# Applied Data Science Project

## L12 – Version Control

# Pillars

Design

Manage

Develop

Communicate

# Version control

- History of the development
  - what has been done 2 days ago
  - what a team member contributed to
  - restore a past version that resulted to be more robust than the last one

- Shared space for collaborators

- Monitor development branches and derive forks to be utilized for spinoff projects

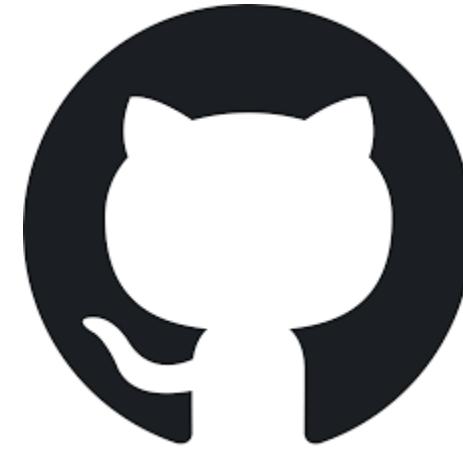- Package the development into tags (releases) that answer project milestones

# Version control





https://colab.research.google.com

Colaboratory natively stores different development versions each labeled either automatically by Colaboratory or defined by the team manually
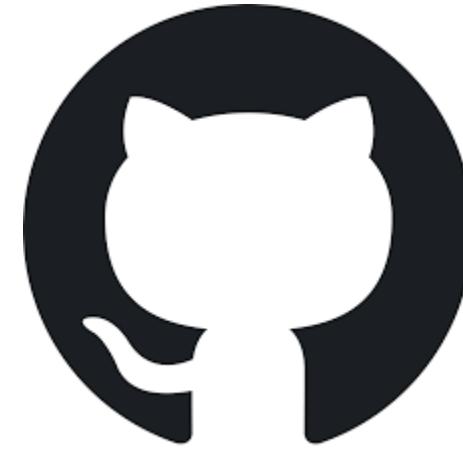
https://github.com

An application that manages local repositories and synchronises with remote repositories utilizing the git protocol. It also offers an intuitive web dashboard for the supervision of the project and analytics

# Two different needs



https://colab.research.google.com

When the project is about one notebook, this is the favorite option in an initial phase (lean option)



https://github.com

When the project grows in terms of files and modules that cannot stay in just a notebook

# Version control



List of differences between the versions indicated in the right panel with the checkboxes (left vs right)

# Version control

# Version control

# Version Control System

Repository is a shared folder where there are saved files necessary for the configuration, development and execution of a project

Workspace is a local folder that developers utilize to work

Both repository and workspace are part of a Version Control System (VCS) that defines the shared development and resolves conflicts

# VCS typologies

local: it is a simple database that tracks changes to files under version

centralized: it is a repository stored on a server, while clients access to individual files

distributed: clients have an integral copy of a repository

# VCS typologies

local: it is a simple database that tracks changes to files under version

centralized: it is a repository stored on a server, while clients access to individual files

distributed: clients have an integral copy of a repository

Git and GitHub

# How a project looks like

.git folder: it is managed automatically and contains the access information to the repository. This folder is created once there is a clone operation from the remote repository

whole working folder: it contains a project version checkout. Files of this folder are computed starting from the indexes present in .git. Those files can be modified by a user locally

stage area: it saves all files that will be included in the request to update the remote repository with the local changes. Files that are modified locally, but not stages, will not generate a request of change remotely

# Flow

linksfoundation.com
COPYRIGHT ©2021 LINKS

# Flow



Untracked | Unmodified | Modified | Staged

Add a file

Edit a file

Stage a file

Remove a file

Commit

# Basic commands

- git status: it gets the status of the files in the workspace
- git add: it stages a file
- git diff: it shows differences between the workspace copy and the one in the repository
- git commit: it does a commit to the workspace
- git mv: it moves a file from a folder to another and the change has an impact to the repository
- git rm: it removes a file with an impact remotely
- git log: it visualizes a log
- git init: it initializes an empty repository
- git clone: it clones a remote repository
- git pull: it downloads the changes done in the remote repository
- git push: it pushes the changes staged to the remote repository

# git clone

remote address of the repository we aim to clone locally

```
$ git clone git@github.com:adsp-polito/adsp-polito.github.io.git
Cloning into 'adsp-polito.github.io'...
remote: Enumerating objects: 56, done.
remote: Counting objects: 100% (56/56), done.
remote: Compressing objects: 100% (53/53), done.
remote: Total 56 (delta 24), reused 0 (delta 0), pack-reused 0
Receiving objects: 100% (56/56), 5.17 MiB | 8.01 MiB/s, done.
Resolving deltas: 100% (24/24), done.
```

# git status

$ git status
On branch main
Your branch is up to date with 'origin/main'.

Untracked files:
  (use "git add <file>..." to include in what will be committed)
        L8 - ADSP - AgileSwDev.pdf        a new file is in the workspace but not in the repository

nothing added to commit but untracked files present (use "git add" to track)

# git add

$ git add L8\ -\ ADSP\ -\ AgileSwDev.pd          add a new file in the workspace
$ git status
On branch main
Your branch is up to date with 'origin/main'.

Changes to be committed:
   (use "git restore --staged <file>..." to unstage)
          new file:   L8 - ADSP - AgileSwDev.pdf     a new file is staged

# git commit

$ git commit L8\ -\ ADSP\ -\ AgileSwDev.pdf  -m "add slides of the lecture made by prof. Marco Torchiano"
[main dd5e042] add slides of the lecture made by prof. Marco Torchiano
 1 file changed, 0 insertions(+), 0 deletions(-)
 create mode 100644 L8 - ADSP - AgileSwDev.pdf
$ git status
On branch main
Your branch is ahead of 'origin/main' by 1 commit.
  (use "git push" to publish your local commits)

nothing to commit, working tree clean

our workspace is ahead to the remote repository

# git push

```
$ git push
Enumerating objects: 4, done.
Counting objects: 100% (4/4), done.
Delta compression using up to 8 threads
Compressing objects: 100% (3/3), done.
Writing objects: 100% (3/3), 1.12 MiB | 9.83 MiB/s, done.
Total 3 (delta 1), reused 0 (delta 0)
remote: Resolving deltas: 100% (1/1), completed with 1 local object.
To github.com:adsp-polito/adsp-polito.github.io.git
   e05b240..dd5e042  main -> main
$ git status
On branch main
Your branch is ahead of 'origin/main' by 1 commit.
  (use "git push" to publish your local commits)

nothing to commit, working tree clean
```

the change is propagated remotely

TORINO

# git pull

```
$ git pull
remote: Enumerating objects: 4, done.
remote: Counting objects: 100% (4/4), done.
remote: Compressing objects: 100% (3/3), done.
remote: Total 3 (delta 1), reused 0 (delta 0), pack-reused 0
Unpacking objects: 100% (3/3), 7.31 MiB | 5.94 MiB/s, done.
From github.com:adsp-polito/adsp-polito.github.io
   dd5e042..f2c6801  main       -> origin/main
Updating dd5e042..f2c6801
Fast-forward
 L9 - ADSP - Scrum.pdf | Bin 0 -> 9285900 bytes
 1 file changed, 0 insertions(+), 0 deletions(-)
 create mode 100644 L9 - ADSP - Scrum.pdf
```

a new file present in the remote repository is also added to the local workspace

# Thank you for your attention.

Questions?

# CONTACTS

Giuseppe Rizzo

Program Manager (LINKS Foundation) and
Adjunct Professor (Politecnico di Torino)

giuseppe.rizzo@polito.it