

GitHub Introduction and Integrating GitHub into your Project Workflow



THE WORLD BANK
IBRD • IDA | WORLD BANK GROUP
Development Economics • Impact



Intro to Git and GitHub

María Reyes Retana

July 15, 2024

Development Impact Evaluation (DIME)
The World Bank

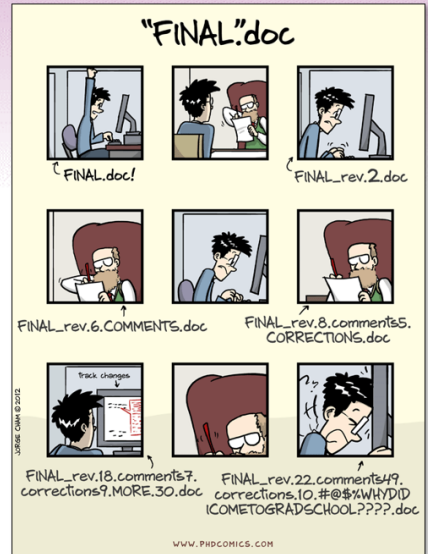
Before the session starts:

1. Do you have a GitHub.com account? If not, go to <https://github.com/join> and sign up
2. Have you shared your GitHub username with the session organizer?
3. Have you installed GitHub Desktop? If not go to <https://desktop.github.com/> and download it.
4. Have you logged in at least once on GitHub Desktop? If not open GitHub Desktop and log in using your GitHub account.
5. Have you been invited to <https://github.com/dime-wb-trainings/lyrics-jul15> and <https://github.com/dime-wb-trainings/GitHub-MockProject>?
6. And have you accepted at <https://github.com/dime-wb-trainings/lyrics-jul15/invitations> and <https://github.com/dime-wb-trainings/GitHub-MockProject/invitations>?

Introduction

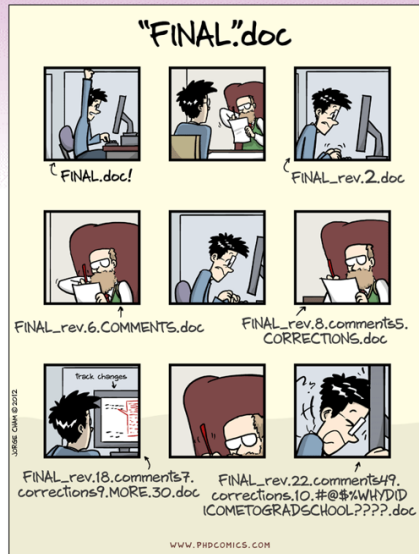
What is Git used for?

- Git solves the *Final.doc* problem



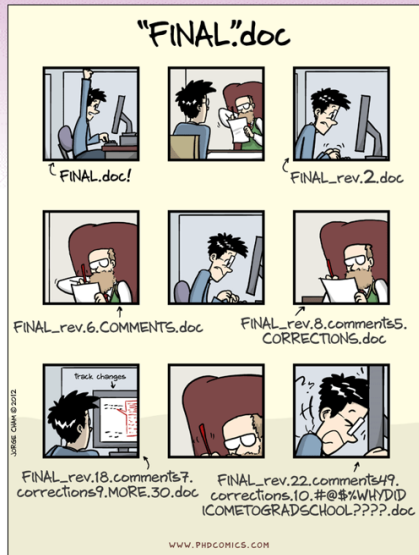
What is Git used for?

- Git solves the *Final.doc* problem
- Common solution to the *Final.doc* problem:
Name all your docs like
YYMMDD_docname_INITIALS.doc



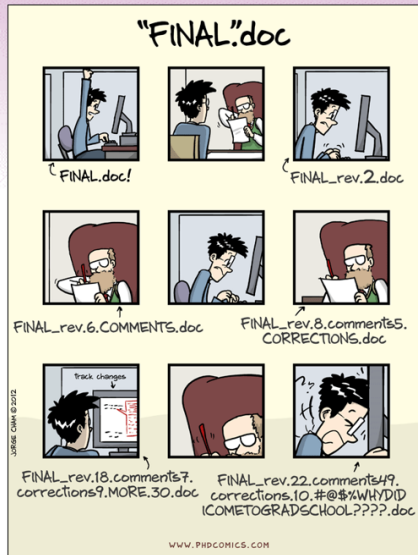
What is Git used for?

- Git solves the *Final.doc* problem
- Common solution to the *Final.doc* problem:
Name all your docs like
YYMMDD_docname_INITIALS.doc
- Git tracks *YYMMDD* and *INITIALS* for all
edits without the user having to remember it

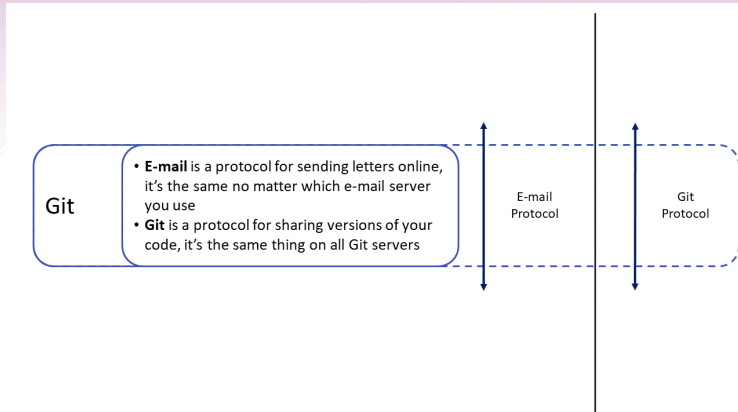


What is Git used for?

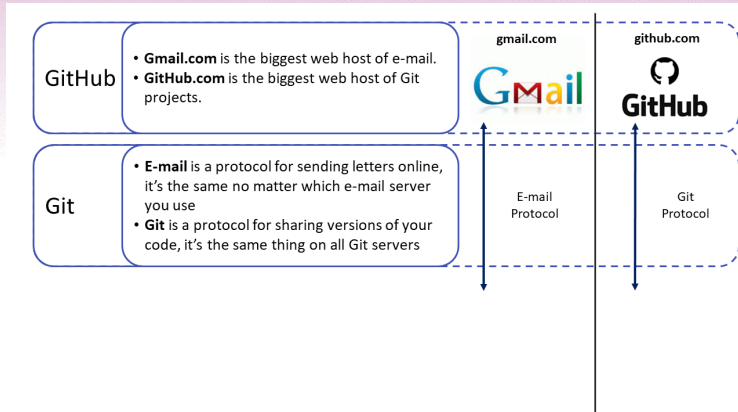
- Git solves the *Final.doc* problem
- Common solution to the *Final.doc* problem:
Name all your docs like
YYMMDD_docname_INITIALS.doc
- Git tracks *YYMMDD* and *INITIALS* for all edits without the user having to remember it
- That's far from everything, Git also solves:
 - Conflicting copy problem (DropBox etc.)
 - I can't re-produce my Baseline report problem
 - Who wrote this code 4 years ago and why?
 - And much much more...



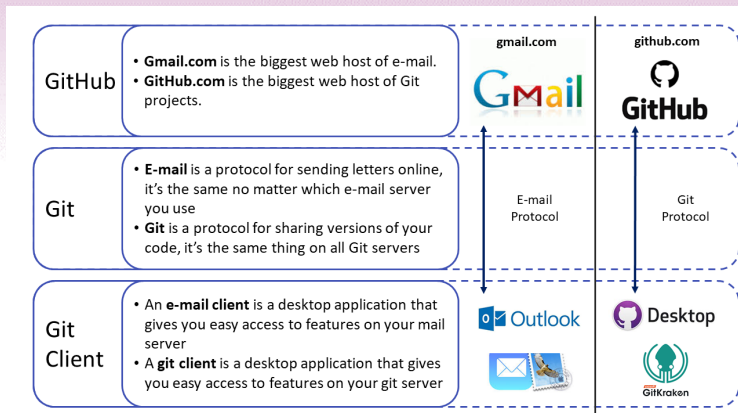
What is Git, GitHub and GitHub Desktop?



What is Git, GitHub and GitHub Desktop?



What is Git, GitHub and GitHub Desktop?



What will we learn?

In **An intro to Git and GitHub - Contributor Role** you will learn to:

- Explore the history of a project folder in GitHub and see what different team members are currently working on
- Download a project folder from GitHub so you can work on it
- Create a space in the project folder where you can make your edits
- Make edits and share those versions with your team. When you are ready, request that your edits are included in the main version

Three Git concepts needed to do this:

- Clone
- Commit
- Branch

Code free training!

No code today!

We will not work with code today.

Code tends to distract people if, for example, they see a command they do not understand.

Instead we will work with lyrics in .txt files that works exactly the same as code files in Git.

How to browse GitHub.Com

Your project folder is called a **repository** in Git, often **repo** for short. When you enter <https://github.com/dime-wb-trainings/lyrics-jul15> you get to what we will call the **repository landing page**.

Go from **GitHub.com** to **repo**

1. From anywhere on github.com click the *octocat* icon in the top left corner.
2. In the menu to your left you see the repositories you are invited to
3. Click any repo to get to the landing page of that repo.

Go from **repo** to **landing page**

1. Click the repo name in [dime-wb-trainings/lyrics-jul15](https://github.com/dime-wb-trainings/lyrics-jul15) at the top of any page within the repo

Clone

What is cloning?

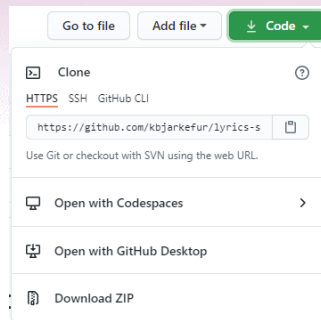
Cloning is similar to downloading a **repository** to your computer.

The difference between cloning and downloading is that **when Git clones a repository it remembers where you downloaded it from**. This is necessary so that Git knows where send any edits you make to the files when sharing them with your team.

How do I clone a repository from GitHub?

How to clone a repository:

1. Go to the **landing page** of <https://github.com/dime-wb-trainings/lyrics-jul15>
2. Click the green *Code* button (see image)
3. Click *Open with GitHub Desktop*
4. Select where on your computer to clone the repository. Do **NOT** clone in a shared folder, like a network drive or in DropBox. Create a *GitHub* folder in non-synced location and clone there. Read more about this in our guide [here](#).



Explore the clone

Explore the clone!

Compare the files and folders you cloned to your computer with those in the repository on

<https://github.com/dime-wb-trainings/lyrics-jul15>

Collaboration on a repository

Collaboration needs two more concepts

In order to collaborate on a repository we need to introduce two topics:

Commits

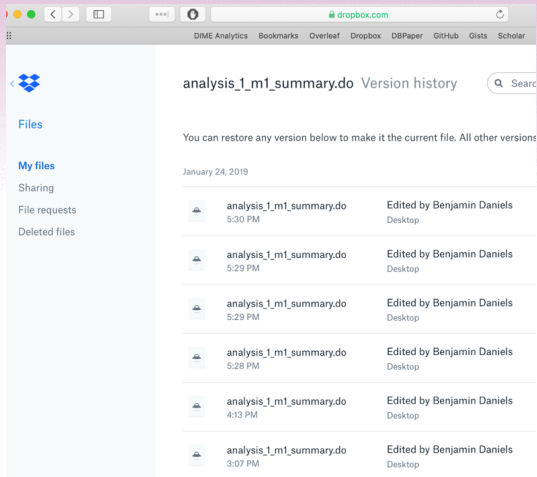
Branches

Commit

What is a version in version control?

First look at how version control works in another software you might be familiar with.

In DropBox each saved version of a file is saved to the version history. This is the only way to do it automatically, but are all these versions meaningful differences?



What is a commit?

Instead of having a list of each saved version of a file, in Git we use **commits to indicate what is each meaningful difference between two versions of our project folder.**

Each commit is a snapshot of all files in the project folder and lists how that snapshot differs from the previous snapshot (the previous commit).

Each commit has a time stamp and tracks who did the commit. This is very similar to the *YYMMDD_docname_INITIALS.doc* solution to the *Final.doc* problem.

How to make a commit

We need to introduce *branches* before we can all commit to the same repository, so for now, let me show you how to make a commit:

1. I add a new lyrics .txt file in the clone
2. I use GitHub desktop to commit the new file to the repository
3. Can you see the new file on your computer?
4. Can you see it if you sync in GitHub Desktop?

Exploring commits

Now that we know what a commit is, we can explore how the <https://github.com/dime-wb-trainings/lyrics-jul15> repository was created.

We will see a list of commits, that at first sight is similar to the the version history in DropBox, but **in Git the version list is more meaningful, as it is a list of only meaningful differences.**

- <https://github.com/dime-wb-trainings/lyrics-jul15/commits>
- This list can also be found in GitHub Desktop in the *History* tab

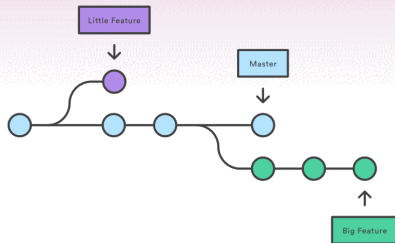
Branches

Introducing branches

Branches is the "killer feature" of Git. This is where Git becomes really powerful as a collaboration tool and not just as version control.

Branches allows you to **create a copy of the code where you can experiment**, if you like the result, **you can very easily merge your experiment with the main version of your code.**

This non-linear version control is much more similar to how we actually work than the strictly linear version control in, for example, DropBox



Looking at branches

One more way to explore the repository:

- Linear progression ->
<https://github.com/dime-wb-trainings/lyrics-jul15/commits>
- Non-linear progression ->
<https://github.com/dime-wb-trainings/lyrics-jul15/network>

Exploring branches

- You can change branch in */commits*. What happens then?
- Go to the landing page, what happens if you change the branch here?
- Which version is in the clone on your computer? They are all actually in your clone, but only one is shown - **checked out** - at the time
- What happens to the files on your computer when you check out another branch in GitHub Desktop?

Working with branches

A typical Git work flow involves multiple branches and there are tools in GitHub to makes that work flow easy, but that is not within today's scope. Although, what you should know after this training is only how to create your own branch and how to commit to it.

Create a branch:

- Go to <https://github.com/dime-wb-trainings/lyrics-jul15> and click the button where it says *Branch: main*.
- Write your name in the field and click *Create branch: your_name*. Make sure it says *from 'main'*.
- See how the button now says *Branch: your_name*

Combining Commit & Branch

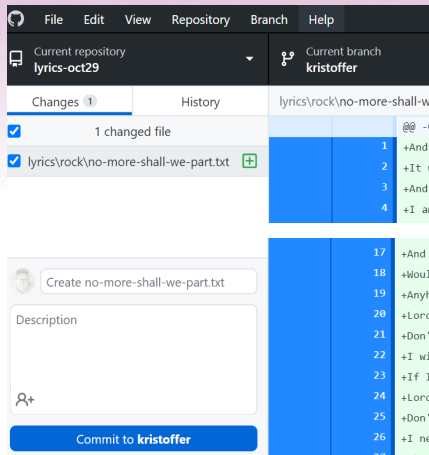
Now it is time to collaborate

Now it is time for you to collaborate:

1. Make sure your branch is checked out in GitHub Desktop.
2. Open a text editor. It could be *Notepad* if you are using Windows, *TextEdit* if you are using Mac, or any other code editor like *Atom*, or *Notepad++* etc.
3. Google the lyrics of your favorite song, and copy the lyrics to a new file in the text editor you just opened
4. Save the lyrics in your local clone according to these instructions:
 - Save the file in *.txt* format (Especially Mac users, this is not always the default)
 - Name the file after the title of the song
 - Save it in the appropriate genre folder (create a new genre folder if needed)

Do your first commit

1. Open the changes tab in GitHub Desktop
2. GitHub Desktop tracks your clone and has noticed that you changed something in it
3. Then you need to do the three steps required to commit a file to the repository:
 - 3.1 Make sure the file you want to add is checked
 - 3.2 Write a commit message
 - 3.3 Click *Commit to your_name* - check out your branch if it says *main*
 - 3.4 Click the sync button



Check your contribution

Check your commit on GitHub:

- Can you find your commit at:
<https://github.com/dime-wb-trainings/lyrics-jul15/network>
- Can you find your commit at:
<https://github.com/dime-wb-trainings/lyrics-jul15/commits>

If you cannot see your commit, make sure that you are looking in the correct branch.

Pull Requests

A feature related to branches is a **pull request**. When the edits you have done in your branch are ready to be merged with the main version of the code, you make a pull request, i.e. you are requesting that your edits are pulled into the main branch.

A pull request can be made either by the person that edited the branch or the repo maintainer.

It is common in GitHub repositories that only the repo maintainer has access to the main branch, and pull requests are then the only way to contribute to the main branch.

Merge your contribution

To make a pull request for your branch:

- Go to <https://github.com/dime-wb-trainings/lyrics-jul15/pulls> and click *New pull request*
- Make sure that the *main* branch is selected as the *base:* branch, and then select your branch as the *compare:* branch
- Scroll down to check the edits you are requesting to be pulled in to the main branch. If it looks ok, then click *Create pull request*
- Then you have the chance to add more instructions if you want, then click *Create pull request* again

Merge your contribution

Can you see your lyrics file in the main branch now? Why not?

Your contribution will not be included until the branch is merged. We have more trainings on the details of merging branches but for now, this is all you need to know:

- Always have someone to review your PR before merging it
- Always delete your branch after it is merged - you can always recreate a new branch with the same name if you want

What have we learned?

In **An intro to Git and GitHub - Contributor Role** you have learned to:

- Explore the history of a project folder in GitHub and see what different team members are currently working on
- Download a project folder from GitHub so you can work on it
- Create a space in the project folder where you can make your edits
- Make edits and share those versions with your team. When you are ready, request that your edits are included in the main version

What have we learned?

In **An intro to Git and GitHub - Contributor Role** you have learned to:

- Explore the history of **repository** and see what different team members are currently working on
- **Clone** a **repository** so you can work on it
- Create a **branch** in the **repository** where you can make your edits
- Make edits and share **commits** with your team. When you are ready, make a **pull request** to the **main branch**

Integrating GitHub into Your Project Workflow

Before We Start the Second Part of the Session

- Now that you know the GitHub basics, we will learn how to integrate it with our project workflow and explore some best practices.
- Make sure you have access to this repository:
<https://github.com/dime-wb-trainings/GitHub-MockProject>.
- We will work on different elements of a project with a good structure. Follow along with the exercises.
- If you get lost, all the instructions are in the repository's README file to guide you.

Importance of Structuring Data Work

- **Saves Time:**
 - Taking time to structure data work at the beginning of a project will save time over the project's full life-cycle.
- **Effective Collaboration:**
 - Facilitates task delegation for senior team members and helps junior team members understand and communicate tasks to their supervisors.
- **Transparent Analytical Decisions:**
 - Organizing data work and setting up effective "inter-generational" communication ensures transparency and ease of understanding.
- **Error Reduction:**
 - A good structure from the beginning makes work clearer and less error-prone, simplifying the implementation of necessary changes.

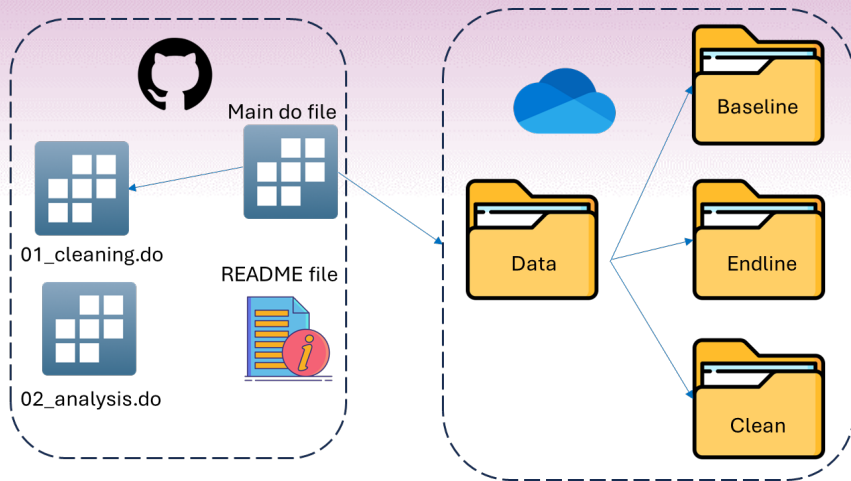
Key Takeaways

In this part of the session, we will cover:

1. **Transparent, Version-Controlled Data Work:** How to use GitHub to manage code files.
2. **Effective Collaboration:** The importance of a standardized folder structure for non-code files.
3. **Implement Reproducibility from the Start:** Create a main script to run all data tasks.
4. **Document Your Work:** Use READMEs and organize all research documentation.

Organizing your data work using these principles will make your research more effective, efficient, and reproducible!

Structure of a good workflow

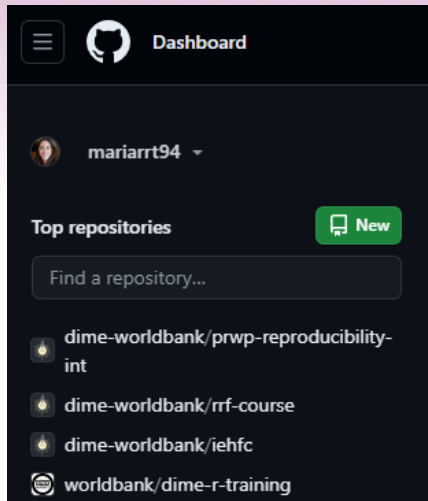


Example of a Good Workflow

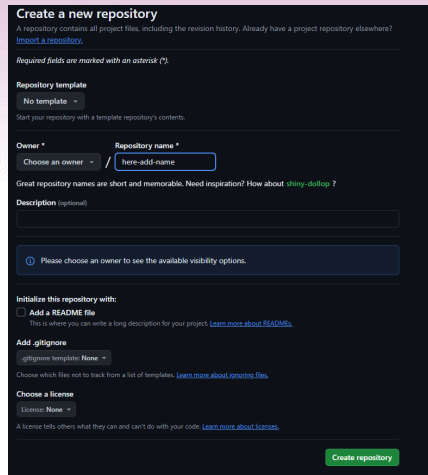
- **Steps for Repository Creation**

1. Go to GitHub and log in to your account.
 2. Click on the "New" button to create a new repository.
 3. Enter a name for your repository.
 4. Set the repository visibility to private if necessary.
 5. Click "Create repository".
- Note: In this case, the repository has already been created, as it normally will in a regular project. This demonstration is for reference.

GitHub: repository creation



Step 1: Create Repository



Step 2: Repository Details

Hands-on: cloning the mock project

Once your project has a repository.

- **Clone the repository**

1. Go to the GitHub repository:

<https://github.com/dime-wb-trainings/GitHub-MockProject>.

2. Click on the green "Code" button and select "Open with GitHub Desktop".
3. Follow the prompts to clone the repository to your local machine.

Hands-on: cloning the mock project

The screenshot shows the GitHub interface for a repository named 'GitHub-MockProject'. The repository is marked as 'Private'. At the top, there are buttons for 'Edit Pins', 'Watch' (0), and 'Code'. The 'Code' button is circled in blue. Below the repository name, there are tabs for 'main', '1 Branch', and '0 Tags'. A search bar 'Go to file' is present. The file list shows several files: 'Code', '.gitignore', 'GitHub-Intro-Workflow-Training.pdf', 'README.md', and 'github-session-jakarta.tex'. The 'README' file is selected. A dropdown menu is open from the 'Code' button, showing options for 'Local' and 'Codespaces'. Under 'Local', there are three options: 'Clone', 'SSH', and 'GitHub CLI'. The 'SSH' option is circled in blue. Below these options, there is a text input field containing the URL 'https://github.com/dime-wb-trainings/GitHub-MockProject'. At the bottom of the dropdown, there is an option 'Open with GitHub Desktop' which is also circled in blue. Other options in the dropdown include 'Download ZIP'.

GitHub-MockProject Private

Edit Pins Watch 0

main 1 Branch 0 Tags

Go to file

Add file <> Code

Local Codespaces

Clone ?

HTTPS SSH GitHub CLI

<https://github.com/dime-wb-trainings/GitHub-MockProject>

Clone using the web URL.

Open with GitHub Desktop

Download ZIP

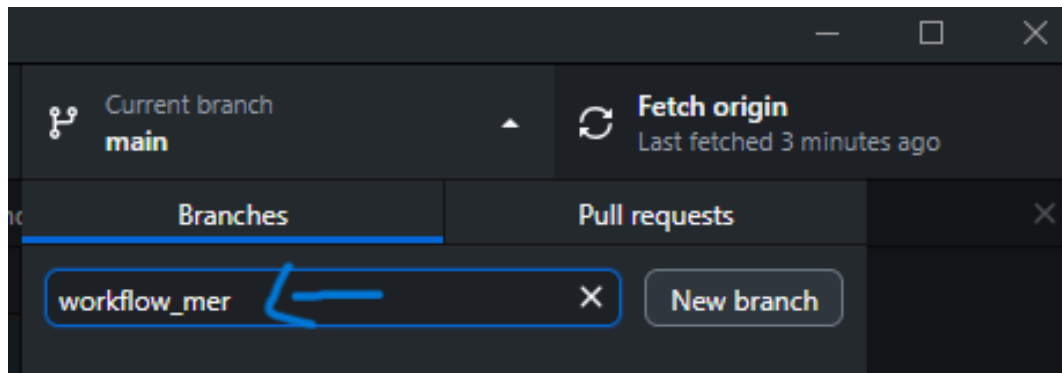
Structure a Good Workflow - GitHub Training Exercise

Hands-on: create a branch in the mock project

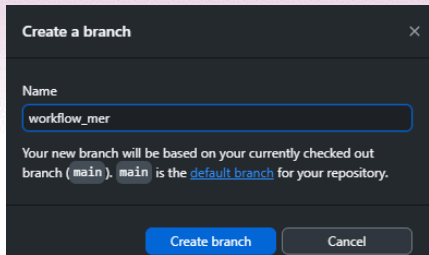
Once you have the repository on your computer:

- **Create a branch**

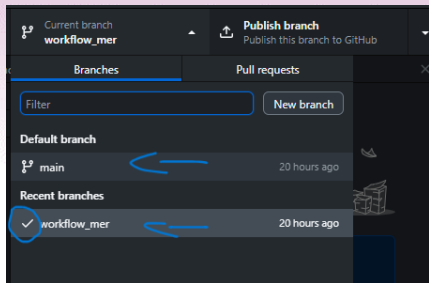
1. As we will be working collaboratively, create a branch named "workflow_" followed by your initials.
2. Switch to that branch to start making changes to the project.



Hands-on: create a branch in the mock project



After you hit new branch, this pop-up will appear.



After you create a branch, GH will move you to that branch, but you can also move between branches.

Branches in a project

When working on a real project, follow the principle: **branch often, merge often.**

In practice:

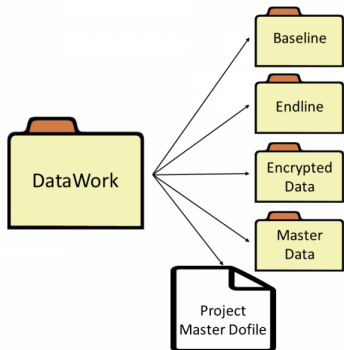
- Create a branch for each task.
- After completing the task, merge the branch back to the main branch by creating a Pull Request (PR).
- Repeat this process for subsequent tasks to ensure smooth and collaborative development.

Folder structure

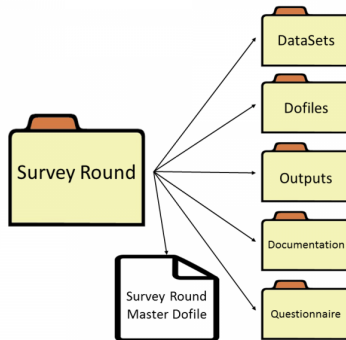
- At the beginning of a research project, organize your data work into a **folder structure that is easy to follow**
- Think about all the sources and rounds of data that your project might have
- Have one separate sub-folder for each source of data, and keep sub-folders for each data source as similar as possible
- Avoid creating individual folders in an ad-hoc manner

A folder structure template

Project Folder:



Data-set/round Folder:



- DIME Analytics created a stata command to automate creation of a standardized folder structure, as well as a template master (main) .do file - `iefolder` - part of the `ietoolkit` package
- Standardized folder structures across our portfolio creates efficiency gains
- `iefolder` was designed for primary data but the same concept can be generalized to secondary data
- You do not need to use `iefolder`, but do think through what folder structure will work for your project and establish and maintain that structure from day 1

Some notes on folder structure

- It's important to use the exact same folder structure on the cloud (where data sits) and GitHub (where code sits)
- Using a similar structure in all your projects will make it easier to move across projects

Hands-on: Folder Structure

- Download the mock data from the provided link.
- Save the data files to the desired location on your local machine.
- See the two roots of your project: `data/` and `code/`.

Hands-on: Folder Structure

- Explore the folder structure:
 - `data/` - Folder for data files (not tracked by Git).
 - `raw/` - Raw data files.
 - `intermediate/` - Intermediate data files.
 - `analysis/` - Data files for analysis.
 - `code/` - Folder for all Stata scripts.
 - `cleaning/` - Scripts for data cleaning.
 - `construction/` - Scripts for data construction.
 - `analysis/` - Scripts for data analysis.
 - `outputs/` - Folder for your outputs, if relevant.
 - `README.md` - Project documentation.
 - `.gitignore` - Specify files and folders to ignore in Git.

- **Main script:** Runs all other scripts in the intended order.
- **Saves time and reduces errors:** Team members don't have to run each file every time code or data changes.
- **Automation is more efficient:** Numbering code files can communicate run order but is inefficient and error-prone.
- **Facilitates understanding and reproducibility:** Main scripts make it easy for all team members to understand and run the code, and to create a reproducibility package for publication.

Main script: additional benefits

- Acts as an ordered **table of contents**
- **Centralizes settings and requirements**
 - Install software packages
 - Set reproducibility parameters
 - Define project-wide variables
- **Facilitates Collaboration:**
 - Bridge between the code and folder structure.
 - Use globals to set root paths dynamically.
 - Simplifies project setup for others.

Main script

Main script

```
* Project folder globals
* -----

global dataWorkFolder      "$projectfolder/DataWork"

*iefolder*1*FolderGlobals*master*****
*iefolder will not work properly if the line above is edited

global mastData            "$dataWorkFolder/MasterData"

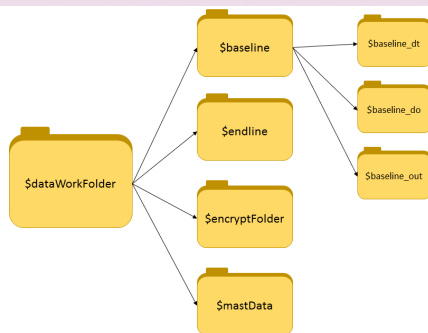
*iefolder*1*FolderGlobals*rawData*****
*iefolder will not work properly if the line above is edited

global encryptFolder       "$dataWorkFolder/EncryptedData"
global masterIdDataSets    "$encryptFolder/IDMasterKey"

*iefolder*1*RoundGlobals*rounds*baseline*****
*iefolder will not work properly if the line above is edited

*baseline folder globals
global baseline            "$dataWorkFolder/baseline"
global baseline_dt         "$baseline/DataSets"
global baseline_do         "$baseline/Dofiles"
global baseline_out        "$baseline/Output"
```

Folder structure



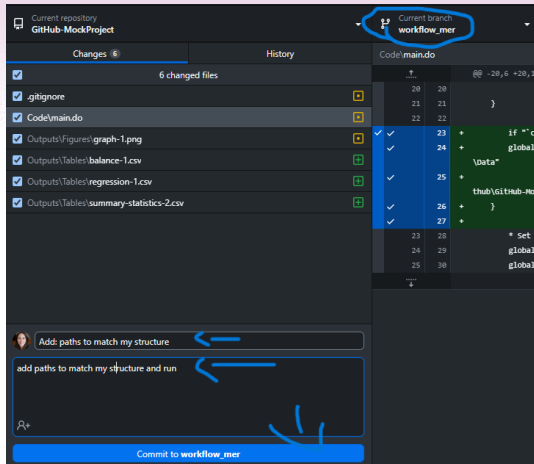
Each significant folder gets a global/object that can be referenced across all scripts

Hands-on: Main Script Setup

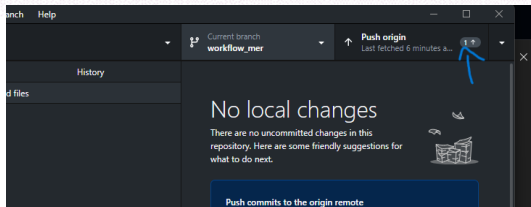
1. **Open the `main.do` file in the mock project folder** (make sure you are in your own branch).
2. Make the necessary modifications in the `main.do` file to **match your project structure**.
 - Add the paths to match your computer and structure.
 - See how the global paths are set dynamically.
3. Use GitHub Desktop to **commit** your changes.
4. **Push/Publish** your changes into github.com

Hands-on: Main Script Setup

Commit changes made in your local machine



Push changes to github.com



Hands-on: Main Script Setup

- Note how this script acts as a bridge between your data in OneDrive and your code.
- To collaborate, a new user would only need to modify the path lines. Everything else would remain the same (assuming you all have the data in OneDrive with the same structure).

- **Quick Project Overview**

- Provides a summary of the project's purpose and objectives.
- Helps new contributors quickly understand the project's scope and goals.

- **Setup Instructions**

- Guides new team members on how to set up the project locally.
- Lists necessary software, dependencies, and installation steps.

- **Documentation of Key Decisions**

- Records important decisions made during the project.
- Provides context for the project's development and structure.

- **Instructions for Use**

- Details how to run and use the project.
- Includes examples and explanations of key functionalities.

- **Improves Collaboration**

- Facilitates onboarding of new team members.
- Ensures consistency and clarity in communication.

- **Explore README.md File**

- Note: The README.md file is already there. You don't need to create.

Other Elements: `.gitignore`

- **Usefulness of `.gitignore` file**

- Prevents tracking of sensitive or unnecessary files.
- Keeps the repository clean and focused.
- Avoids conflicts from environment-specific files.
- Binary files are not well tracked in GitHub.

- **Explore `.gitignore` file**

- Note: The `.gitignore` file is already there. You don't need to create it, but you might need to modify it for your project.

Useful links: GitHub

- All DIME Analytics GitHub trainings: <https://osf.io/e54gy/>
- Other DIME Analytics GitHub resources:
<https://github.com/worldbank/dime-github-trainings>. For example:
 - DIME Analytics GitHub Templates (for example .gitignore):
<https://github.com/worldbank/dime-github-trainings/tree/master/GitHub-resources/DIME-GitHub-Templates>
- Markdown cheat sheet (how to format text on GitHub.com):
<https://www.markdownguide.org/cheat-sheet/>

Useful links: Workflow

- Template Main do file: <https://github.com/worldbank/dime-data-handbook/blob/main/code/box-2-4-stata-master-dofile.do>
- Template README
https://github.com/worldbank/wb-reproducible-research-repository/blob/main/resources/README_Template.md.
- Example package and README.
<https://reproducibility.worldbank.org/index.php/catalog/148>.
- Presentation on Structure Datawork. <https://osf.io/pa8y5>.