Problems and Caution
ooooo

What is Git?
ooooo

Using Git
oo

Caveats
o

# Data Management for Reproducable Research

Thomas R. Cook

23-OCT 2015

Problems and Caution
○○○○○

What is Git?
○○○○○

Using Git
○○

Caveats
○

Problems and Caution
○○○○○

What is Git?
○○○○○

Using Git
○○

Caveats
○

Long-term reproducability and Mysterious Data:

- Common Scenario:
- Get novel data file
    - Make some changes to it
    - Save over original file

Long-term reproducability and Mysterious Data:

- Common Scenario:
- Get novel data file
    - Make some changes to it
    - Save over original file

Six months later. . .

- Return to project. . .
    - What does log_inerv_1234.b mean? How did I get it? Why is it driving my results?
- Even worse if someone asks for your replication data
    - You need to be able to explain how you arrived at a given variable/model/etc
- Moar worse if your co-author created log_inerv_1234.b

# Six months later...

- Return to project...
  - What does `log_inerv_1234.b` mean? How did I get it? Why is it driving my results?
- Even worse if someone asks for your replication data
  - You need to be able to explain how you arrived at a given variable/model/etc
- Moar worse if your co-author created `log_inerv_1234.b`

Six months later...

- Return to project...
  - What does log_inerv_1234.b mean? How did I get it? Why is it driving my results?
- Even worse if someone asks for your replication data
  - You need to be able to explain how you arrived at a given variable/model/etc
- Moar worse if your co-author created log_inerv_1234.b



- Make R script, DO file or other script that generates data

# Two big challenges with managing data

① Track file changes over tiem
- Long-term reproducibility
- Version management

② Collaboration with others

Common Solutions:

⑤ Edit data in-place (!)

⑥ Dropbox

⑦ Track Changes/time-machine

⑧ Email

⑨ New folder per version

Problems and Caution
OOOO●

What is Git?
OOOOO

Using Git
OO

Caveats
O

That's a start...

### But what about these other nightmare scenarios:

- Someone asks for old version of replication data

# That's a start...

## But what about these other nightmare scenarios:

- Someone asks for old version of replication data
- Coauthor deletes files in his/her dropbox folder

That's a start. . .

### But what about these other nightmare scenarios:

- Someone asks for old version of replication data
- Coauthor deletes files in his/her dropbox folder
- You and co-author try to work on data at the same time (the dreaded `conflited copy`)

That's a start. . .

> ### But what about these other nightmare scenarios:
>
> - Someone asks for old version of replication data
> - Coauthor deletes files in his/her dropbox folder
> - You and co-author try to work on data at the same time (the dreaded `conflited copy`)
> - Need to identify how project today differs from version 6 months ago

That's a start...

## But what about these other nightmare scenarios:

- Someone asks for old version of replication data
- Coauthor deletes files in his/her dropbox folder
- You and co-author try to work on data at the same time (the dreaded `conflited` copy)
- Need to identify how project today differs from version 6 months ago
- Need to identify prior/abandoned approaches to analysis

That's a start. . .

> ### But what about these other nightmare scenarios:
>
> - Someone asks for old version of replication data
> - Coauthor deletes files in his/her dropbox folder
> - You and co-author try to work on data at the same time (the dreaded `conflited` copy)
> - Need to identify how project today differs from version 6 months ago
> - Need to identify prior/abandoned approaches to analysis
> - The list goes on. . .

That's a start. . .

### But what about these other nightmare scenarios:

- Someone asks for old version of replication data
- Coauthor deletes files in his/her dropbox folder
- You and co-author try to work on data at the same time (the dreaded `conflited copy`)
- Need to identify how project today differs from version 6 months ago
- Need to identify prior/abandoned approaches to analysis
- The list goes on. . .
- Git can help resolve all of these

Git is. . .

- Distributed Version Control System

Problems and Caution
00000

What is Git?
●0000

Using Git
00

Caveats
0

Git is. . .

- Distributed Version Control System
  - Distributed $\rightarrow$ decentralized, many users

Problems and Caution
○○○○○
What is Git?
●○○○○
Using Git
○○
Caveats
○

Git is. . .

- Distributed Version Control System
    - Distributed → decentralized, many users
    - Version Control → track changes, enable rollbacks,

Git is. . .

- Distributed Version Control System
    - Distributed → decentralized, many users
    - Version Control → track changes, enable rollbacks,
    - System → System

Git is. . .

- Distributed Version Control System
  - Distributed $\rightarrow$ decentralized, many users
  - Version Control $\rightarrow$ track changes, enable rollbacks,
  - System $\rightarrow$ System
- Think: "Track changes" on steroids

# Repos

- Git tracks sets of files – multiple files at once
- Folder with set of files tracked by Git: Repository
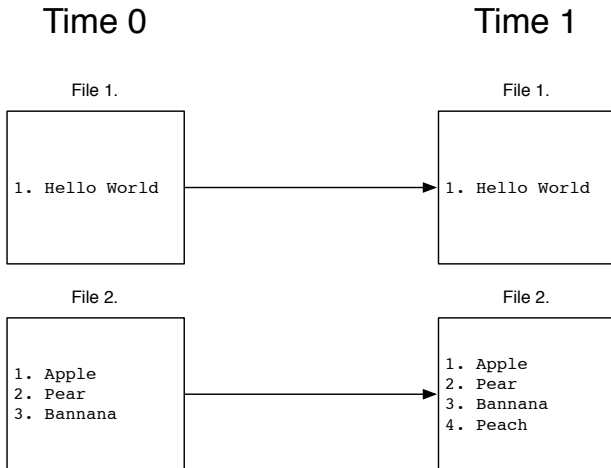    - Generally, a Git repo looks and works just like a folder
- Think: Repo  project

## Commits

- A snapshot of (specified) files tracked by Git
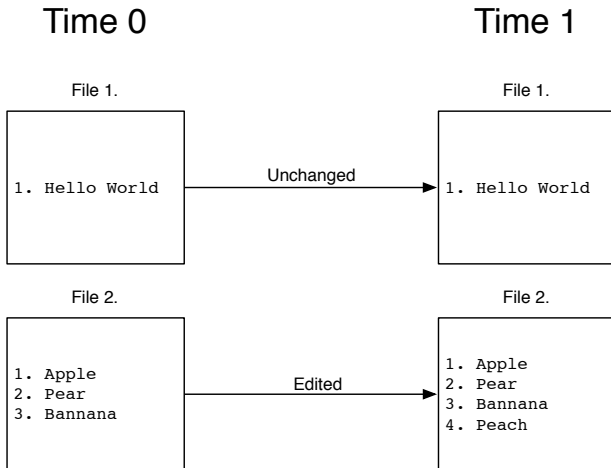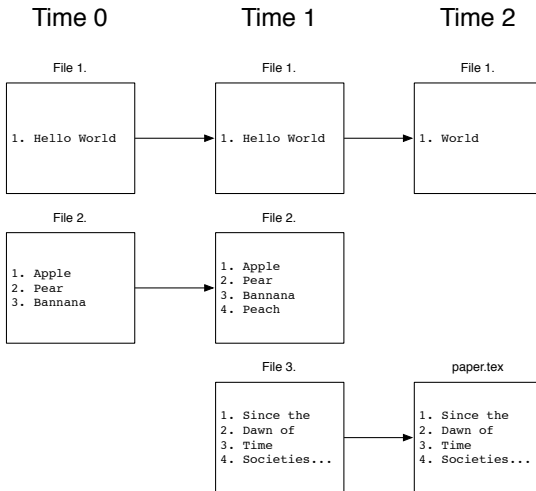  - Captures *changes* in specified files (since last commit)

## Commits

# File Perspective

### Time 0                                    ### Time 1

File 1.                                       File 1.

```
1. Hello World
```
→
```
1. Hello World
```

File 2.                                       File 2.

```
1. Apple
2. Pear
3. Bannana
```
→
```
1. Apple
2. Pear
3. Bannana
4. Peach
```

## Commits

## File Perspective

### Time 0                                    ### Time 1

File 1.                                        File 1.

1. Hello World  —— Unchanged ——>  1. Hello World

File 2.                                        File 2.

1. Apple                                       1. Apple
2. Pear       —— Edited ——>                    2. Pear
3. Bannana                                     3. Bannana
                                               4. Peach

Problems and Caution
ooooo

What is Git?
oo●oo

Using Git
oo

Caveats
o

## Commits

## Git Perspective
## (Diff Perspective)

### Time 0
(commit1)

File 1.

file added ⟶ No changes

File 2.

File 2.

file added ⟶ + 4. Peach

### Time 1
(commit 2)

File 1.

## Commits

- A snapshot of (specified) files tracked by Git
  - Captures *changes* in specified files (since last commit)
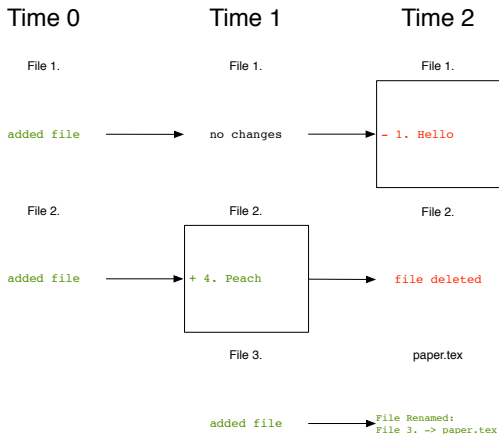  - Captures Files Added/Removed/Moved

## Commits

### File Perspective

Problems and Caution
○○○○○

What is Git?
○○●○○

Using Git
○○

Caveats
○

# Commits

Problems and Caution
00000

What is Git?
00●00

Using Git
00

Caveats
0

Commits

- A snapshot of (specified) files tracked by Git

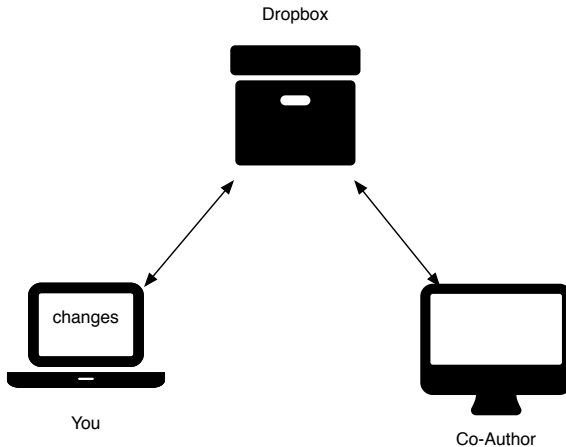  - Upshot: Can track file changes very closely over time

# Distribution/Collaboration

- Git enables Collaboration – it is a distributed system.
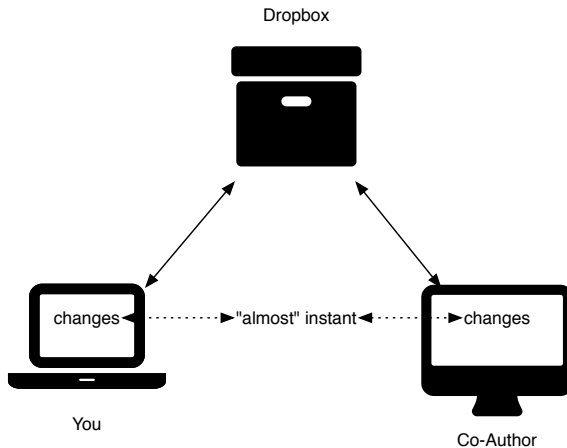  - Contrast to Dropbox

Problems and Caution
OOOOO
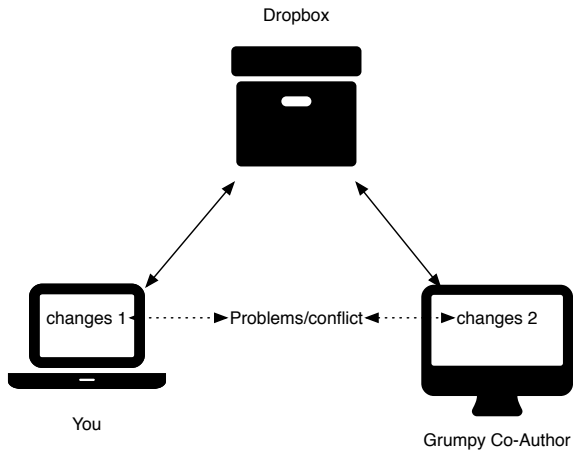
What is Git?
OOO●O

Using Git
OO

Caveats
O

# Distribution/Collaboration

Dropbox

Synchronized files

Synchronized files

You

Co-Author

# Distribution/Collaboration



Dropbox

changes

You

Co-Author

Problems and Caution
○○○○○

What is Git?
○○○●○

Using Git
○○

Caveats
○

## Distribution/Collaboration



Dropbox

changes ┄┄┄► "almost" instant ◄┄┄ ► changes

You

Co-Author

# Distribution/Collaboration

Dropbox



changes 1 ········►Problems/conflict◄········ ►changes 2

You

Grumpy Co-Author

# Distribution/Collaboration

- Git enables Collaboration – it is a distributed system.

  - Download repo to local computer
  - Make changes and commit
  - Push changes to server when ready
  - Pull changes from server when ready

# Distribution/Collaboration
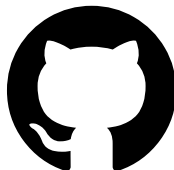
Github
(or other service -- bitbucket, aws, etc. )



You

Co-Author

# Distribution/Collaboration

Github
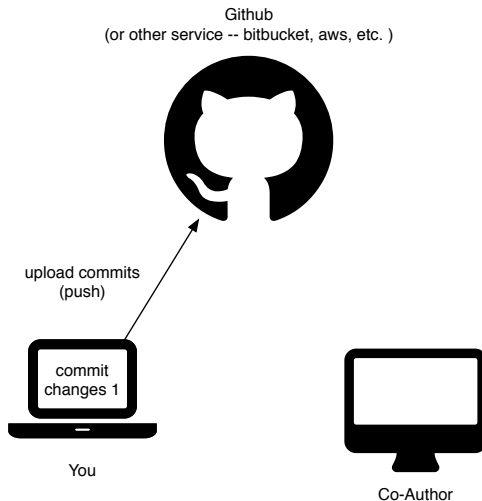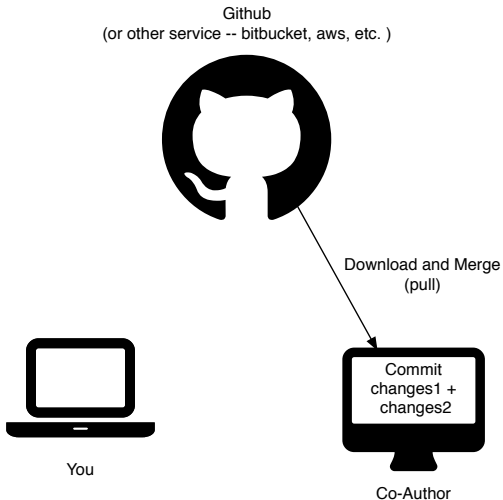(or other service -- bitbucket, aws, etc. )



changes 1

You

Asynchronous

changes 2

Co-Author

# Distribution/Collaboration



Github
(or other service -- bitbucket, aws, etc. )

upload commits
(push)

commit
changes 1

You

Co-Author

# Distribution/Collaboration

Github
(or other service -- bitbucket, aws, etc. )

Download and Merge
(pull)

Commit
changes1 +
changes2

You

Co-Author

# Other things Git does

- Roll-back to previous versions
- Branch development/management
- Integration in to lots of software
- Best way to explore: start using git

Today

- Sourcetree
  - Setup Repo
  - Clone Repo
  - Checkout
  - Commit
  - Pull

## Where to get help

- Easy help
  - Lots of places
  - Stackoverflow.com
  - Google
  - Github youtube channel
  - Sourcetree help

- Punching deck and interactive learning:
  - try.github.io
  - www.codeschool.com/courses/git-real
  - A great course at lynda.com www.lynda.com/Git-tutorials/Git-Essential-Training/100222-2.html

- Deep Dive
  - pro-git book by Scott Chacon and Ben Straub. – Free online http://git-scm.com/book/en/v2

Things Git is bad at

- Tracking binary files – word files, images, etc. It will track them, but it's not ideal

# Merge Conflicts

- Git is good at fixing conflicts
- When it can't you need to fix them
- Diff, resolve using 'mine'/'theirs'

## Problem:

- Software versions change over time

# That's actually sort of a hard question to answer

- Virtualization software, but not totally
- But totally if on pc/mac