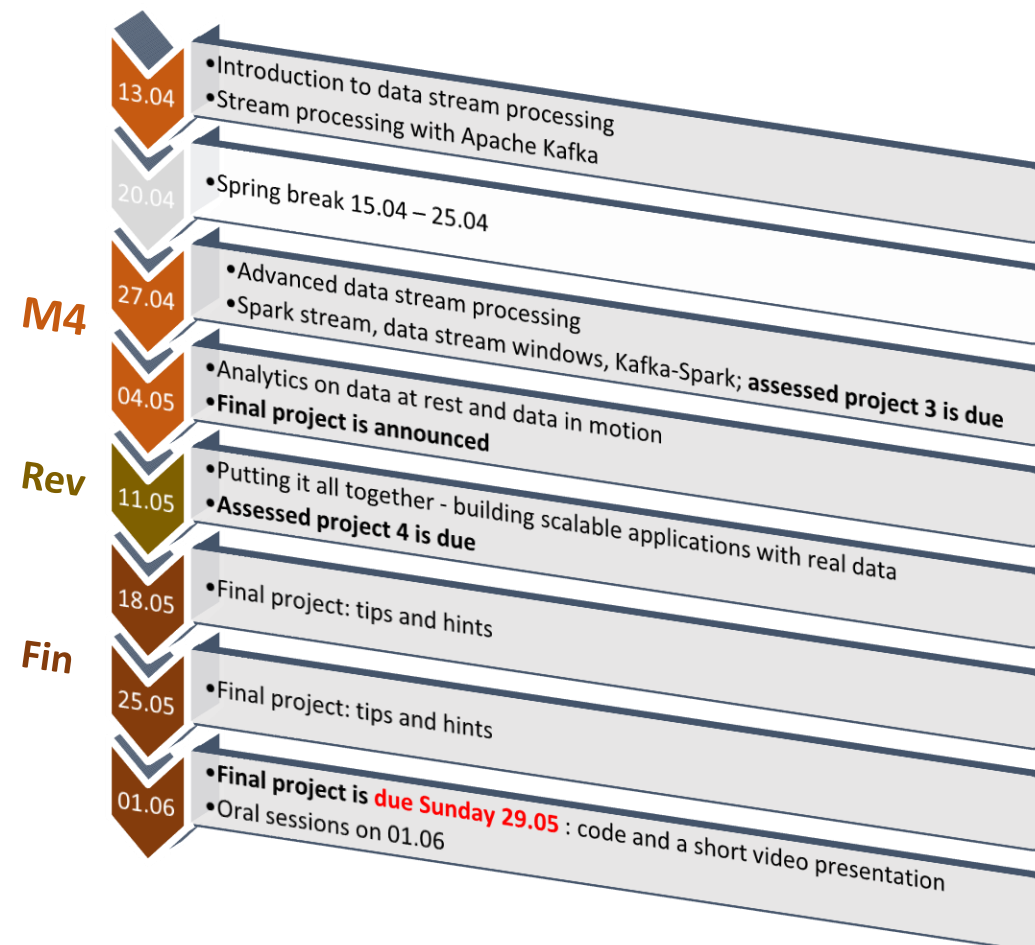
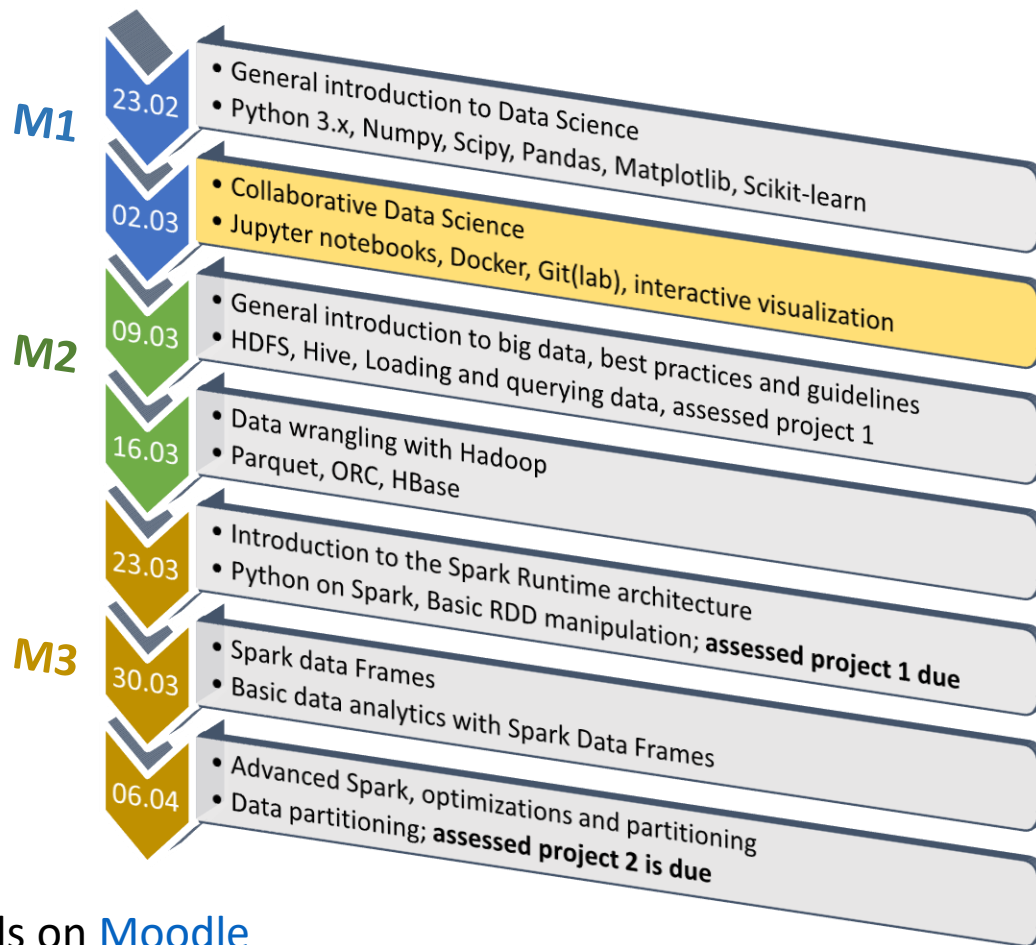
THE DATA SCIENCE LAB

Elements of Collaborative Data Science

COM 490 – Spring 2022

Week 2

Agenda Spring 2022

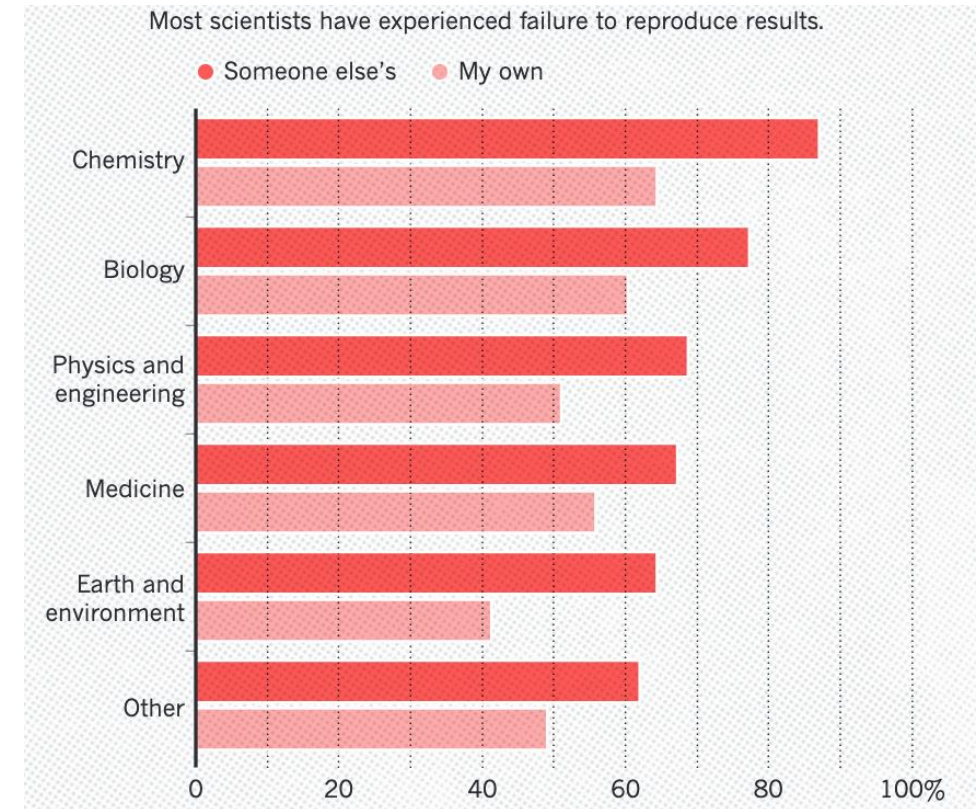


*Details on [Moodle](#)

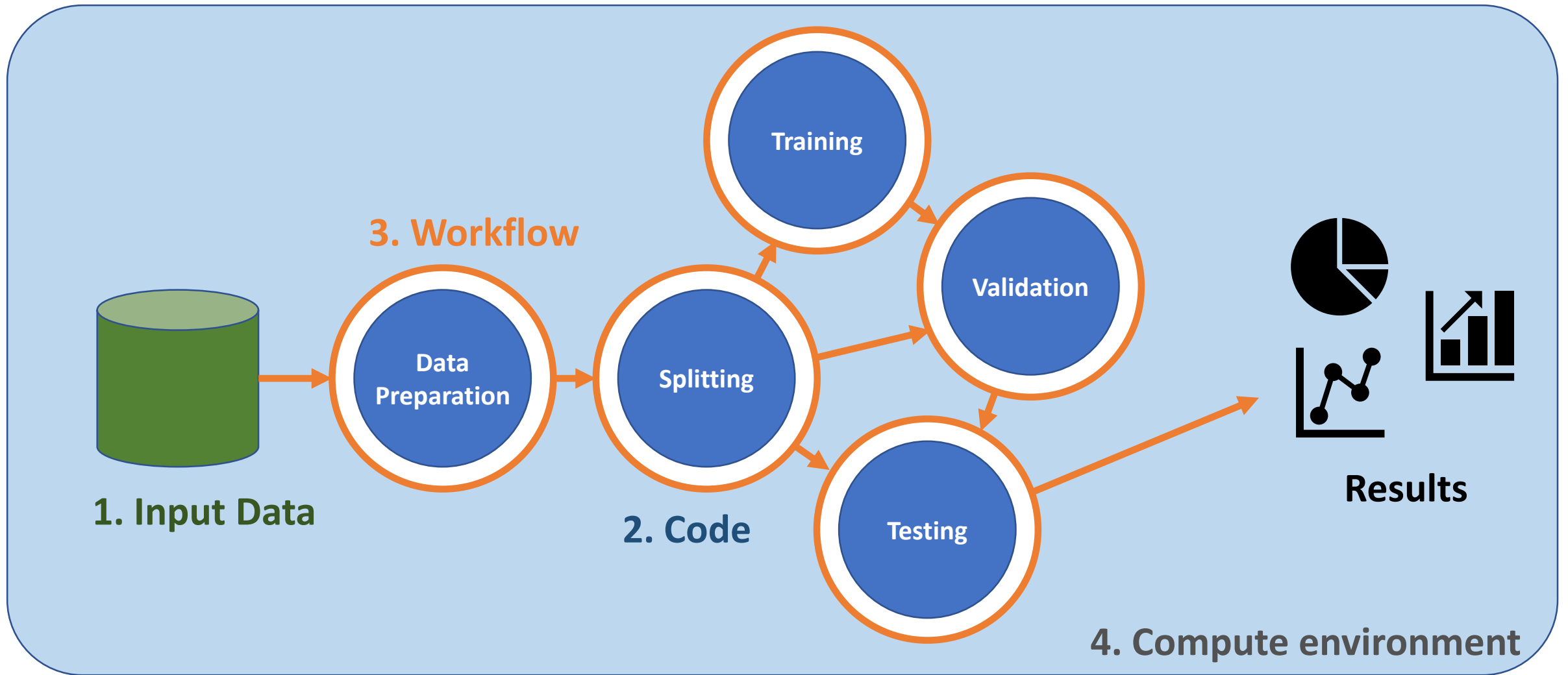
Week 1 – Questions?

Forewords – Scientific Reproducibility “Crisis”

- It is real and cannot be ignored
 - [1,500 scientists lift the lid on reproducibility](#)
(Nature, Aug. 2016)



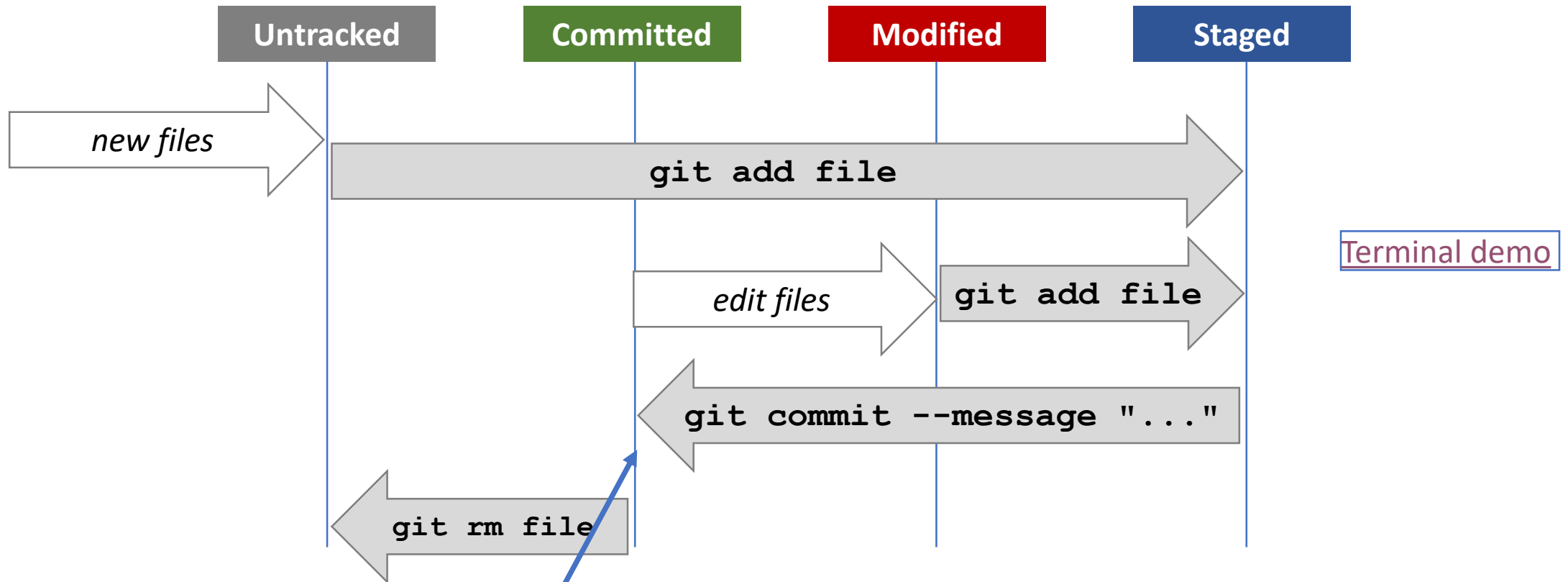
Ingredient of Collaborative Data Science



Collaborative Data Science - Tools

- **Code versioning**
 - Version Control Software: [git](#), github or gitlab
 - *Alternatives: SVN, Mercurial, Bazaar, BitKeeper, CVS, RCS, ...*
- **Data versioning**
 - Large file versioning: git extension for large file systems (git lfs)
- **Compute environment**
 - Runtime environment
 - [Docker](#) images and containers
 - *Alternatives: LXC, VMWare, VirtualBox, ...*
 - Package managers (Python)
 - [conda](#), pip
- **Workflow executor (not part of this course)**
 - [Snakemake](#) for Python, Apache [Oozie](#) for Hadoop
 - *Alternatives: make (and clones), nextflow, Cromwell/Toil (CWL), ...*

Git – File Lifecycle

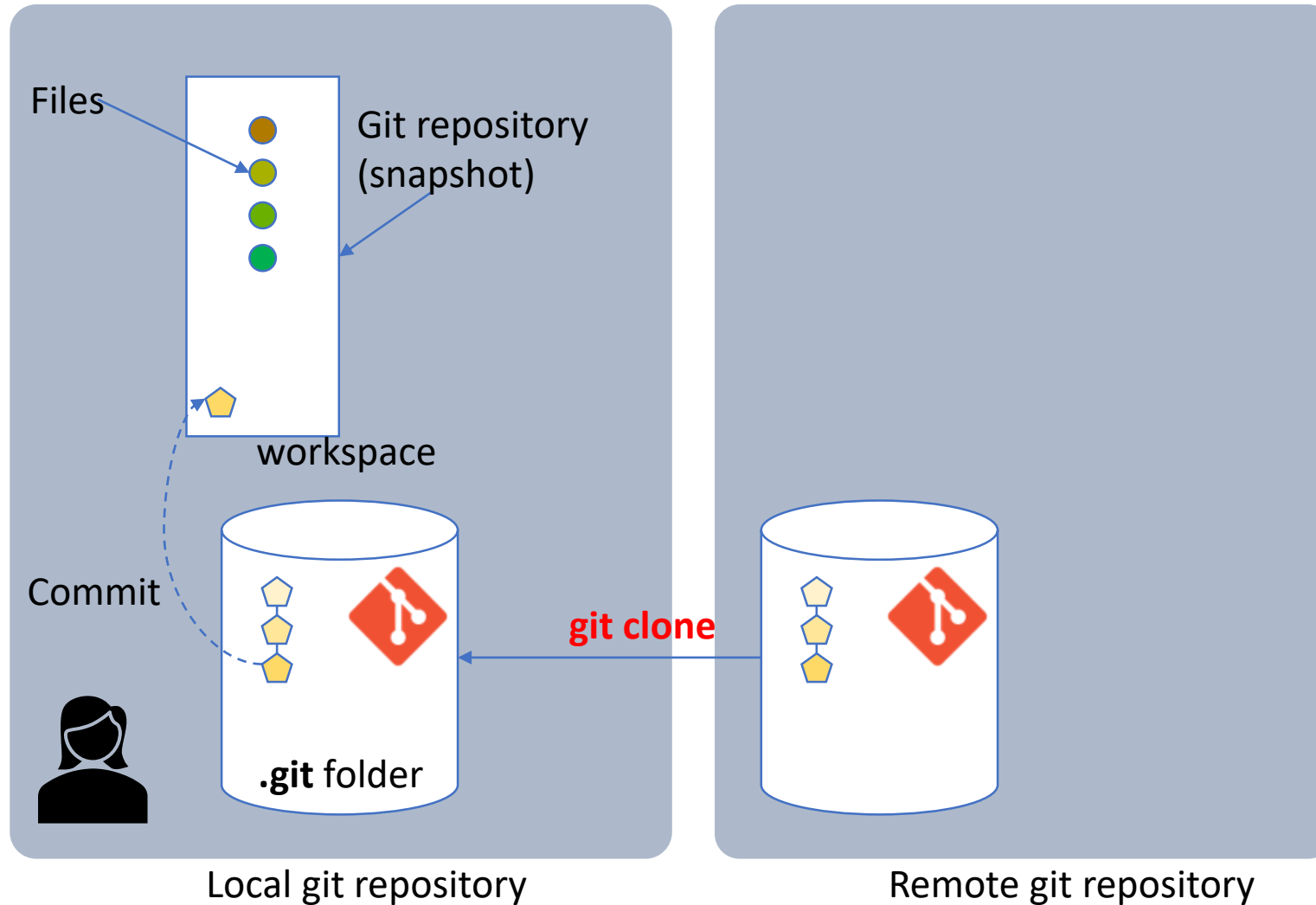


[Terminal demo](#)

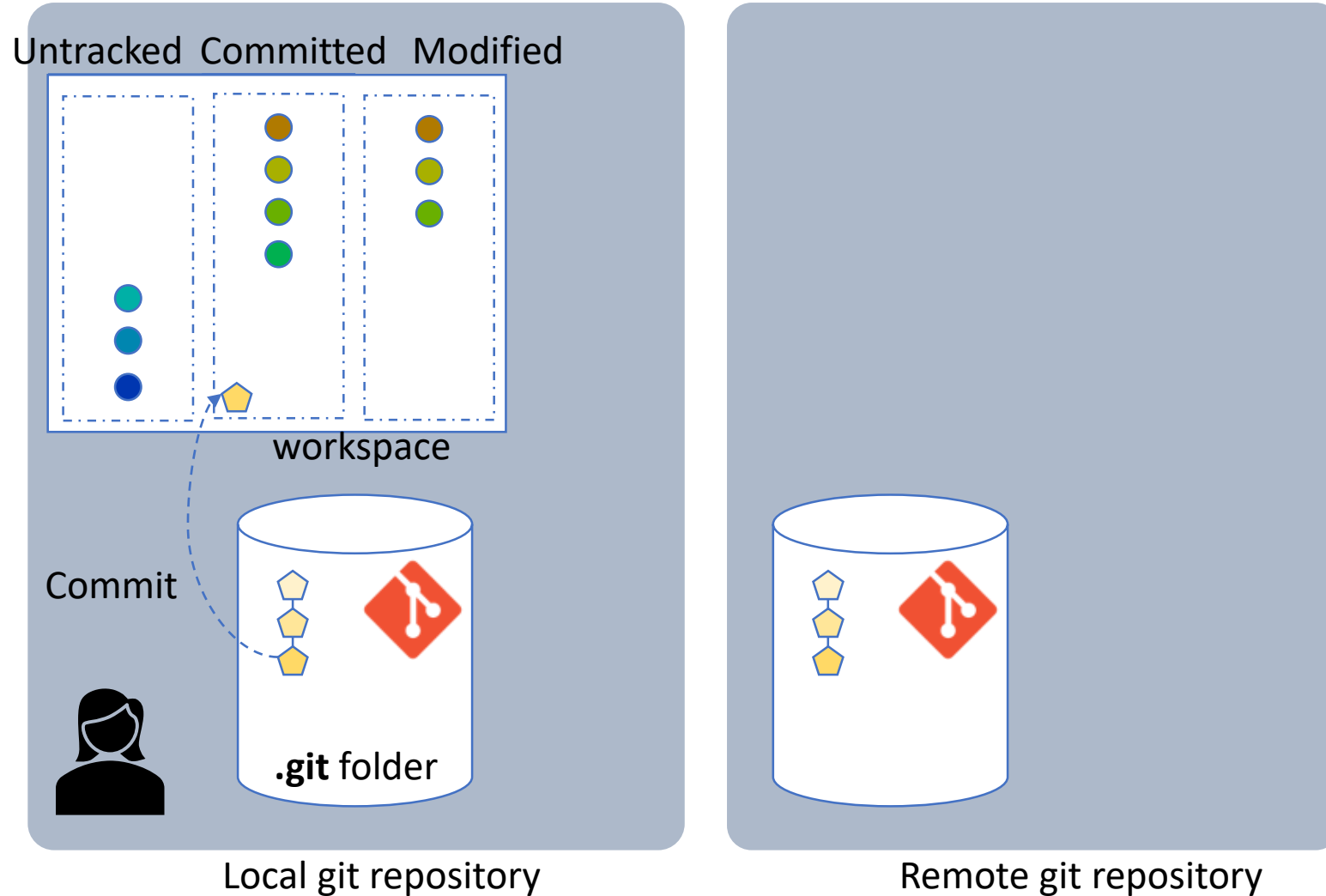
a new commit is created in the project repository

source: git-scm.com

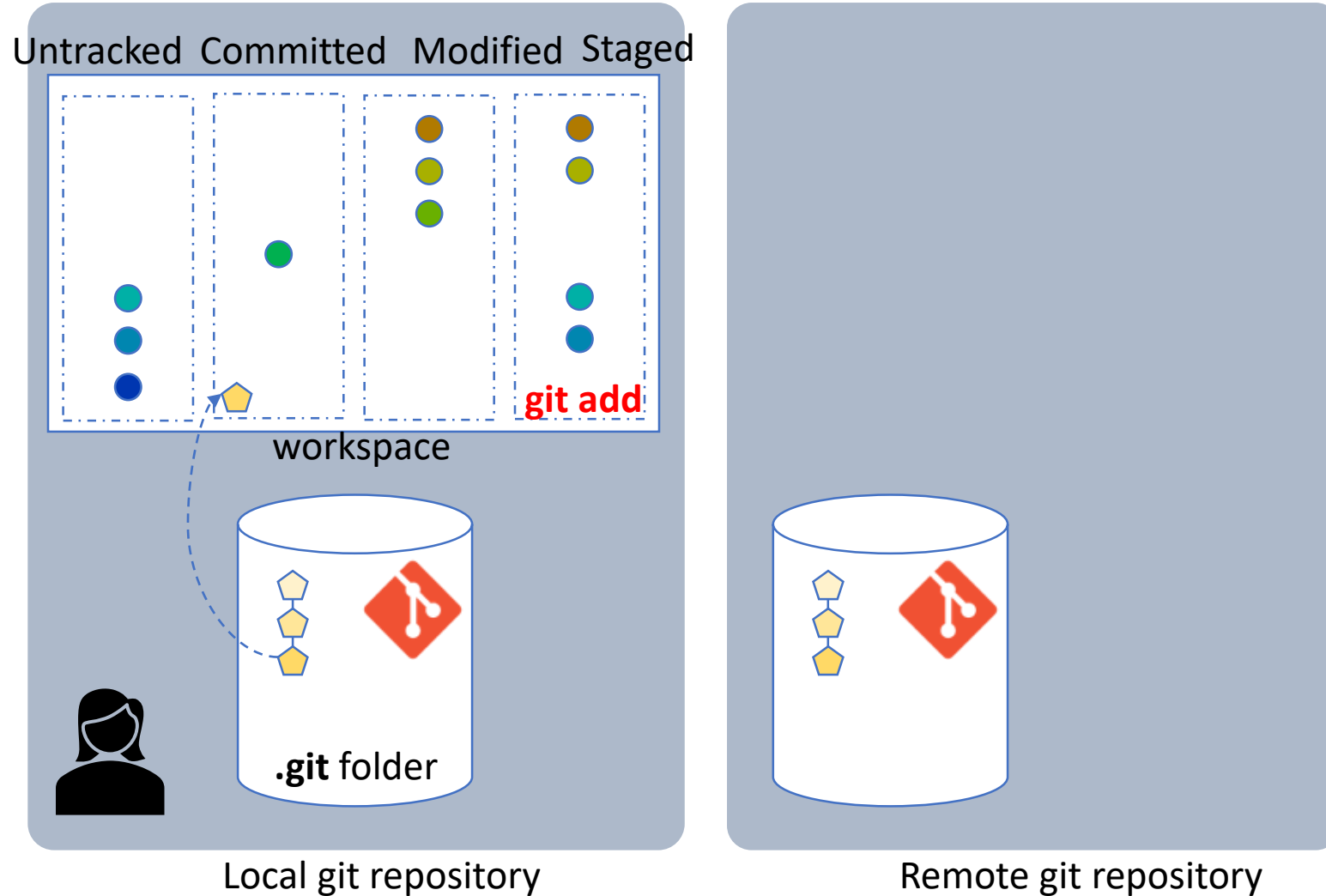
Git – Decentralized VCS



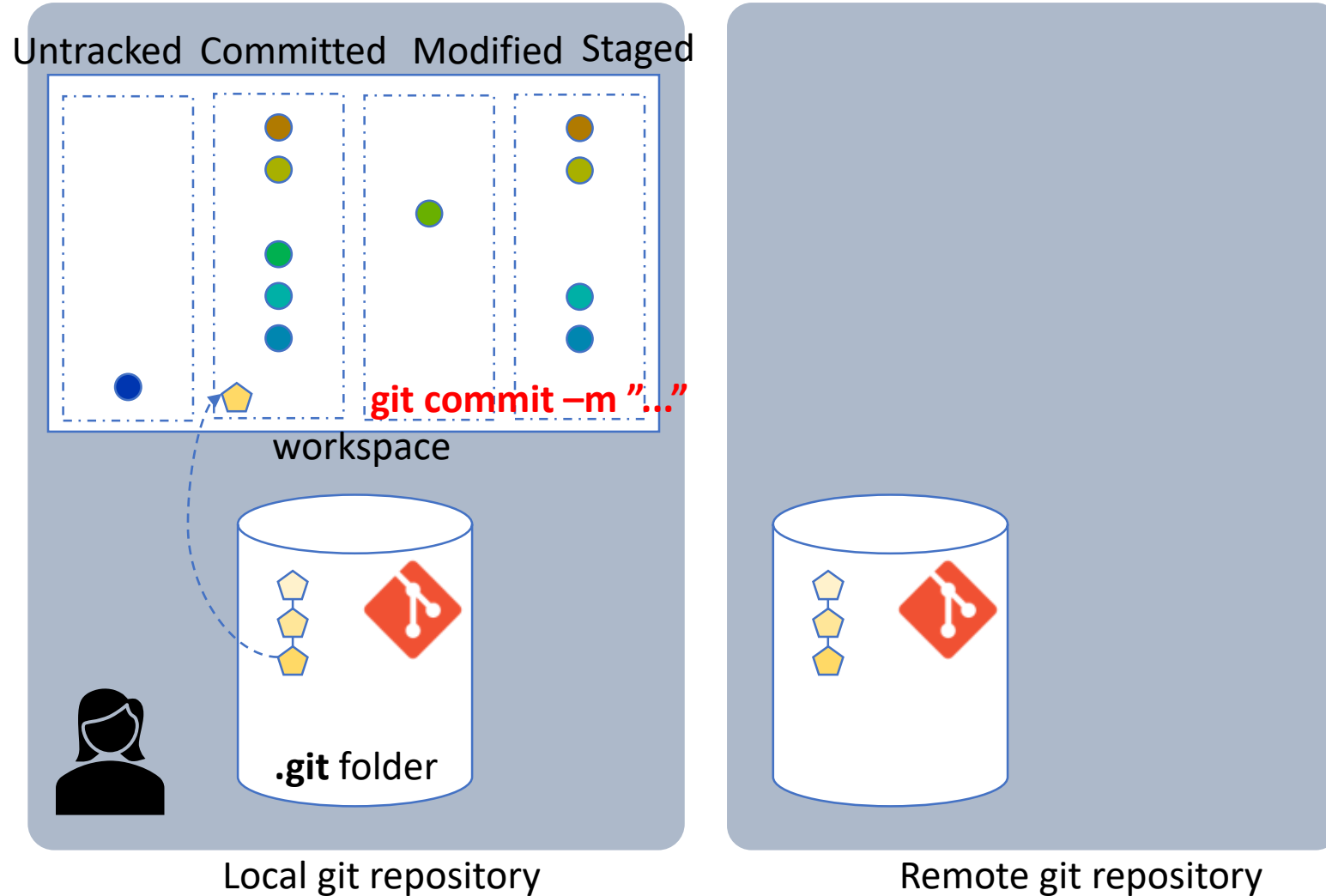
Git – Decentralized VCS



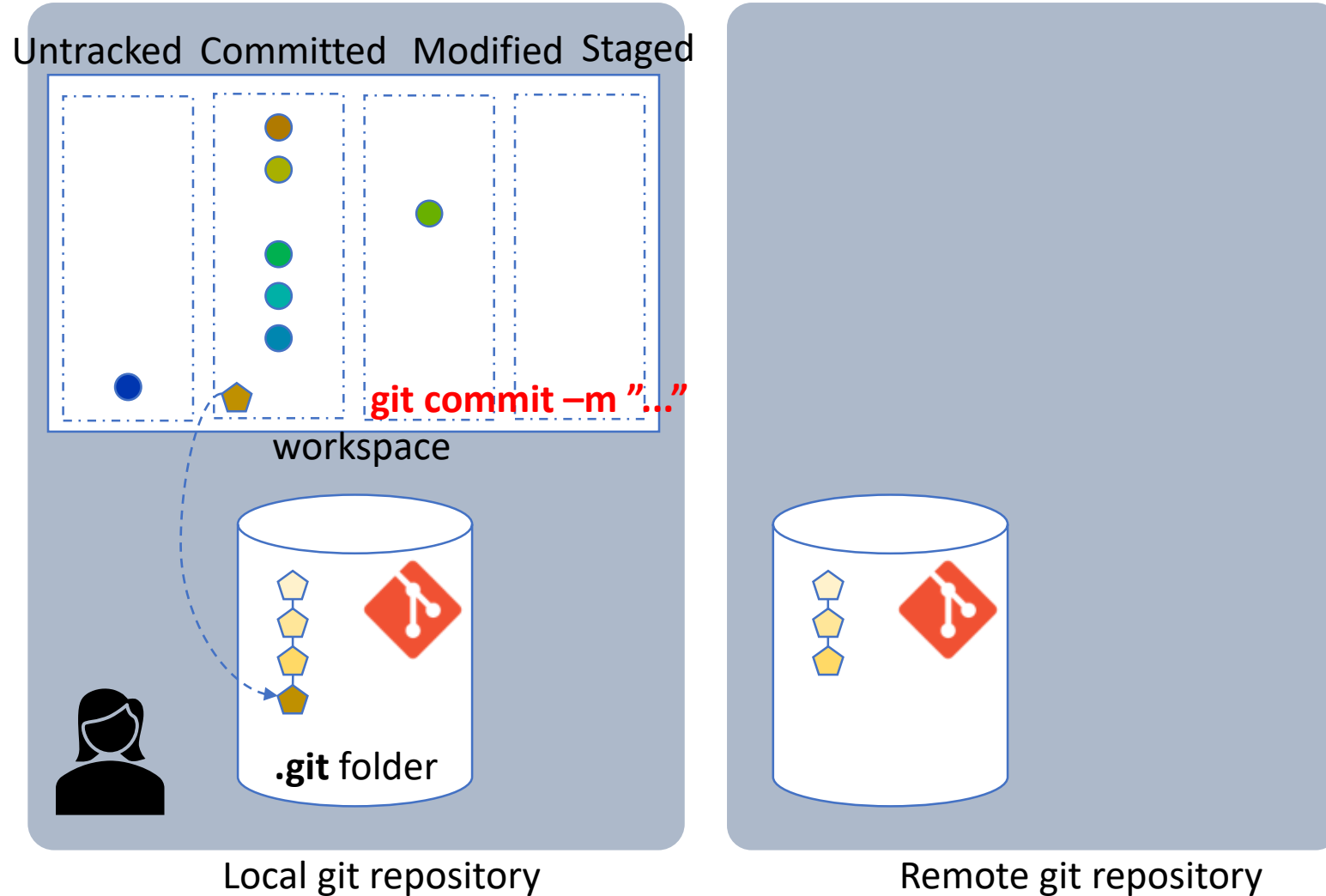
Git – Decentralized VCS



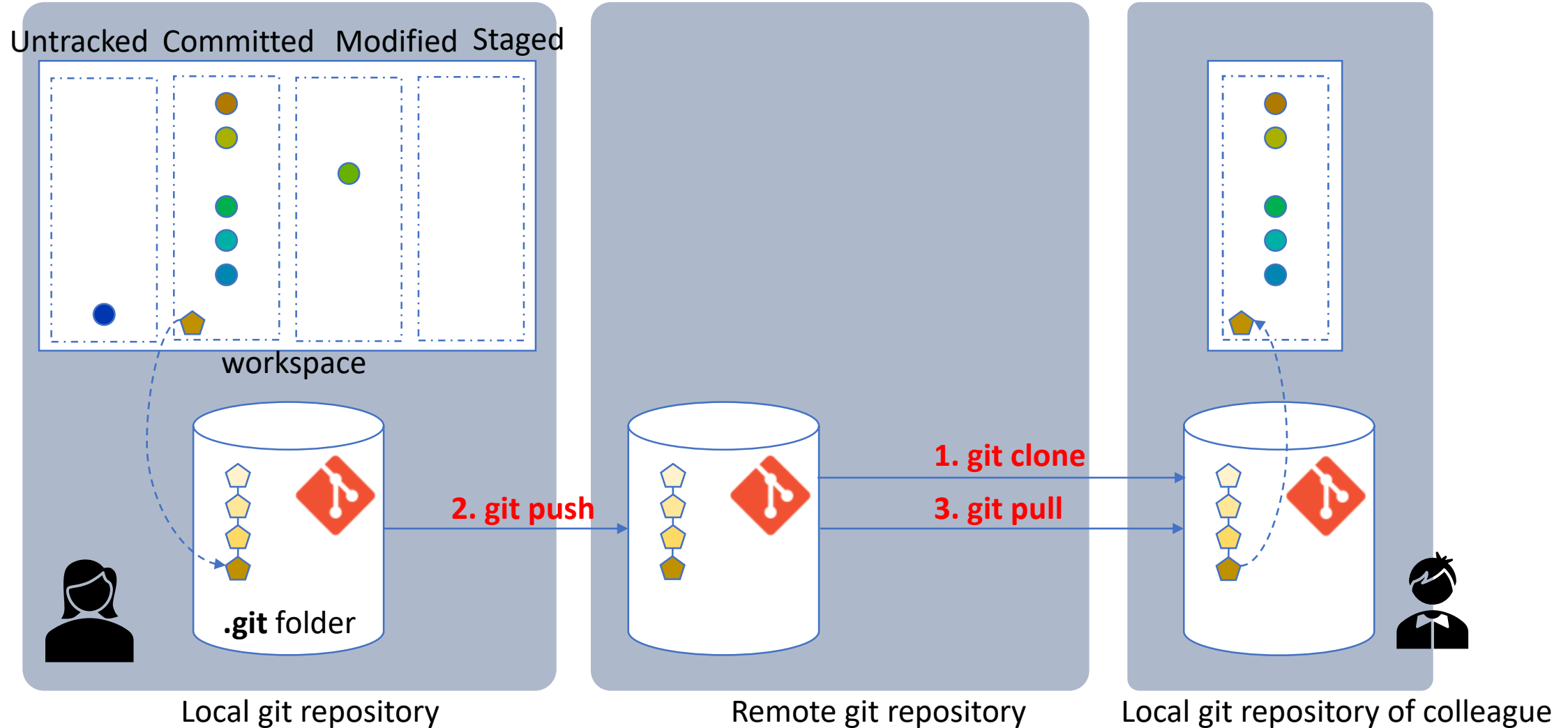
Git – Decentralized VCS



Git – Decentralized VCS



Git – Decentralized VCS



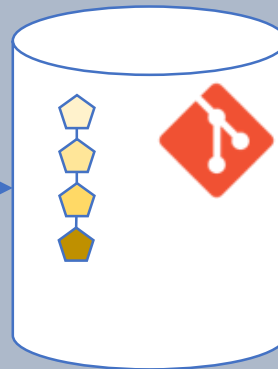
Git – Decentralized VCS



Remote Git Repository Manager

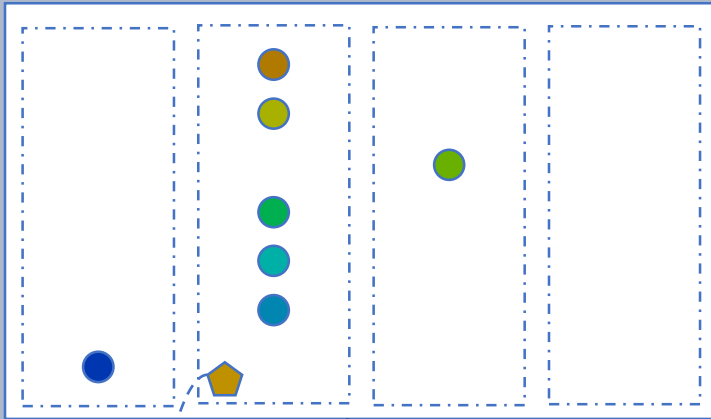


- Git repository management
- Access rights management
- Issue tracking
- Continuous integration, delivery



Remote git repository

Untracked Committed Modified Staged



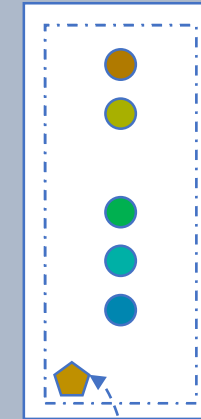
workspace



Local git repository

git clone
git push
git pull

git clone
git push
git pull



Local git repository of colleague

Git – Decentralized VCS



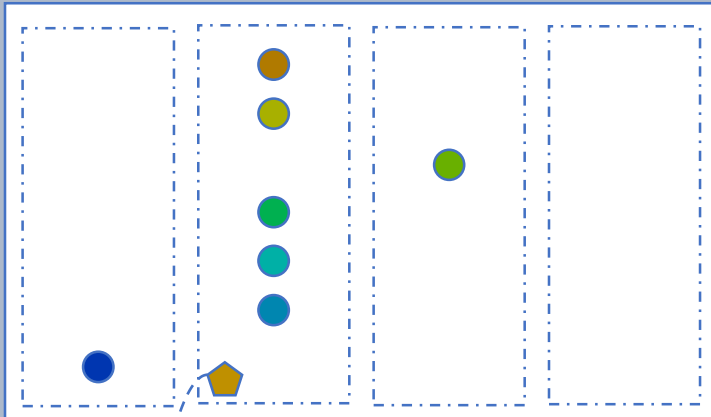
Remote Git Repository Manager



- Git repository management
- Access rights management
- Issue tracking
- Continuous integration, delivery

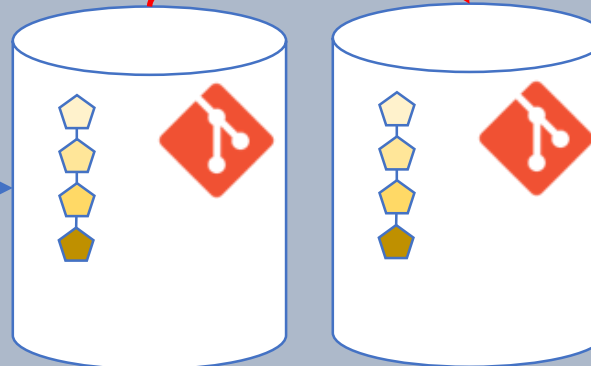
Fork

Untracked Committed Modified Staged



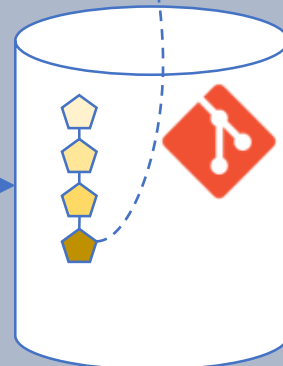
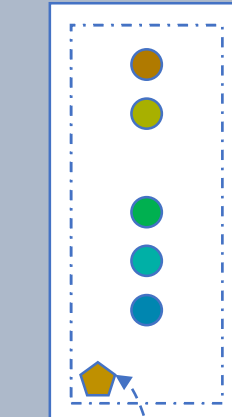
Local git repository

git clone
git push
git pull



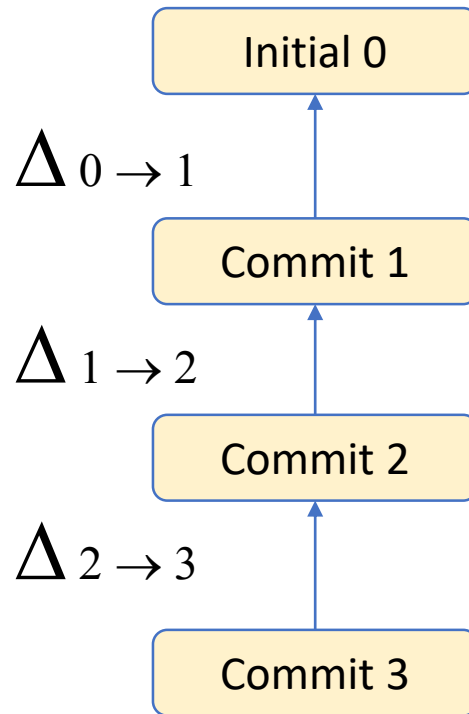
Remote git repository

git clone
git push
git pull



Local git repository of colleague

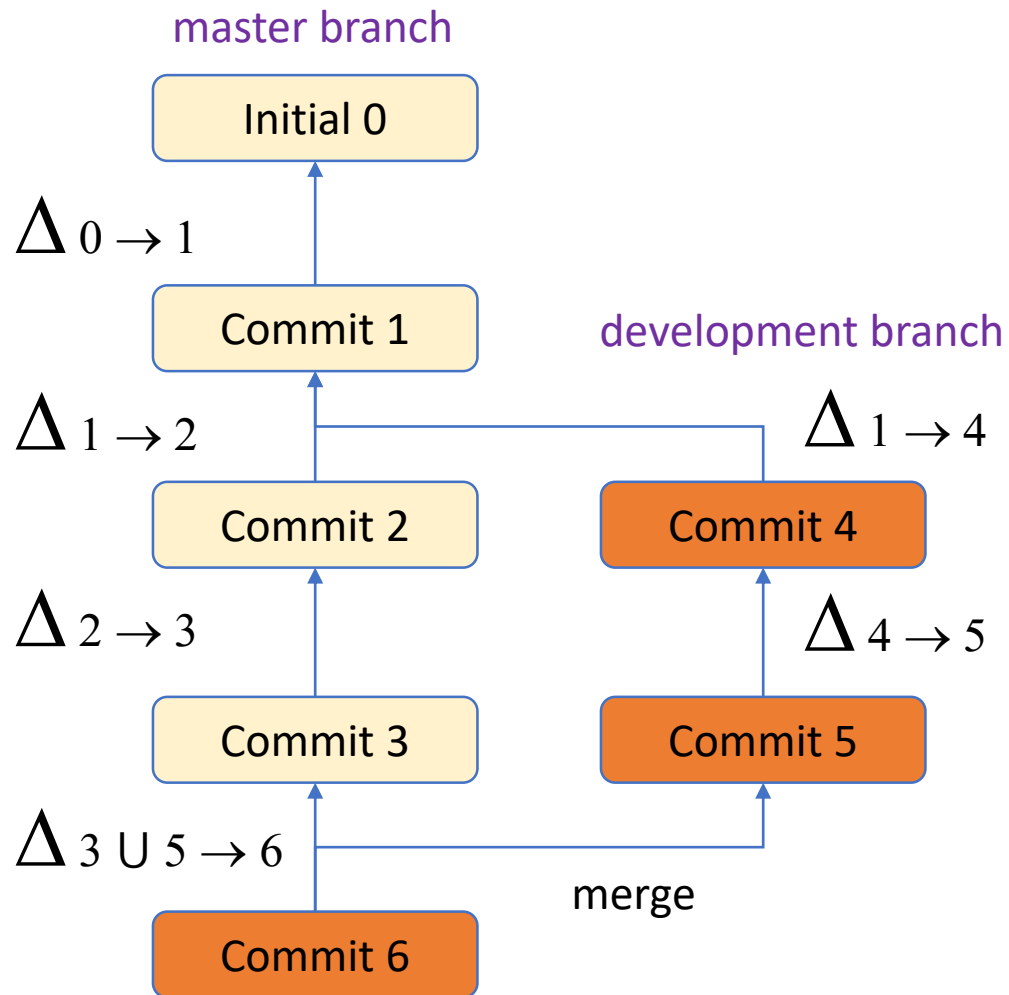
Git – Commit History (git log)



- A commit is a single point in the commit history. It is a snapshot of all the tracked files at that point.
- By default, a succession of commits follow a linear evolution

source: git-scm.com

Git – Branching



- A commit is a single point in the commit history. It is a snapshot of all the tracked files at that point.
- By default, a succession of commits follow a linear evolution
- Using Git, it is also possible to work in parallel on separate branches.

Git branching strategies

- GitFlow (complex)
- Github Flow (easy)
- Gitlab Flow (easy)
- OneFlow (medium)
- ...

source: git-scm.com

Git – Common Commands



- Create a new repository

```
git init
```

- Or, copy a repository locally from a remote origin – linked to origin (pull/push)

```
git clone https://dslab2021-renku.epfl.ch/gitlab/com490-pub/w2-exercises.git
```

- Verify state of repository (untracked, modified, staged files, commits ahead)

```
git status
```

- Stage files for a grouped commit

```
git add files+
```

```
git add .
```

- Commit files from staged area

```
git commit --message "description courte de la validation"
```


Git – Common Commands



- Display commit log (project history)

```
git log [--all][--graph]
```

- Checkout earlier commit

```
git checkout {commit-hash|branch-name}
```

- Propagate new commits to remote (origin) repository

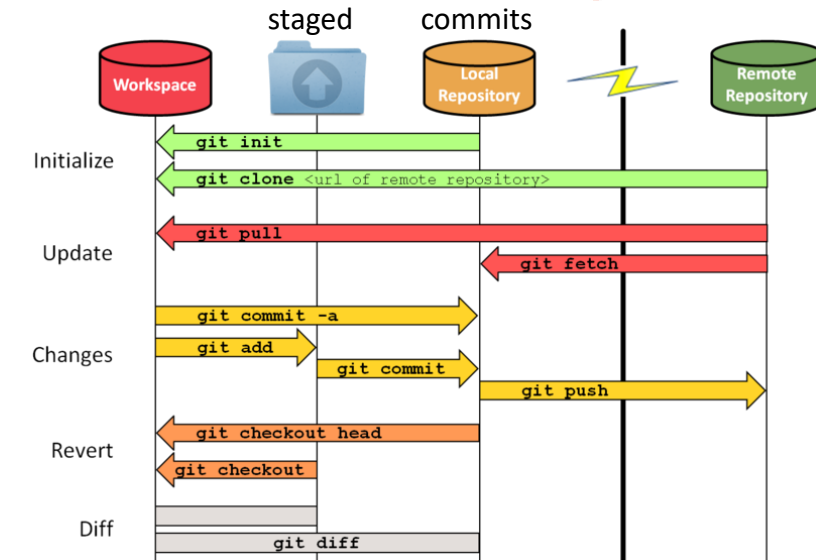
```
git push
```

- Retrieve latest commits from remote (origin) repository

```
git pull
```

- Get help

```
git command --help
```



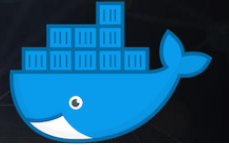
Git LFS - Large File Systems

- Git LFS is used to work with big files (e.g. data files)
 - Replace large files with tiny text files inside git (hash of file content), pointing to content stored outside git
 - Pointer files are swapped out for the real files at checkout (using git lfs smudge and clean)
- Basic commands
 - git lfs track *files+*
 - git lfs ls-files
 - git lfs fetch
- Most common mistakes
 - You believe you are reading the real file, but it is a pointer file (e.g. you forgot to do a git lfs fetch)
 - You have an existing file, track it with git lfs and replace it with a tiny pointer file, but the large file is still in the git history
- More advanced commands
 - git lfs migrate # convert history of git repository to use Git LFS (when large files are accidentally committed)
- Refs:
 - <https://git-lfs.github.com>
 - <https://github.com/git-lfs/git-lfs/wiki/Tutorial>

Runtime Environment – Containerization

- Containerization of applications: isolation and portability
 - "It's all about encapsulating or packaging the software code and all its dependencies so that it can run consistently and coherently on any infrastructure. ”
 - Popular containerization systems: docker (LXC, ...)

Docker containers



- Docker
 - Used by developers to build, share and run software in areas called containers
- Some Docker concepts
 - **Docker image**: a binary containing the executables and dependencies of the software.
 - **Docker container**: an instance of the image that resides in memory and uses the CPU of the host machine.
 - Containers are isolated from each other and from the machine for security reasons
 - A user can safely (in theory) be granted admin right inside a container, without being admin on the host machine

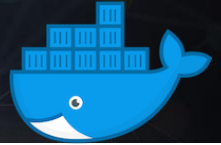
Dockerfile - Demo

- Dockerfiles are ‘recipes’
 - Provide instructions to build an image
 - Install software, include files, ...
- Docker on Renku
 - Each project has a Dockerfile
 - Defines the environment of sessions
 - You can change it !

Dockerfile 1.97 KB

```
1  # For finding latest versions of the base image see
2  # https://github.com/SwissDataScienceCenter/renkulab-docker
3  ARG RENKU_BASE_IMAGE=renku/renkulab-py:3.9-0.10.1
4  FROM ${RENKU_BASE_IMAGE}
5
6  # Uncomment and adapt if code is to be included in the image
7  # COPY src /code/src
8
9  # Uncomment and adapt if your R or python packages require extra linux (ubuntu) software
10 # e.g. the following installs apt-utils and vim; each pkg on its own line, all lines
11 # except for the last end with backslash '\' to continue the RUN line
12 #
13 # USER root
14 # RUN apt-get update && \
15 #     apt-get install -y --no-install-recommends \
16 #     apt-utils \
17 #     vim
18 # USER ${NB_USER}
19
20 # install the python dependencies
21 COPY requirements.txt environment.yml /tmp/
22 RUN conda env update -q -f /tmp/environment.yml && \
23     /opt/conda/bin/pip install -r /tmp/requirements.txt && \
24     conda clean -y --all && \
25     conda env export -n "root"
26
27 # RENKU_VERSION determines the version of the renku CLI
28 # that will be used in this image. To find the latest version,
29 # visit https://pypi.org/project/renku/#history.
30 ARG RENKU_VERSION=0.16.2
31
32 #####
33 # Do not edit this section and do not add anything below
34
```


Docker – Common Dockerfile Directives



Directive	Description
FROM alpine:latest	Base image. Have a look at https://hub.docker.com/explore/
USER root	Set the user for the following commands
RUN apt-get update && apt-get install -y renku	Run arbitrary shell commands Each RUN adds an overlay layer!
COPY ./requirements.txt /tmp/requirements.txt	Put a copy of a local file into the docker image, it also add a new layer!
EXPOSE 80 443	Expose a port (also -p port1:port2 of docker run)
VOLUME ["/var/log"]	Specify mount points for external volumes, i.e. expose external volumes inside docker container (see also -v or --mount of docker run)
ENTRYPOINT ["/bin/bash","-i"]	The command that is run when container is started (see also CMD)

Common mistakes:

- Each directive must be on one line, use \ to break lines
- Do not assume bash for the shell (it's a simple shell)
- Directives add up. If a file is added in a layer and is removed in the next – it becomes hidden, and is still using space.

References

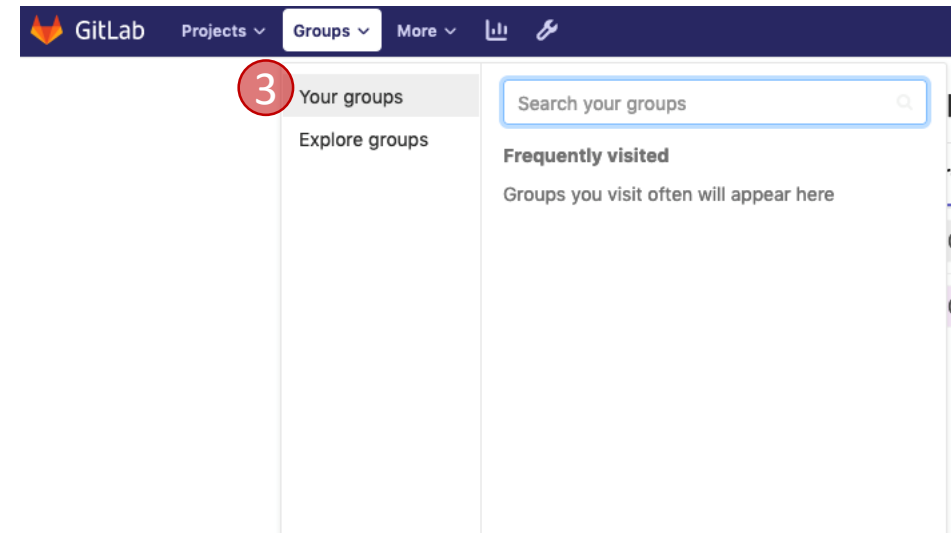
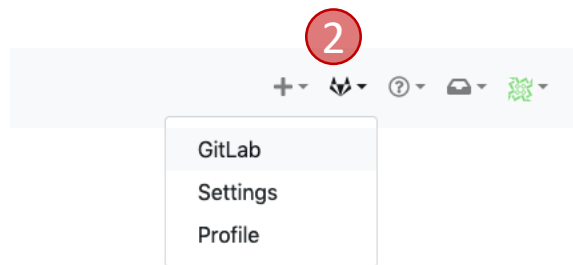
- Git, Gitlab
 - [git documentation](#)
 - [git lfs documentation](#)
 - [git lfs tutorial](#)
 - [gitlab basics](#)
- Docker
 - [Documentation](#)
- Conda
 - [Command reference](#)

Start your engines

<https://dslab2022-renku.epfl.ch/projects/com490/lab-course>

RENKU – Creating a Gitlab group

1. Login on <https://dslab2022-renku.epfl.ch>
2. Open the gitlab view (top right corner)
3. Under Groups select Your Groups or Explore Groups
4. Select **New group** (top right)



RENKU – Creating a Gitlab group

5. In New Group enter a name for your group
6. You can leave the description empty
7. Make sure your group is private
8. Click **Create group**.

New group

[Groups](#) allow you to manage and collaborate across multiple projects. Members of a group have access to all of its projects.

Groups can also be nested by creating [subgroups](#).

Projects that belong to a group are prefixed with the group namespace. Existing projects may be moved into a group.

Group name

My Awesome Group

Group URL

https://dsilab2021-renku.epfl.ch/gitlab/ my-awesome-group

Group description (optional)

Group avatar

Choose file... No file chosen

The maximum file size allowed is 200KB.

Visibility level

Who will be able to see this group? [View the documentation](#)

- ☒ Private
The group and its projects can only be viewed by members.
- ☐ Internal
The group and any internal projects can be viewed by any logged in user.
- ☐ Public
The group and any public projects can be viewed without any authentication.

Create group

Cancel

RENKU – Creating a Gitlab group

9. In the newly created group, list your collaborators under “Invite member” tab
10. Give your collaborators, the Developer role permission at a minimum
11. **Invite**.
12. You can later add/remove/edit the members of the group.
13. Remember to keep the group private!
14. See [gitlab groups](#) documentation

Invite member

GitLab member or Email address

9

Eric Pierre Bouillet Bouillet Sofiane Sarni Sarni Pamela Delgado Tao Sun

Choose a role permission

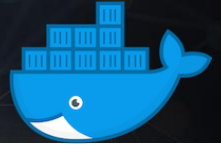
✓ Guest
Reporter
Developer
Maintainer
Owner

Expiration date

11 Invite

Additional slides

Docker - Overview



Dockerfile(s)

```
FROM ubuntu:latest
```

```
USER root
```

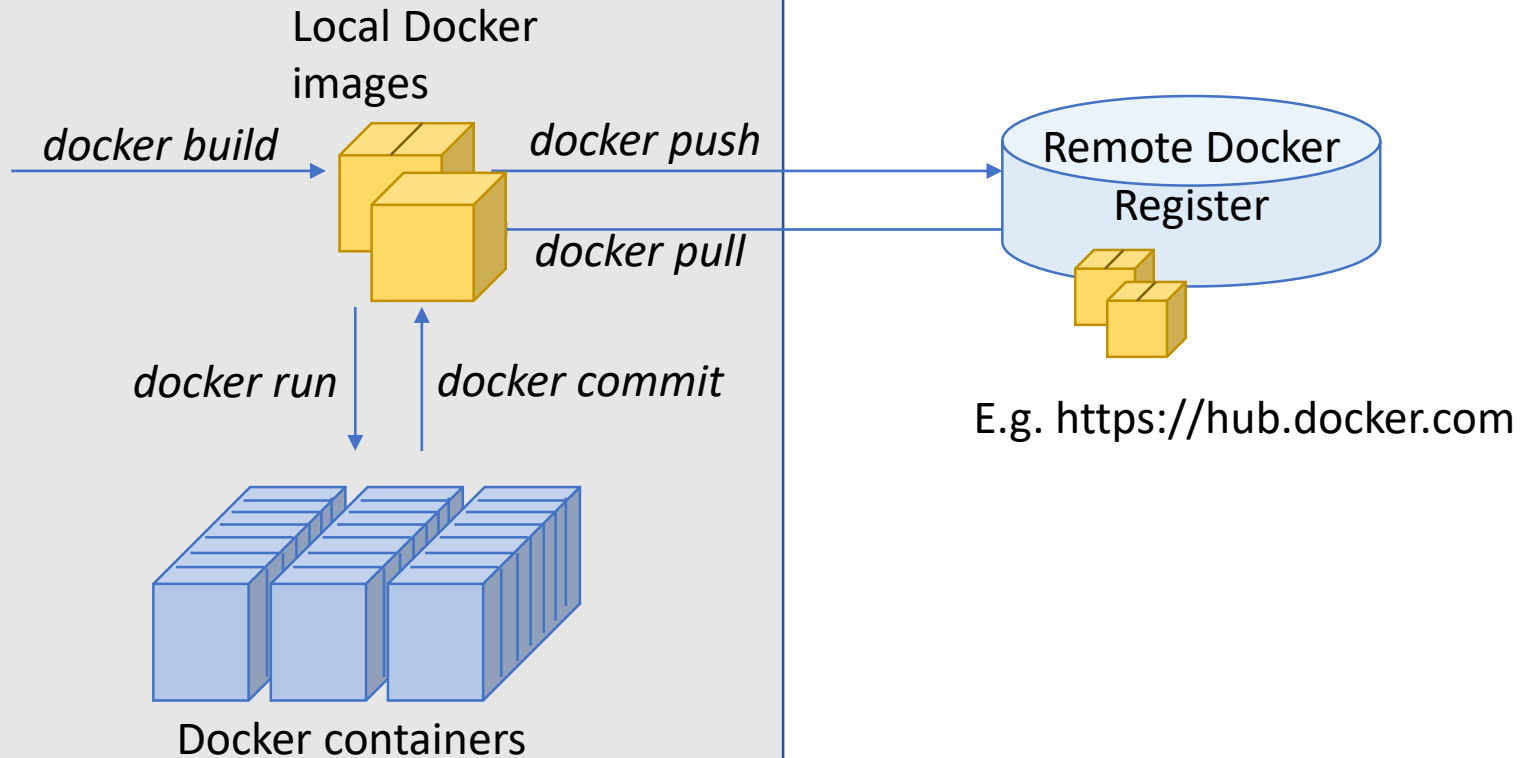
```
RUN apt-get update && \  
    apt-get install curl
```

```
USER jovyan
```

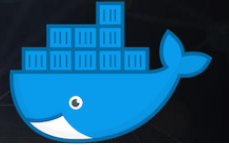
```
RUN curl -sSf -o Miniconda3.sh \  
    'https://repo.anaconda.com/miniconda/Miniconda3-  
    -latest-Linux-x86_64.sh'
```



Docker Engine running on server



Docker – Common Commands



- Create a Docker instance from an image, download image if not available locally

```
docker run [-it][--rm][-p port:port] [docker-image[:version]] [command]
```

- List docker instances

```
docker ps --all
```

- Stop and start docker instance

```
docker stop instance-id
```

```
docker start instance-id
```

- Kill and remove instance

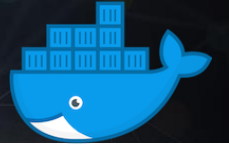
```
docker kill instance-id
```

```
docker rm instance-id
```

- Attach (stdin, stdout, stderr) to a running instance

```
docker attach instance-id
```

Docker – Common Commands



- Run command inside a running Docker instance (possibly with extended privileges)

```
docker exec [-it][--privileged][--user uid:gid] instance-id command
```

- Build image from a *Dockerfile*.

```
docker build .
```

- List images cached locally

```
docker images
```

- Remove local image

```
docker rmi instance-id
```

- Get help for other commands and additional arguments

```
docker [command] --help
```