# Lecture 3: Processing Linguistic Data, Git/GitHub

LING 1340/2340: Data Science for Linguists

Na-Rae Han

# Objectives

▸ To-do 1: What linguistic data did you find?

▸ HW1: What did you process?

▸ GitHub: completing the fork triangle

▸ DataCamp tutorials

▸ Tools:

  ◆ Git and GitHub

  ◆ Jupyter Notebook

**You should be taking NOTES!**

# First thing to do every class

# To-do #1

▶ What linguistic data sets did you look at?

◆ Corpus data?

◆ Non-corpus data?

▶ What makes a dataset a corpus?

# Back to Class-Exercise-Repo

https://github.com/Data-Science-for-Linguists-2021/Class-Exercise-Repo
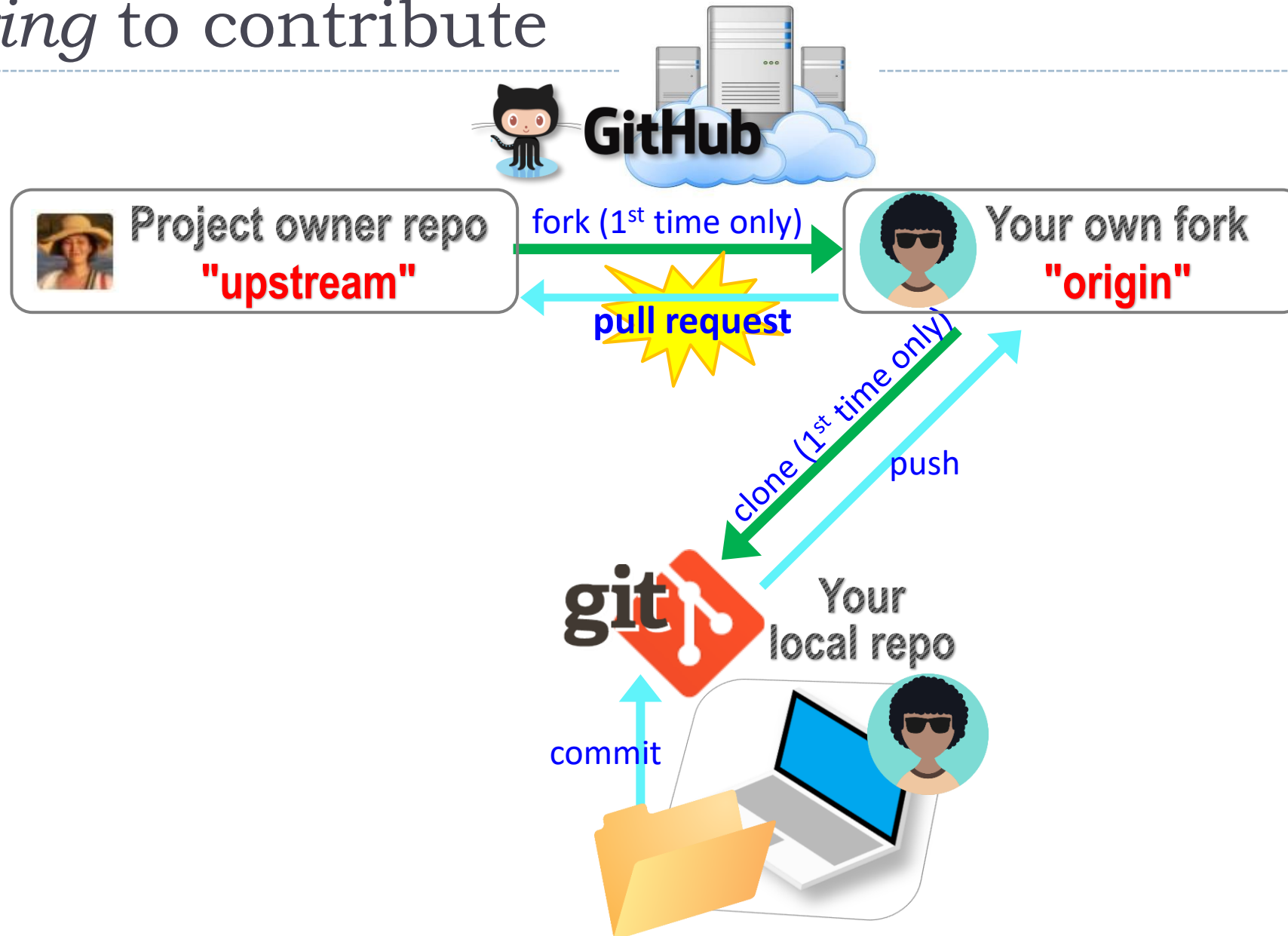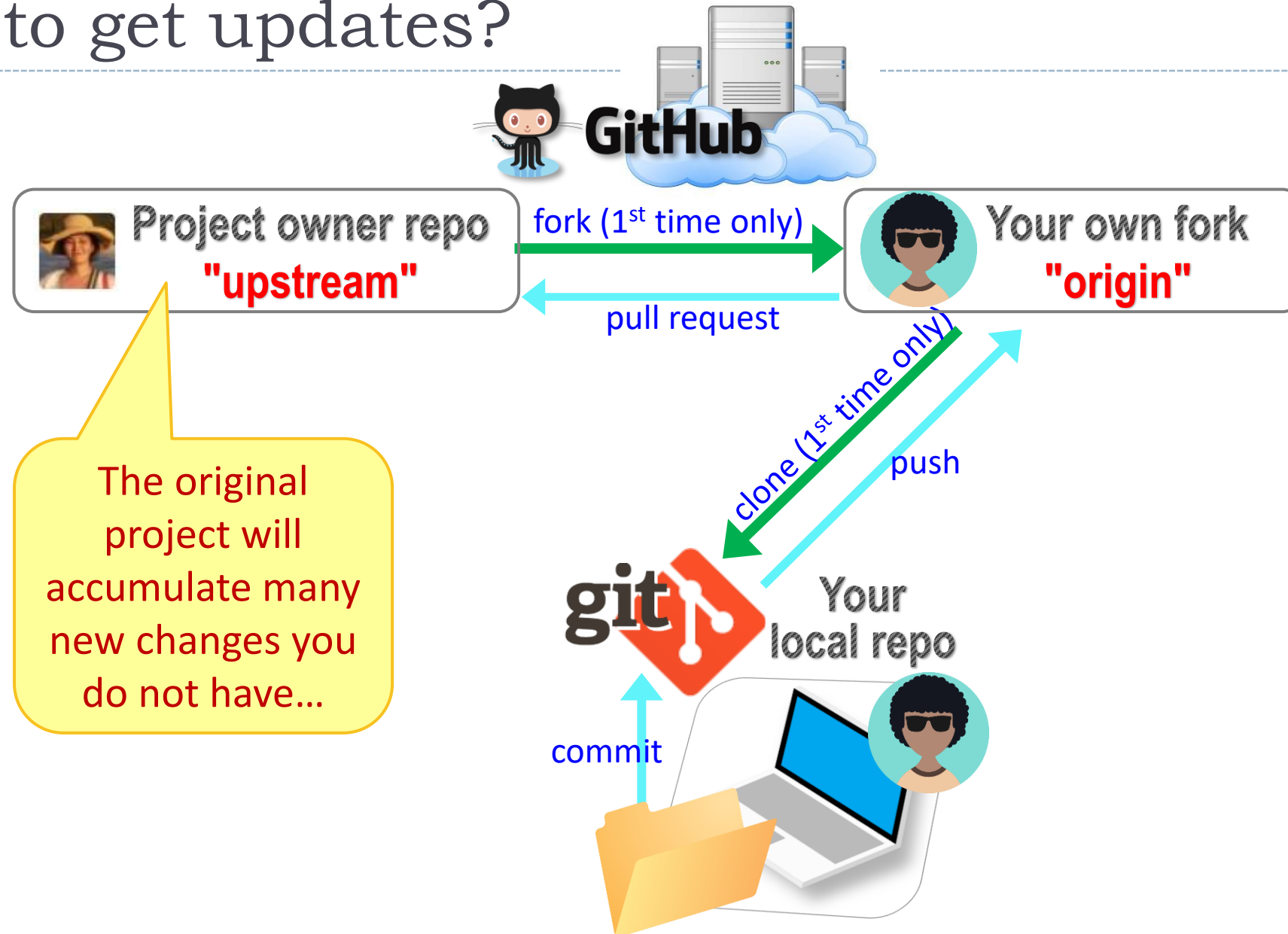
▶ Todo1

    ◆ Your To-do 1 submissions


▶ Lots of files -- I have merged in everyone's contributions.

▶ **But! Your own fork does not have those.**
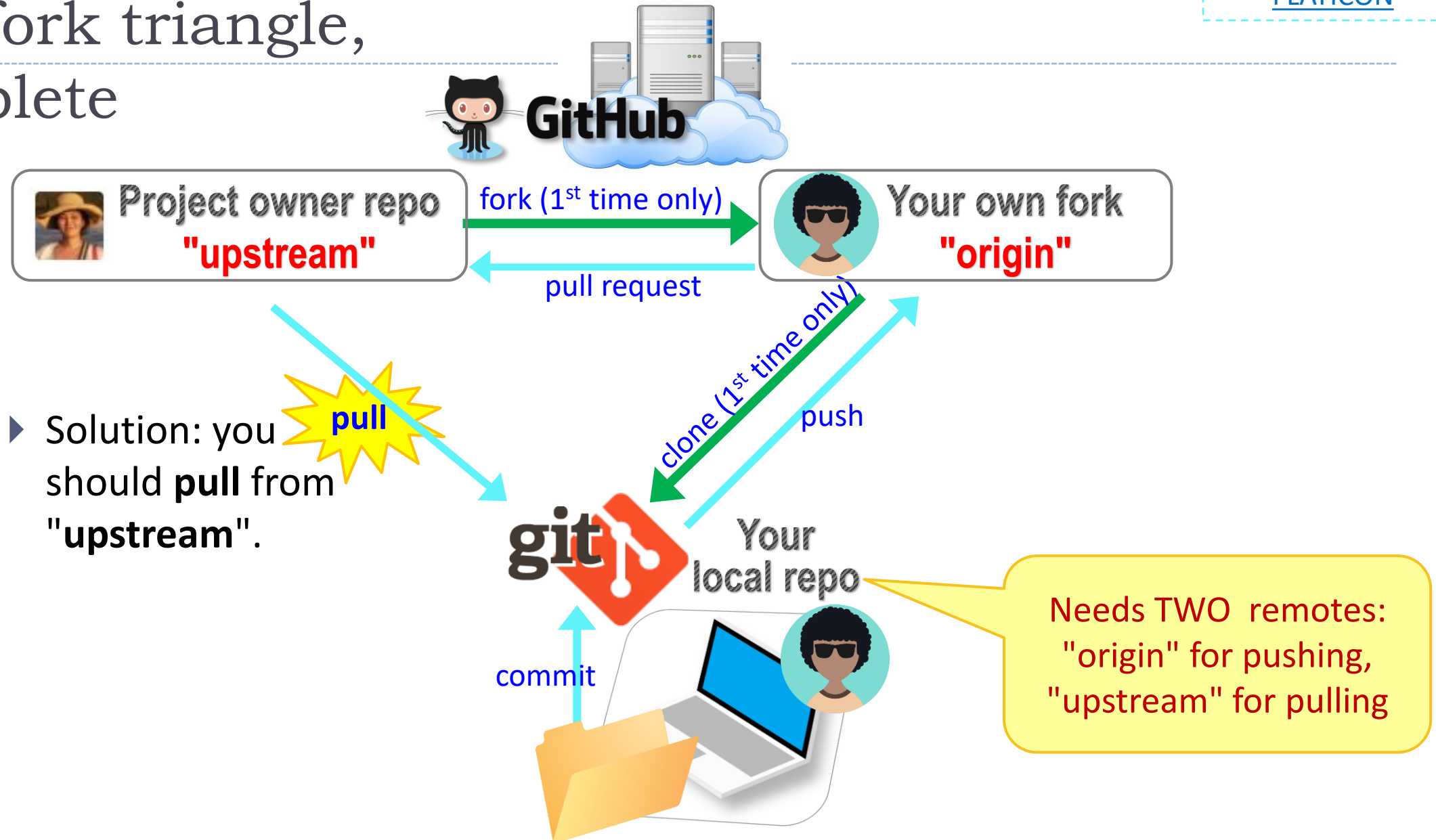
# *Offering* to contribute

**GitHub**

Project owner repo
**"upstream"**

fork (1st time only)

Your own fork
**"origin"**

pull request

clone (1st time only)

push

**git**

Your
local repo

commit

# How to get updates?

The original project will accumulate many new changes you do not have...

# The fork triangle, complete

GitHub

Project owner repo **"upstream"** — fork (1st time only) → Your own fork **"origin"**

← pull request

▸ Solution: you should **pull** from "**upstream**".

**pull**

clone (1st time only) / push

git — Your local repo

commit

Needs TWO remotes: "origin" for pushing, "upstream" for pulling

# Keeping your fork up-to-date

▸ The original repo ("upstream") will keep changing.

    ◆ How to keep your copies (GitHub fork and local repo) up-to-date?

▸ Cloning already configured your GitHub fork as "origin":
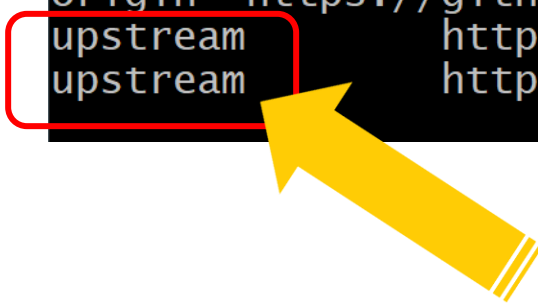
```
narae@T480s MINGW64 ~/Documents/Data_Science/Class-Exercise-Repo (main)
$ git remote -v
origin  https://github.com/narae-student/Class-Exercise-Repo.git (fetch)
origin  https://github.com/narae-student/Class-Exercise-Repo.git (push)
```
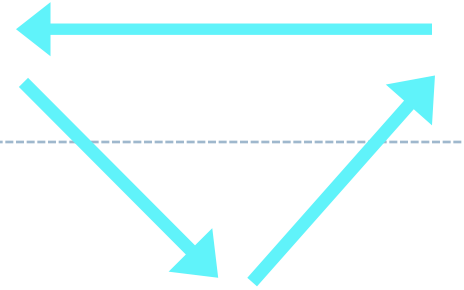
▸ Configure the original repo as another remote: "upstream"

    ◆ `git remote add upstream <GitHub-repo-URL.git>`

▸ When it's time to sync, pull from upstream:

    ◆ `git pull upstream main`

▸ Pushing should be done to your GitHub fork ("origin").

    ◆ `git push origin main` ⟵ You might be able to leave out "origin main".

# Two remotes: "origin", "upstream"

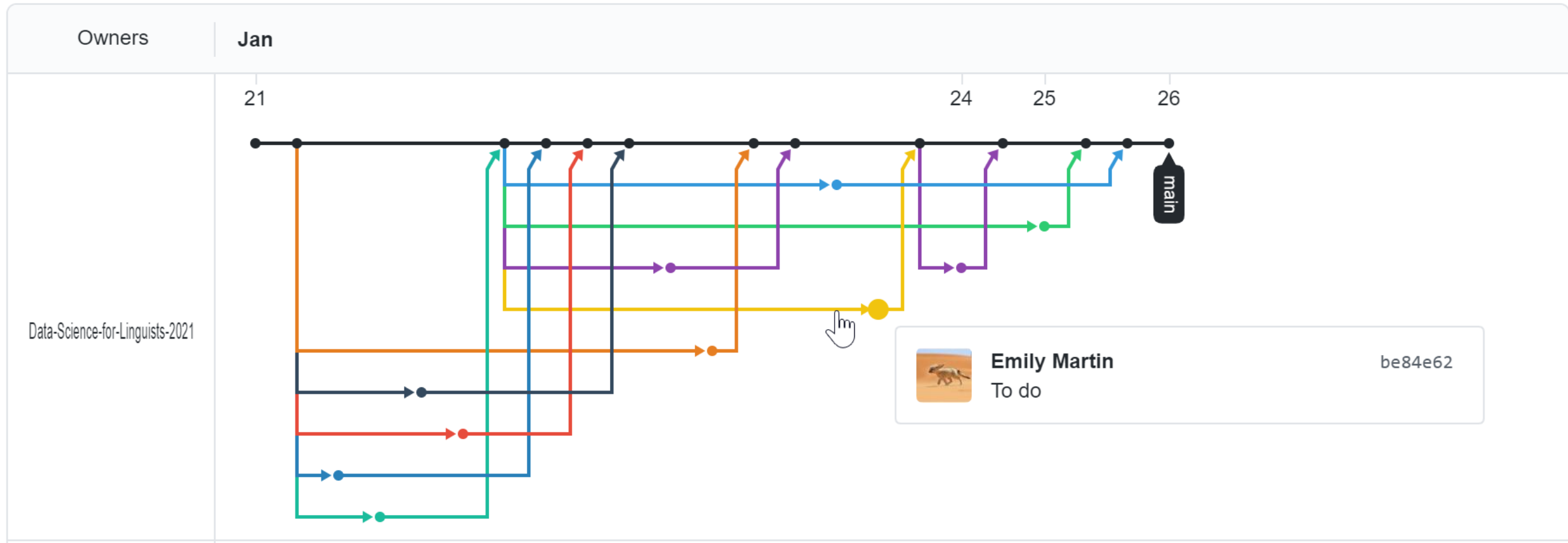# The fork triangle: workflow

▶ ## On your **laptop**

1.  Check your local repo's status: `git status`. Get it to a clean state.
2.  Pull from "upstream", syncing your local repo: `git pull upstream main`. Your local repo now has all latest changes.
    - If there is a merge conflict, you will need to resolve it. (fingers crossed)
3.  Do your work! New files, edits, etc.
4.  Do your usual local Git routine: `git add` and `git commit`.
5.  Push new versions to your own GitHub fork ("origin"): `git push origin main`

▶ ## On **GitHub**

1.  Check your forked repo. It should have your new work.
2.  Create a **pull request** for the original repo ("upstream") owner.
3.  Give it some time, and check back on the status of your pull request.

# Many forks and merges

- https://github.com/Data-Science-for-Linguists-2021/Class-Exercise-Repo/network

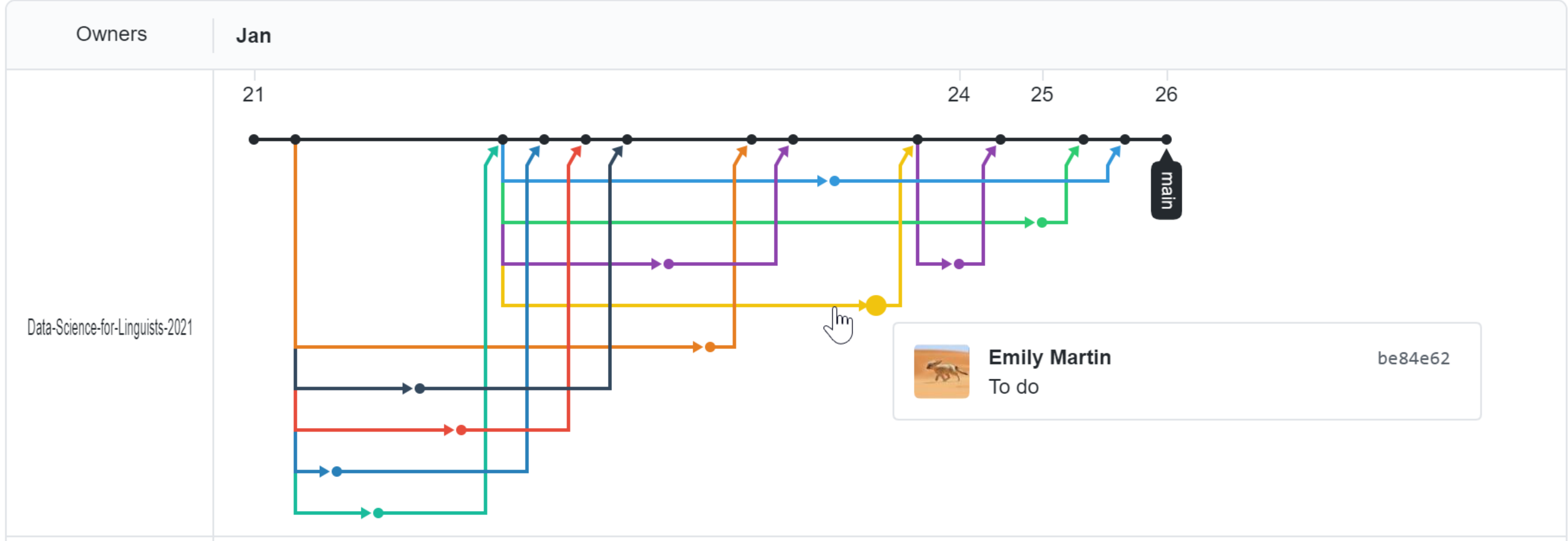

# HW1: processing pull request, merging

▸ With everyone working on their own files/folders, merging is conflict-free:

# HW1: sync your HW1-Repo

1. Configure "upstream" remote:

   `git remote add upstream https://github.com/Data-Science-for-Linguists-2021/HW1-Repo.git`

2. Pull from upstream:

   `git pull upstream main`

3. Push to your GitHub fork:

   `git push origin main`

| | |
|---|---|
| Everyone's repos are synced. | Now, everyone has everyone's homework submission. |

# HW1: Review

▸ What did you all work on?

▸ You wish list: what new skills would you like to learn?

▸ What is the `.gitignore` file?

▸ Why did we exclude data files from Git?

▸ What is up with that "your_file_here.txt" blank file? What is `git rm`?

▸ Jupyter Notebook: do you like it?

# Git and GitHub are complicated.

▶ They are powerful tools.

▶ There are a lot of abstract, high-level concepts involved.

▶ Concepts do not make sense before you get hands-on.

▶ You cannot get hands-on without the right context.

← We will learn slowly, learning various pieces as we go.

← You need to be patient, careful and methodical. Make sure you don't rush, and follow instructions.

# Git and GitHub are complicated.

▸ We will follow some ground rules.

▸ DO NOT EDIT A REPOSITORY'S CONTENT THROUGH GITHUB.



▸ Don't accidentally commit a file! Be mindful of what you add.  Avoid using:
  ◆ `git add .`
  ◆ `git add *`

▸ For now, do not **delete** or **re-name** any previously committed file.
  ◆ If you must: use `git rm` and `git mv`.

# Course Group on DataCamp

▸ Video-based, interactive tutorials



We get FREE access this semester -- all you can learn! Use **Pitt email address** to sign up.

# How to use DataCamp



- ▶ Topics for the next couple of weeks:
  - ◆ numpy library
  - ◆ pandas library
  - ◆ visualization libraries such as `matplotlib`

- ▶ The video tutorials are linked as "assignments"
  - ◆ Great learning resource, but not mandatory.
  - ◆ They *complement* the textbook nicely.

- ▶ Online exercise interface needs some getting used to.
  - ➔ next slide

https://campus.datacamp.com/courses/intro-to-python-for-data-science/chapter-2-python-lists?ex=7

# Your text editor in shell

▸ You should be able to launch your text editor from shell and create a new text file in the directory.



```
narae@T480s MINGW64 ~
$ cd Documents/Data_Science/

narae@T480s MINGW64 ~/Documents/Data_Science
$ which atom
/c/Users/narae/AppData/Local/atom/bin/atom

narae@T480s MINGW64 ~/Documents/Data_Science
$ atom newfile.txt

narae@T480s MINGW64 ~/Documents/Data_Science
$ ls
Class-Exercise-Repo/  languages/  newfile.txt

narae@T480s MINGW64 ~/Documents/Data_Science
$
```

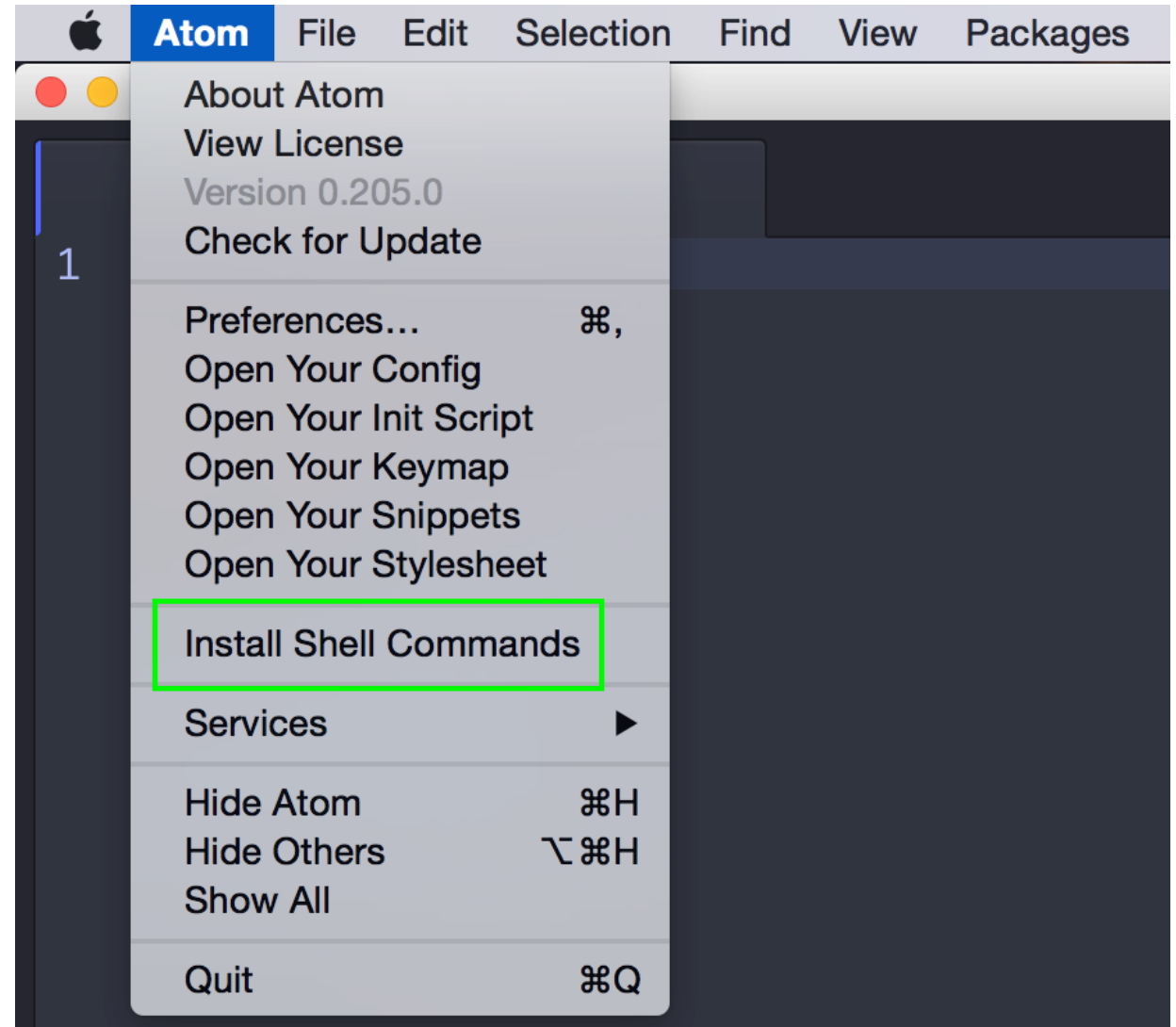Your command-line environment knows where atom is installed

Atom launches in a new window. I type in some stuff and save file.

New file has been created!

# Mac users: configure Atom for shell

▸ https://stackoverflow.com/questions/22390709/how-to-open-atom-editor-from-command-line-in-os-x

    ▸ "Install Shell Commands"

    ▸ After this, you can launch atom directly from your Terminal (bash or zsh shell).

# Mac only: Bash vs. Zsh

‣ **Windows** folks are using Git-bash, which has nice colorized Git output

‣ **Mac**: new default shell is zsh, older versions will have bash

　◆ In your terminal, execute echo $0



Bash shell (older)

Zsh shell (newer, recommended)

# Wrapping up

▶ **To-do #2 is out: due Thu.**

  ◆ Study numpy, make your own study notes as JNB. Submit via Class-Exercise-Repo.

▶ **Try out DataCamp tutorials!  Especially the numpy chapter.**

▶ **Learn:**

  ◆ Git, GitHub

  ◆ Jupyter Notebook

  ◆ numpy

  ◆ pandas