Aurélie Névéol, LIMSI-CNRS

# Tools for
# Reproducibility
# &
# Collaborative Information Management

# Outline

- **Challenges**
  - Version management
  - Synchronous editing
  - … Beware of partial/relative information
- **Solutions**
  - Metadata
  - Cloud and collaborative editing
  - Versioning file systems
  - Version control systems

# Tracking the Whole Information
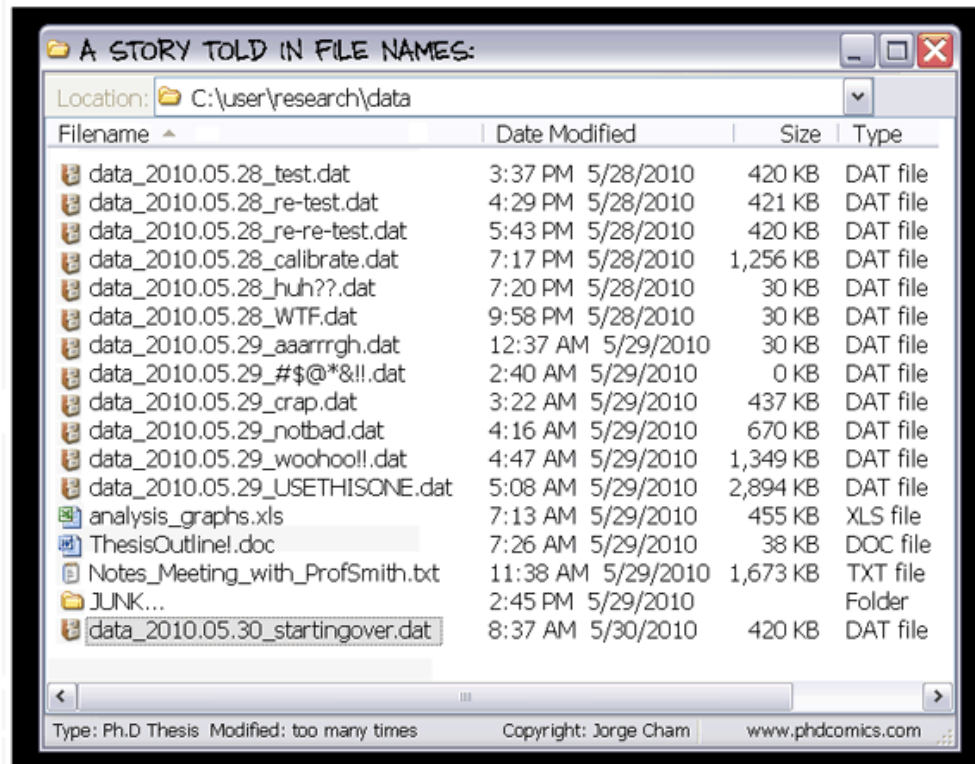
## Information integrity:

- Information can be distributed over multiple document sources
  - Reference to other documents
    - « see attached », .bib file in a latex document
    - library dependancies in code
  - Use of document or other IDs in corpus collections
    - Separate corpus file and gold standard key
    - (Database management systems, e.g. EMS, Workbench)

↳ Need to preserve integrity in updates
↳ Need to question integrity over time

# Metadata

## Common approaches to data management…
### (from PhD Comics: A Story Told in File Names, 28.5.2010)

# Metadata management

## Simple work practices
– File and document editing policies

## With advantages
– Quick and easy, no learning curve

## … and drawbacks
– Metadata piles up over time and user input
– Version control is hazardous
– Stability over time is questionable

# Versioning file systems

A versioning file system is any computer file system which allows a computer file to exist in several versions at the same time

- – Sample tool: RCS
- – Different from backup systems

# Versioning file systems

## Advantages

– Easy and transparent to the user

– Changes are dated and old versions available

## Drawbacks

– Some training required to use the tool

# Cloud, collaborative web

A series of tools
- – Doodle, Dropbox, Framapads, GoogleDocs, Skype, etc.

With advantages
- – Ubiquitous access (from multiple places and devices)
- – Multiple users can edit, merge seamless

… and drawbacks
- – Connexion required
- – Key aspects are provider dependent
  - • Security: data is physically stored by a provider
  - • Privacy: who has access to the data?
  - • Stability over time

# Version control systems

A *version control system* is a software tool that manages changes to documents, computer programs, large web sites, and other collections of information.

- Changes are identified by a number or letter code.
- Each revision is associated with a timestamp and the person making the change.
- Revisions can be compared, restored, and with some types of files, merged.
- Sample tools: sccs, cvs, rcs, svn, git

# Version control systems

- ## Advantages
    - Ubiquitous access (from multiple places and devices)
    - Multiple users can edit via local/distant copies
    - Does not require a connexion at all times
    - Key aspects user dependent: security, privacy, stability
- ## Drawbacks
    - Technical complexity: need a system administrator
    - Github ~cloud
    - Merge – dealing with concurrent editing sometimes tricky

# In practice, what to use when?

- **By yourself - rush, small scale need):** metadata
- **By yourself - all needs:** RCS, GIT, SVN
- **Multiple users - low privacy or security requirements:** cloud
- **Multiple users - all needs:** GIT/SVN

- **GIT/SVN all the way…**
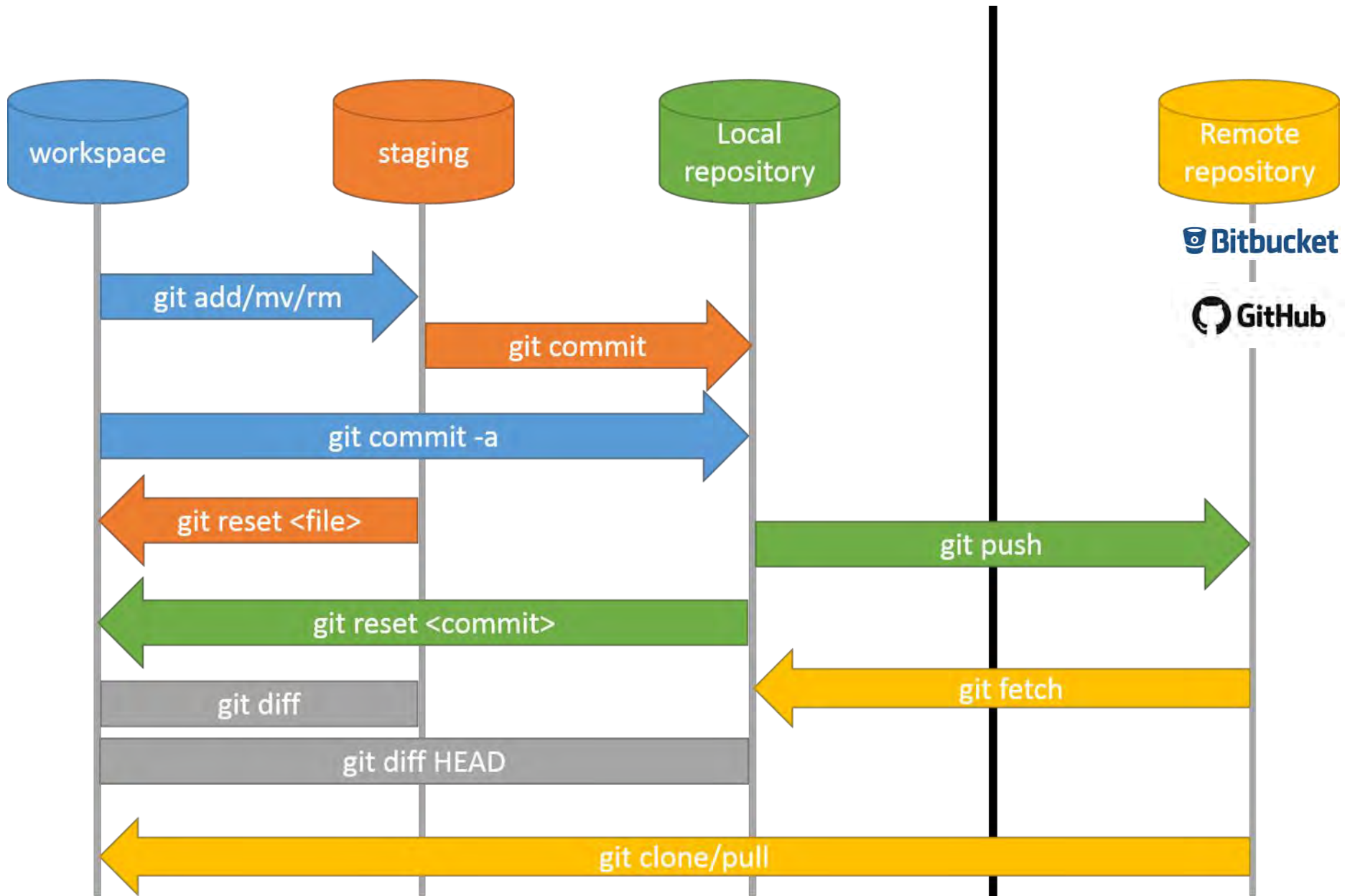  - More details now
  - Pointers to tools & how to install

# A Brief Introduction
# to Version Control



Created by Linus Torvalds, in 2005

git (noun) : [*british informal*] An unpleasant
or contemptible person

http://www.moxie.io/images/git-operations.png

# Download Git

[https://git-scm.com/download](https://git-scm.com/download)

Available for all major platforms
OS X, Windows, Linux

# Git Graphical User Interfaces

https://git-scm.com/download/guis

# Git organises snapshots for you



**Snapshots of the same file over time**

# How does it work?

**git** **command** **options**

(or equivalent button clicking
in your favorite GUI)

# Repository / « Repo »

A folder where Git is tracking changes

# Make a new repository

```
cd my_projcet
git init
```

# *commit*

- Create a snapshot of your repository
- Commit the changes you have made

# Steps to *commit*

`git status` (which files changed?)

`git diff` (which lines changed?)

`git add my_report.tex` (I want this file in my next commit)

`git add pic.jpg papers.bib` (**these** files too!)

`git commit` (OK, save a snapshot of what I just added)

`git log` (Show me the commit history)

# *Commit* message should have

# a concise summary.

Put it on the first line. 70 characters or less.

*Commit* message should have

a detailled explanation.

Think lab notebook. Wrap your text at 70 characters

# Informative message

< 70 characters

Updated article classification features

The function compute_features now produces additional features related to token characteristics and external clusters.

# Beware of commit message drifting



| | COMMENT | DATE |
|---|---|---|
| ○ | CREATED MAIN LOOP & TIMING CONTROL | 14 HOURS AGO |
| ○ | ENABLED CONFIG FILE PARSING | 9 HOURS AGO |
| ○ | MISC BUGFIXES | 5 HOURS AGO |
| ○ | CODE ADDITIONS/EDITS | 4 HOURS AGO |
| ○ | MORE CODE | 4 HOURS AGO |
| ○ | HERE HAVE CODE | 4 HOURS AGO |
| ○ | AAAAAAAA | 3 HOURS AGO |
| ○ | ADKFJSLKDFJSDKLFJ | 3 HOURS AGO |
| ○ | MY HANDS ARE TYPING WORDS | 2 HOURS AGO |
| ○ | HAAAAAAAAANDS | 2 HOURS AGO |

AS A PROJECT DRAGS ON, MY GIT COMMIT MESSAGES GET LESS AND LESS INFORMATIVE.

https://m.xkcd.com/1296/

# A History of your *commits*

```
┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│   60d1ef3    │◄─────│   4f6b996    │◄─────│   2c538a0    │
└──────────────┘      └──────────────┘      └──────────────┘
```

SHA-1
identification key
*(hash key)*

You messed something up and you want to go back

# git checkout - my_file

Go back to my_file version per most recent *commit*

# git checkout `60d1ef3`

Go back to the *commit* labeled as 60d1ef3

(and then you can branch out from there – coming up)

You broke something and you want to change history

**git reset --hard 60d1ef3**

Revert everything to the commit labelled 60d1ef3

# git reset --hard

Revert everything to the most recent commit

# What if you wanted that code but not at that moment?
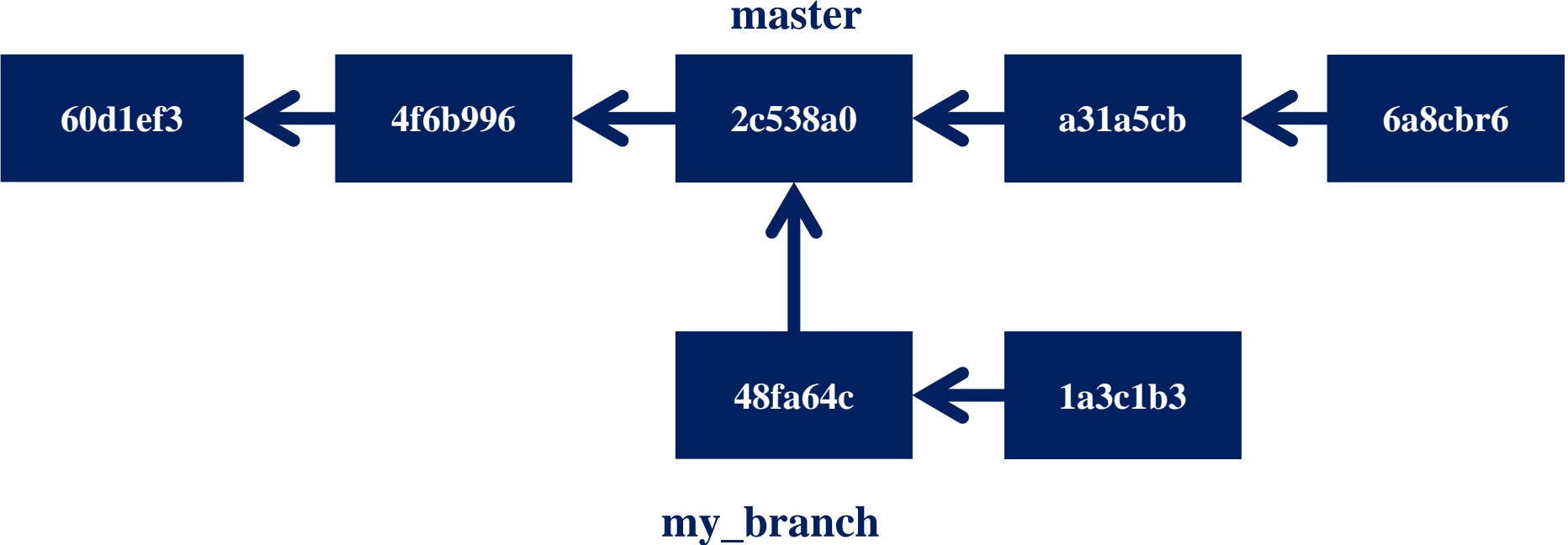
# Branching

Track separate versions of your code

# **Making a branch**

`git branch` (list the branches of the repo)

`git branch` `my_branch` (make a branch called `my_branch`)

`git checkout` `my_branch` (switch to `my_branch`)

Now you can *commit* changes to that branch.

master

60d1ef3 ← 4f6b996 ← 2c538a0 ← a31a5cb ← 6a8cbr6

48fa64c ← 1a3c1b3

my_branch

# Take Home Messages

# Git helps you organize snapshots of your projects

# These snapshots are called *commits*

If you mess up,
you can always go back
as long as there's been a commit.

Branches let you
try out new ideas
without losing access to the
version that works.

**90 % of the time :** `status / log / commit / push / pull`

**8 % :** `checkout / merge`

**2 % : other commands**

# Further documentation

___

**gitimmersion.com** An interactive tutorial

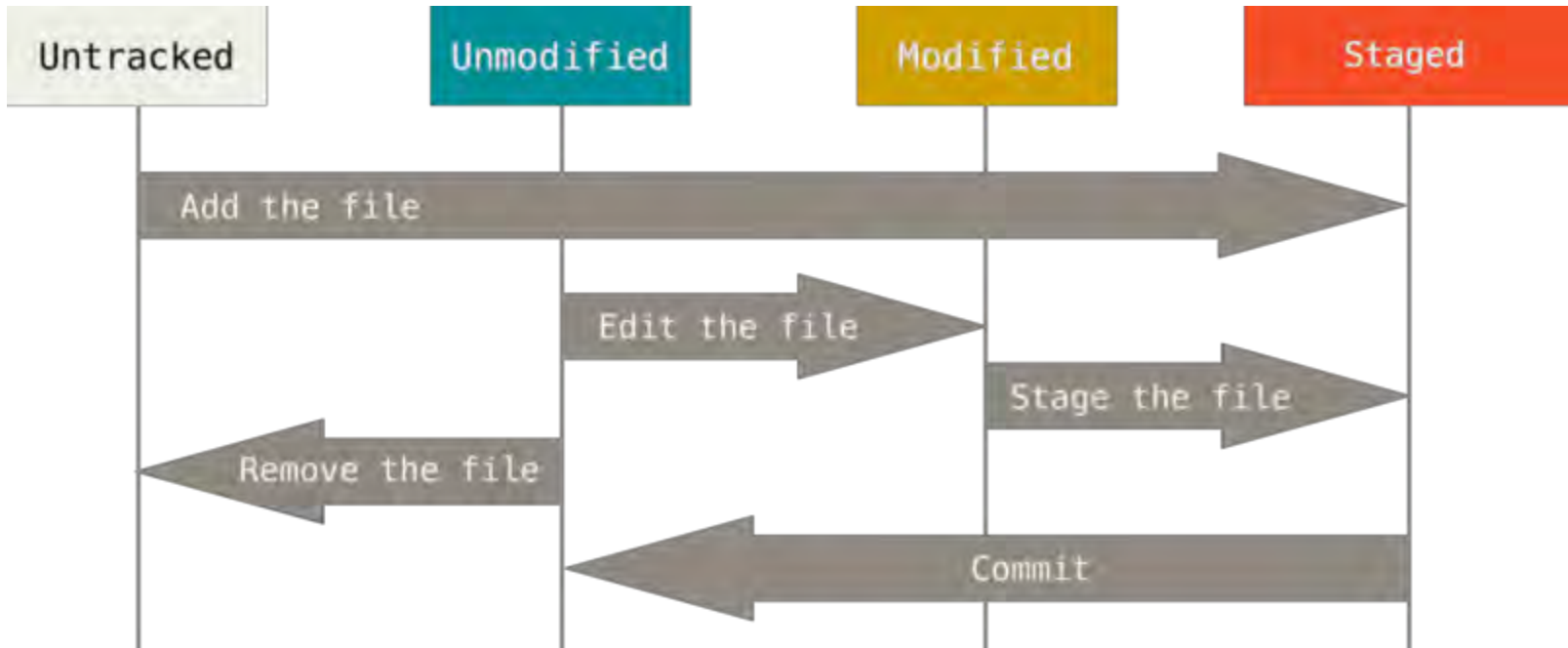**book.git-scm.com** A detailed text book on Git

**think-like-a-git.net** Advance use of Git

**nvie.com/posts/a-successful-git-branching-model/** on using branches

# Acknowledgements

- Kevin Chen (Columbia University)
- Nicolas Grenier, Thomas Lavergne (LIMSI-CNRS)
- Patrick Paroubek (LIMSI-CNRS)
- Pierre Zweigenbaum (LIMSI-CNRS)

https://git-scm.com/book/en/v2/Git-Basics-Recording-Changes-to-the-Repository

# Methods in Research on Research

**A presentation delivered at the**

**first MiRoR training event**
**October 19-21, 2016**
**Ghent, Belgium**

www.miror-ejd.eu          @MirorProject