

Guide for Accelerating Computational Reproducibility in the Social Sciences



BERKELEY INITIATIVE FOR TRANSPARENCY
IN THE SOCIAL SCIENCES

ACRE Team

2020-12-19

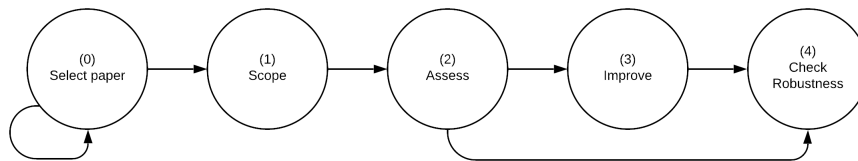
Contents

	5
Introduction	7
Beyond binary judgments	8
Stages of the exercise	8
Reproduction Strategies	11
1 Selecting a paper	13
1.1 From candidate to declared paper	13
1.2 Identify your relevant timeline.	15
2 Scoping	19
2.1 Read and summarize the paper	19
2.2 Record scope of the exercise	20
2.3 Setup your own revised reproduction package.	21
3 Assessment	23
3.1 Describe the inputs.	24
3.2 Connect display items to all its inputs	27
3.3 Assign a reproducibility score.	29
4 Improvements	35
4.1 Display item improvements	35
4.2 Paper-level improvements	39
5 Checking for Robustness	41
5.1 Feasible robustness checks: increasing the number of feasible spec- ifications	43
5.2 Reasonable robustness check: justifying and testing.	45
6 Concluding the Reproduction	47
6.1 Outputs	47
6.2 Anonymity and data sharing	48

7	Guidance for a Constructive Exchange Between Reproducers and Original Authors	51
7.1	For reproducers contacting the authors of the original study . . .	52
7.2	For original authors Rrsponding to requests from reproducers . .	61
7.3	Harassment and/or discrimination	63
8	Examples of Reproduction Trees	65
8.1	Stylized examples	65
8.2	Examples from real reproduction attempts	68
9	Tips and Resources for Reproducible Workflow	129
9.1	Reproducible workflows:	129
9.2	Links to resources, organizations and people for reproducible work	130
10	Contributions	133
10.1	Contributing feedback on these guidelines	133
10.2	List of Contributors: Guidelines content and source code:	133
10.3	Suggested citation format	134
10.4	Acknowledgments	134
11	Definitions	135
11.1	Concepts in reproducibility	135
11.2	Concepts in the ACRE exercise and the platform	138

See a full list of contributors

A beta version of the **Social Science Reproduction Platform** is now available! **Sign up here** if you would like to be part of our beta testing in the spring.



(0) Select

(1) Scoping

(2) Assessment

(3) Improvement

(4) Robustness

Display-item-level

Paper-level

Select paper

Read paper

Describe inputs

+ Raw data

+ Version control

Analytical choices

Search materials

Identify claims

Reproduction diagrams

- + Analysis data
- + Documentation
- Type of choice
- Check SSRP
- Declare estimates
- Reproduction score
- + Analysis code
- + Dynamic document
- Choice value
- Declare/Discard paper
- + Cleaning code
- + File structure
- Justify and test alternatives
- Declare estimates
- Debug analysis code
- Debug cleaning code

This work is licensed under the Creative Commons Attribution-NonCommercial 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/4.0/>.

Introduction

Computational reproducibility is defined as the degree to which it is possible to obtain consistent results using the same input data, computational methods, and conditions of analysis (of Sciences, 2019). In 2019, the American Economic Association updated its Data and Code Availability Policy to require that the AEA Data Editor verify the reproducibility of all papers before they are accepted by an AEA journal. Similar policies have been adopted in political science, particularly at the *American Journal of Political Science*. In addition to the requirements laid out in such policies, specific recommendations were produced by data editors of social science journals to facilitate compliance. This change in policy is expected to improve the computational reproducibility of all published research going forward, after several studies showed that rates of computational reproducibility in the social sciences range from somewhat low to alarmingly low (Galiani et al., 2018; Chang and Li, 2015; Kingi et al., 2018).

Replication, or the process by which a study’s hypotheses and findings are re-examined using different data or different methods (or both) (King, 1995) is an essential part of the scientific process that allows science to be “self-correcting.” *Computational reproducibility*, or the ability to reproduce the results, tables, and other figures of a paper using the available data, code, and materials, is a necessary condition for replication. Computational reproducibility is assessed through the process of *reproduction*. At the center of this process is the *reproducer* (you!), a party rarely involved in the production of the original paper. Reproductions sometimes involve the *original author* (whom we refer to as “the author”) in cases where additional guidance and materials are needed to execute the process.

This Guide is meant to be used in conjunction with the **Social Science Reproduction Platform** (SSRP), an open-source platform that crowdsources and catalogs attempts to assess and improve the computational reproducibility of published social science research. Though in its current version, the Guide is principally intended for reproductions of published research in economics, it may be used in other social science disciplines, and we welcome contributions that aim to “translate” any of its parts to other social science disciplines (learn how you can contribute here). The purpose of this document is to provide a common approach, terminology, and standards for conducting reproductions. The goal

of reproductions, in general, is to assess and improve the computational reproducibility of published research in a way that promotes a better understanding of research and facilitates additional robustness checks, extensions, collaborations, and replications.

This Guide and the SSRP were developed as part of the Accelerating Computational Reproducibility in Economics (ACRE) project, which aims to assess, enable, and improve the computational reproducibility of published economics research. The ACRE project is led by the Berkeley Initiative for Transparency in the Social Sciences (BITSS)—an initiative of the Center for Effective Global Action (CEGA)—and Dr. Lars Villhuber, Data Editor for the journals of the American Economic Association (AEA). This project is supported by the Laura and John Arnold Foundation.

View slides used for the presentation “How to Teach Reproducibility in Classroom”

Beyond binary judgments

Assessments of reproducibility can easily gravitate towards binary judgments that declare an entire paper “reproducible” or “non-reproducible.” These guidelines suggest a more nuanced approach by highlighting two realities that make binary judgments less relevant.

First, a paper may contain several scientific claims (or major hypotheses) that may vary in computational reproducibility. Each claim is tested using different methodologies, presenting results in one or more display items (outputs like tables and figures). Each display item will itself contain several specifications. Figure 1 illustrates this idea.

Second, for any given specification there are several levels of reproducibility, ranging from the absence of any materials to complete reproducibility starting from raw data. And even for a specific claim-specification, distinguishing the appropriate level can be far more constructive than simply labeling it as (ir)reproducible.

Note that the highest level of reproducibility, which requires complete reproducibility starting from raw data, is very demanding to achieve and should not be expected of all published research — especially before 2019. Instead, this level can serve as an aspiration for the field of economics at large as it seeks to improve the reproducibility of research and facilitate the transmission of knowledge throughout the scientific community.

Stages of the exercise

A reproduction attempt is divided into five stages, corresponding to the first five chapters of these guide:

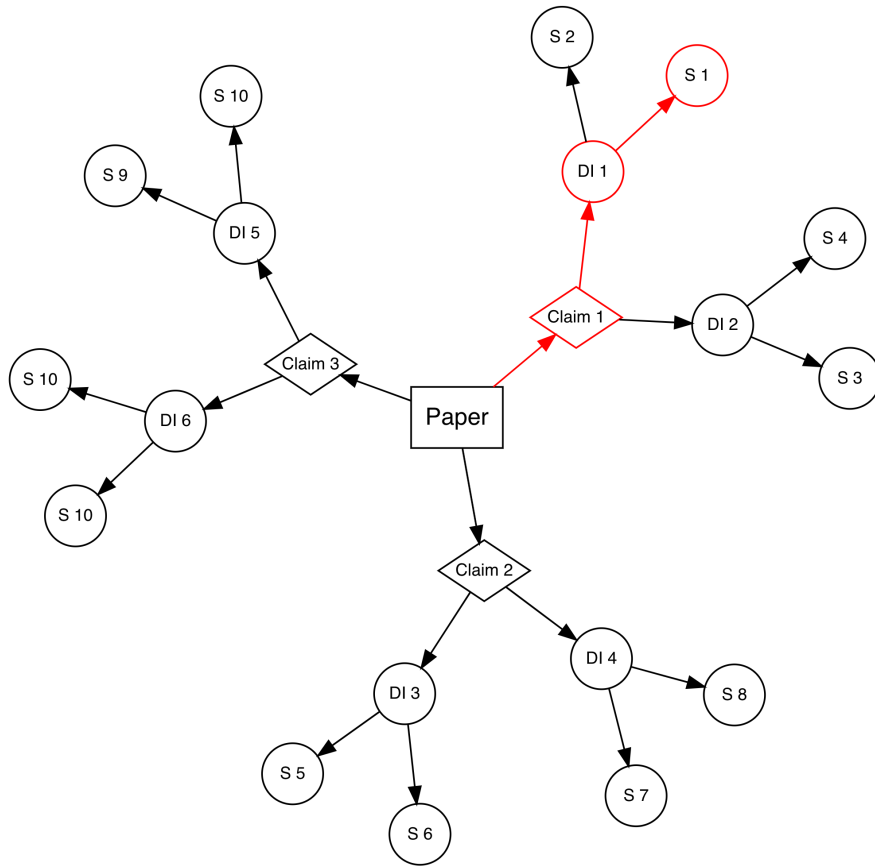


Figure 1: One paper has multiple components to reproduce.
 DI: Display Item, S: Specification

0. **Paper selection**, where you (the reproducer) will select a candidate paper and verify the availability of a reproduction package. Depending on availability, you will declare the paper and start the exercise, or select a new candidate paper (after leaving a short record);
1. **Scoping**, where you will define the scope of the exercise recording which claims, display items, and specifications you will focus for the remainder of the exercise;
2. **Assessment**, where you will review and describe in detail the available reproduction package, and assess the current level of computational reproducibility of the selected display items;
3. **Improvement**, where you will modify the content and/or the organization of the reproduction package to improve its reproducibility;
4. **Robustness checks**, where you will identify feasible robustness checks and/or assess the reasonableness of specific variations in analytical choices.

These guidelines do not include a possible fifth stage of **extension**. Here you may extend the current paper by including new methodologies or data. If you were to extend the same methodology and research question into a different sample, that would bring you closer to a *replication*.

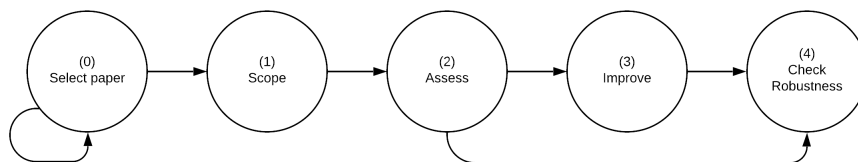


Figure 2: Four stages of a reproduction attempt

This process need not be chronologically linear. For example, you may realize that the scope of a reproduction is too ambitious and switch to a less intensive one. Later in the exercise, you can also begin testing different specifications for robustness while also assessing a paper's level of reproducibility. The only stage that should go first, and may not be edited further once finished, is the scoping stage, as it defines the scope of the exercise.

This guide, and the SSRP platform, will change the key unit of analysis as you progress through each stage. As Figure 3 shows, the scoping stage will be centered around the key scientific claims selected for reproduction. Once those have been identified, the next two stages will guide you on how to assess and improve the reproducibility of the display items supporting those claims. In the final stage the unit of analysis is once again at the claim level.

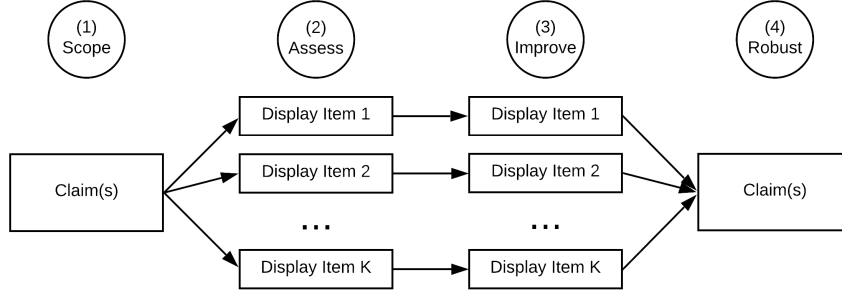


Figure 3: Relevant unit of analysis at each stage of a reproduction attempt

Reproduction Strategies

Generally, a reproduction will begin with a thorough reading of the study being reproduced. However, subsequent steps may follow from a *reproduction strategy*. For example, a reproduction may closely follow the order of the steps outlined above. This might entail the reproducer first choosing a set of results whose production they are interested in assessing or understanding, completely reproducing these results to the extent possible, and then making modifications to the reproduction package. Another potential strategy could be for the reproducer to develop potential robustness checks or extensions while reading the study, which would lead to the definition of a set of results to be assessed via reproduction. Yet another reproduction strategy may be for the reproducer to seek out a paper that uses a particular dataset to which they have access or an interest in using, reproducing the results that use that dataset as an input, then probing the robustness of the results to various data cleaning decisions.

The various uses of reproduction makes the number of potential reproduction strategies quite large. In choosing or designing a reproduction strategy, it is helpful to clearly identify the goal of the reproduction. In all of the examples laid out in the paragraph above, the order in which the steps of the reproduction exercise are taken is at least partially determined by what the reproducer hopes to get from the exercise. The structure provided in these guidelines, together with a clear reproduction goal, can facilitate the implementation of an efficient reproduction strategy.

Chapter 1

Selecting a paper

The goal of this stage is to help you define the scope of your exercise by declaring a paper and the specific output(s) on which you will focus. You might first consider multiple papers without analyzing them more closely (we refer to these as **candidate papers**) before moving forward with your **declared paper**.

The main difference between a candidate and a declared paper is the availability of a reproduction package. A **reproduction package** is the collection of materials that make it possible to reproduce a paper. This package may contain data, code, or documentation. If you are unable to independently locate the reproduction package for your paper, you can ask the paper's author for it (find guidance on this in Chapter 5) or simply choose another candidate paper. If you still want to explore the reproducibility of a paper with no reproduction package, these guidelines provide instructions for requesting materials from authors to create a public reproduction package, or if this proves unsuccessful, for building your reproduction package from scratch.

To avoid duplicating the efforts of others who may be interested in reproducing one of your candidate papers, **we ask that you record your candidate papers in the SSRP database** (currently under development).

Note that in this stage, *you are not expected to review the reproduction materials in detail*, as you will dedicate most of your time to this in later stages of the exercise.

1.1 From candidate to declared paper

At this point of the exercise, you are *only validating the availability* of (at least) one reproduction package and not assessing the quality of its content. Follow the steps below to verify that a reproduction package is available, and stop whenever you find it (this may mean that you have found your declared paper).

1. Check whether previous reproduction attempts have been recorded in the SSRP Database for the paper (more on the SSRP Database in the next section).
2. Check the journal or publisher's website, looking for materials named "Data and Materials," "Supplemental Materials," "Reproduction/Replication Package/Materials," etc.
3. Look for links in the paper (review the footnotes and appendices).
4. Review the personal websites of the paper's author(s).
5. Contact the author(s) to request the reproduction package using this email template. In this and future interactions with authors, we encourage you to follow our guidance outlined in Chapter 5.
6. Deposit the reproduction package in a trusted repository (e.g., Dataverse, Open ICPSR, Zenodo, or the Open Science Framework) under the name **Original reproduction package for - Title of the paper**. You will be asked to provide the URL of the repository in Survey 1.

In case you need to contact the authors, make sure to *allocate sufficient time for this step* (we suggest at least three weeks before the date you plan to start the reproduction). Instructors should also plan to accordingly (e.g., if the ACRE exercise is expected to take place in the middle of the semester, students should review candidate papers and (if applicable) contact the authors in the first few weeks of the semester).

Review the decision tree (Figure 1.1) below for a more detailed overview of this process. Remember, *if at any step of the process you decide to abandon the paper, make sure to record the candidate paper in the ACRE database* before moving on to another candidate paper. Once you have obtained the reproduction package, the *candidate paper* becomes your *declared paper* and you can move forward with the exercise! Do not invest time in doing a detailed read of any paper until you are sure that it is your declared paper.

1.1.1 Candidate paper entries in the SSRP Database

If the SSRP database contains previous reproduction attempts of the paper, you will see a report card with the following information:

Box 1: Summary Report Card for ACRE Paper Entry
Title: Sample Title
Authors: Jane Doe & John Doe
Original Reproduction Package Available: Yes (link)/No. [If "No"] **Contacted Authors?:** Yes/No
 [If "Yes(contacted)"] **Type of Response:** Categories (6).
Additional Reproduction Packages: Number (eg., 2)
Authors Available for Further Questions for ACRE Reproductions: Yes/No/Unknown

If after taking steps 1-5 above (or for some other reason) you are unable to locate the reproduction package, record your candidate paper (and if applicable, the outcome of your correspondence with the original authors) in the SSRP database following the example above.



Figure 1.1: Decision tree to move from candidate to declared paper

1.2 Identify your relevant timeline.

Before you begin working on the four main stages of the reproduction exercise (Scoping, Assessment, Improvement, and Robustness), it is important to manage your own expectations and those of your instructor or advisor. Be mindful of your time limitations when defining the scope of your reproduction activity. These will depend on the type of exercise chosen by your instructor or advisor and may vary from a weeklong homework assignment, to a longer class project that may take a month to complete or a semester-long project (an undergraduate thesis, for example).

Table 1 shows an example distribution of time across three different reproduction formats. The Scoping and Assessment stages are expected to last roughly the same amount of time across all formats (lasting longer for the semester-long activities, and acknowledging that less experienced researchers, such as undergraduate students, may need more time). Differences emerge in the distribution of time for the last two main stages: Improvements and Robustness. For shorter exercises, we recommend avoiding any possible improvements to the raw data (or cleaning code). This will limit how many robustness checks are possible (for example, by limiting your ability to reconstruct variables according to slightly different definitions), but it should leave plenty of time for testing different specifications at the analysis level.

2 weeks (~10 days)

1 month (~20 days)

1 semester (~100 days)

analysis data

raw data

analysis data

raw data

analysis data

raw data

Scoping

10% (1 day)

5% (1 day)

5% (5 days)

Assessment

35%

25%

15%

Improvement

25%

0%

40%

20%

30%

Robustness

25%

5%

25%

25%

Chapter 2

Scoping

The goal of this stage is to help you define the scope of your reproduction by identifying key scientific claims and the specific display items on which you will focus.

Once you have identified your declared paper, get familiarized with it and choose the specific display items on which you will focus for the remainder of the exercise.

2.1 Read and summarize the paper

Depending on how much time you have, we recommend that you write a short (1-2 page) summary of the paper. This will help remind you of the key elements to focus on for the reproduction, and demonstrate your understanding of the paper (for yourself and others like your instructor or advisor).

When reading or summarizing the paper, try to answer the following questions:

- Would you classify the paper's scientific claims as mainly focused on estimating a causal relationship, estimating/predicting a descriptive statistic of a population, or something else?
- How many scientific claims (descriptive or causal) are investigated in the paper?
- What is the population for which the estimates apply?
- What is the population that is the focus of the paper as a whole?
- What are the main data sources used in the paper?
- How many display items are there in the paper (tables, figures, and inline results)?
- What is the main statistical or econometric method used to examine each claim?

- What is the author’s preferred specification (or yours, if the authors are not clear)?
- What are some robustness checks for the preferred specification?

2.2 Record scope of the exercise

By now you should have a fairly good understanding of the paper’s content. You do not, however, need to have spent any time reviewing the reproduction package in detail.

At this point, you should clearly specify which part of the paper will be the main focus of your reproduction. Focus on specific estimates, represented by a unique combination of claim-display item-specification as represented in figure 1. If you plan to scope more than one claim, *we strongly recommend starting with just one* and recording your results. You can then initiate another record in ACRE later for the second (or third, fourth, etc.) claim to reproduce using the materials and knowledge you developed in the first exercise. You can, however, reproduce more than one claim if you are already familiar with the paper.

In the Assessment stage, the reproduction will be centered around the display item(s) that contain the specification you indicate at this point.

Declare specific main estimates to reproduce.

Identify a scientific claim and its corresponding preferred specification, and record its magnitude, standard error, and location in the paper (page, table #, and table row and column). If the authors did not explicitly chose a particular estimate, you will be asked to select one. In addition to the preferred estimate, reproduce up to five estimates that correspond to alternative specifications of the preferred estimate.

Declare possible robustness checks for main estimates (optional).

After reading the paper, you might wonder why the authors did not conduct a specific robustness test. If you think that such analysis could have been done *within the same methodology* and *using the same data* (e.g., by including or excluding a subset of the data like “high-school dropouts” or “women”), please specify a robustness test that you would like to conduct before starting the Assessment stage.

These are the elements you will need for the Scoping stage. **You now have all the elements necessary to complete Survey 1.**

2.3 Setup your own revised reproduction package.

As part of the scoping stage, you have now identified the reproduction package of the original authors. In addition to these materials, we recommend that you create your own copy and called it a **revised reproduction package**. As you work through the next stages (assessment, improvement, and robustness) it is likely that you will modify some of the contents of the reproduction package. Keeping a record of those changes will help you as reproducers (eg. to document your assessments, or to communicate with original authors) and will likely help future reproducers allowing them to build on top of any contribution that they have created.

Please deposit the starting copy of your revised reproduction package (i.e. just the copy of the original reproduction package) in a trusted repository. Examples of trusted repositories include: Dataverse, openICPSR, Figshare, Dryad, Zenodo, Open Science Framework and others. We encourage reproducers to use version control software (Git) during their reproductions, and suggest they submit the entire repository (or link it to their GitHub repository).

Chapter 3

Assessment

In this stage, you will review and describe in detail the available reproduction materials, and assess levels of computational reproducibility for selected display items, as well as practices that increase the reproducibility of the overall paper. This stage is designed to record as much of the learning process behind a reproduction as possible to facilitate incremental improvements, and allow future reproducers to pick up easily where others have left off.

In the *Scoping* stage, you declared a paper, identified the specific claims you will reproduce, and recorded the main estimates that support the claims. In this stage, you will decide if you are interested in assessing the reproducibility of that entire display item (e.g., “Table 1”), or will assess only a pre-specified estimates (e.g., “rows 3 and 4 of Table 1”). Additionally, you can include other display items of interest.

The goal of this stage is to have an assessment of the current stage of the reproduction package, before you suggest any improvements. By the end of this section you will have a very specific description of the current stage of the reproduction package. With this you can try to improve its level of reproducibility and, potentially, approach the authors with specific questions to obtain missing materials or fix possible issues.

In the first section of this stage, you will provide a detailed description of the entire reproduction package. Second, you will connect the display items you’ve chosen to reproduce with their corresponding inputs. With these elements in place, you can score the level of reproducibility of each display item, and report on paper-level dimensions of reproducibility.

Tip: We recommend that you first focus on one specific display item (e.g., “Table 1”). After completing the assessment for this display item, you will have a much easier time assessing others.

3.1 Describe the inputs.

This section explains how to list the input materials found or referred to in the reproduction package. At this point it is hard to identify the materials that correspond to your selected claims and display items, so we recommend to list *all* the files in the reproduction packages. However if there are too many files and this step becomes intractable, you can focus on the materials required to reproduce a specific display item (this can be done by reviewing the code scripts in the section below).

First, you will identify data sources and connect them with their raw data files (when available). Second, you will locate and provide a brief description of the analytic data files. Finally, you will locate, inspect, and describe the analytic code used in the paper.

The following terms will be used in this section:

- **Raw data** – Unmodified data files obtained by the authors from the sources cited in the paper. Data from which personally identifiable information (PII) has been removed are *still considered raw*. All other modifications to raw data make it *processed*.
- **Analysis/Analytic data** – Data used as the final input in a workflow in order to produce a statistic displayed in the paper (including appendices).
- **Cleaning code:** A script associated primarily with data cleaning. Most of its content is dedicated to actions like deleting variables or observations, merging data sets, removing outliers, or reshaping the structure of the data (from long to wide, or vice versa).
- **Analysis code:** A script associated primarily with analysis. Most of its content is dedicated to actions like running regressions, running hypothesis tests, computing standard errors, and imputing missing values.

3.1.1 Describe the data sources and raw data.

In the paper you chose, find references to all *data sources* used in the analysis. A data source is usually described in narrative form. For example, if in the body of the paper, or in the appendix, you see text like “...for earnings in 2018 we use the Current Population Survey...”, the data source is “Current Population Survey 2018”. If it is mentioned for the first time on page 1 of the Appendix, its location should be recorded as “A1”. Do this for all the data sources mentioned in the paper. Each row represents a unique data source

Data sources also vary by unit of analysis, with some sources matching the same unit of analysis used in the paper (as in previous examples), while others are less clear (e.g., “our information on regional minimum wages comes from the

Table 3.1: Raw data information

data_source	page	data_files	directory	known_missing
"Current Population Survey 2018"	A1	cepr_march_2018.dta	data/	
"Provincial Administra- tion Reports"	A4	coast_simplepoint2.csv; rivers_simplepoint2.csv; RAIL_dummies.dta; rail- ways_Dissolve_Simplify_point2.csv	Data/maps/ data/to_clean/	
"2017 SAT scores"	4	Not available	data/to_clean/	
...

Bureau of Labor Statistics." This should be recorded as "regional minimum wages from the Bureau of Labor Statistics").

Next, look at the reproduction package and map the *data sources* mentioned in the paper to the *data files* in the available materials. Record their folder locations relative to the main reproduction folder¹. In addition to looking at the existing data files, we recommend that you review the first lines of all code files (especially cleaning code), looking for lines that call the datasets. Inspecting these scripts may help you understand how different data sources are used, and possibly identify any files that are missing from the reproduction package. Whenever a data source contains multiple files, enter them on the same cell, separated by semicolon (;).

If you cannot find the files or file name corresponding to a specific data source, enter "missing" in the filename column.

Record this information in this standardized spreadsheet (download it or make a copy for yourself), using the following structure:

Note: lists of files in the `data_files` and `known_missing` columns should have entries separated by a semi-colon to for the spreadsheet to be compatible with the ACRE Diagram Builder.

3.1.2 Describe the analytic data sets.

List all the analytic files you can find in the reproduction package, and identify their locations relative to the main reproduction folder. Record this information in the standardized spreadsheet.

¹a relative location takes the form of `folder_in_rep_materials/sub_folder/file.txt`, in contrast to an absolute location that takes the form of `/username/documents/projects/repros/folder_in_rep_materials/sub_folder/file.txt`

Table 3.2: Analysis data information

analysis_data	location	description
final_data.csv	analysis/fig1/	data for figure1
all_waves.csv	final_data/v1_april/	data for region-level analysis
...

Table 3.3: Code files information

file_name	location	inputs	outputs	description	primary_type
output_table1.csv	code/analysis/analysis_data/all_part1	code/analysis/analysis_data/all_part1	code/analysis/analysis_data/all_part1	produces first part of table 1 (unformatted)	analysis
data_cleaning02.R	code/cleaning/admin_01raw	analysis_data/all_waves	analysis_data/all_waves	removes outliers and missing vals from raw admin data	cleaning
...

As you progress through the exercise, add to the spreadsheet a one-line description of each file’s main content (for example: **all_waves.csv** has the simple description **data for region-level analysis**). This may be difficult in an initial review, but will become easier as you go along.

The resulting report will have the following structure:

3.1.3 Describe the code scripts.

List all code files that you found in the reproduction package and identify their locations relative to the master reproduction folder. Review the beginning and end of each code file and identify the inputs required to successfully run the file. Inputs may include data sets or other code scripts that are typically found at the beginning of the script (e.g., **load**, **read**, **source**, **run**, **do**). For each code file, record all inputs together and separate each item with “;”. Outputs may include other datasets, figures, or plain text files that are typically at the end of a script (e.g., **save**, **write**, **export**). For each code file, record all outputs together and separate each item with “;”. Provide a one-line description of what each code file does. Record all of this information in the standardized spreadsheet, using the following structure:

As you gain an understanding of each code script, you will likely find more inputs and outputs – we encourage you to update the standardized spreadsheet. Once finished with the reproduction exercise, classify each code file as *analysis* or *cleaning*. We recognize that this may involve subjective judgment, so we suggest that you conduct this classification based on each script’s main role.

Note: If a code script takes multiple inputs and/or produces multiple outputs they should be listed as semicolon separated lists in order to be compatible with the ACRE Diagram Builder.

3.2 Connect display items to all its inputs

Using the information collected above, you can trace your display items to their primary sources. Upload the standardized spreadsheets from above (sections 3.1.1, 3.1.2 and 3.1.3) to your reproduction at the SSRP and use the Diagram Builder. This will generate a reproduction diagram tree that represents the information available on the workflow behind a specific display item.

If the reproduction package does not organize the code around display items, you will be asked to identify all the outputs that contain the results used in a specific display item.

3.2.1 Complete workflow information

If you were able to identify all the relevant components in the previous section, the SSRP Diagram Builder will produce a tree diagram that looks similar to the one below.

```
table1.tex
  [code] analysis.R
    analysis_data.dta
      [code] final_merge.do
        cleaned_1_2.dta
          | [code] clean_merged_1_2.do
          | merged_1_2.dta
          | [code] merge_1_2.do
          | cleaned_1.dta
          | | [code] clean_raw_1.py
          | | raw_1.dta
          | cleaned_2.dta
          | [code] clean_raw_2.py
          | raw_2.dta
        cleaned_3_4.dta
          [code] clean_merged_3_4.do
            merged_3_4.dta
              [code] merge_3_4.do
                cleaned_3.dta
```

```

|      [code] clean_raw_3.py
|      raw_3.dta
cleaned_4.dta
|      [code] clean_raw_4.py
|      raw_4.dta

```

This diagram, built with the information you provided, is an important contribution to understanding the necessary components required to reproduce a specific display item. It also summarizes key information to allow for more constructive exchanges with original authors or other reproducers. For example, when contacting the authors for guidance, you can use the diagram to point out specific files you need. Formulating your request this way makes it easier for authors to respond and demonstrates that you have a good understanding of the reproduction package. You can also add this diagram into a new version of the readme file in your revised reproduction package.

3.2.2 Incomplete workflow information

In many cases, some of the components of the workflow will not be easily identifiable (or missing) in the reproduction package. Here the Diagram Builder will return a partial reproduction tree diagram.

For example, let's start from a simpler complete tree like the following:

```

table1.tex  analysis.R  analysis_data.dta  final_merge.do  cleaned_1_2.dta
clean_merged_1_2.do  merged_1_2.dta

```

For this case, if the file `final_merge.do` is missing from the previous diagram, the SSRP Diagram Builder will produce the following diagram:

```

table1.tex  analysis.R  analysis_data.dta
cleaned_1_2.dta  clean_merged_1_2.do  merged_1_2.dta

```

In this case, you can still manually combine this partial information with your knowledge from the paper and own judgement to produce a “candidate” tree diagram (which might lead to different reproducers recreating different diagrams). This may look like the following:

```

table1.tex  analysis.R  cleaned_1_2.dta  MISSING_CODE_FILE_1
cleaned_1_2.dta  clean_merged_1_2.do  merged_1_2.dta

```

To leave a record of the reconstructed diagrams, you will have to amend the input spreadsheets using placeholders for the missing components. In the example above, you should add the following entries to the code description spreadsheet:

As in the cases with complete workflows, these diagrams (fragmented or reconstructed trees) provide important information for assessing and improving the reproducibility of specific outputs. Reproducers can compare reconstructed trees and/or contact original authors with highly specific inquiries.

Table 3.4: Adding rows to code spreadsheet

file_name	location	inputs	outputs	description	primary_type
...
MISSING CODE	CODE FILE	cleaned_1_2	cleaned_1_2	missing code	unknown

For more examples of diagrams connecting final outputs to initial raw data, see [here](#).

3.3 Assign a reproducibility score.

Once you have identified all possible inputs and have a clear understanding of the connection between the outputs and inputs, you can start to assess the output-specific level of reproducibility.

Take note of the following concepts in this section:

- **Computationally Reproducible from Analytic data (CRA):** The output can be reproduced with minimal effort starting from the *analytic* datasets.
- **Computationally Reproducible from Raw data (CRR):** The output can be reproduced with minimal effort from the *raw* datasets.
- **Minimal effort:** One hour or less is required to run the code, not including computing time.

3.3.1 Levels of Computational Reproducibility for a Specific Output

Each level of computational reproducibility is defined by the availability of data and materials, and whether or not the available materials faithfully reproduce the output of interest. The description of each level also includes possible improvements that can help advance the reproducibility of the output to a higher level. You will learn in more detail about the possible improvements.

Note that the assessment is made *at the output level* – a paper can be highly reproducible for its main results, but suffer from low reproducibility for other outputs. The assessment includes a 10-point scale, where 1 represents that, under current circumstances, reproducers cannot access any reproduction package, while 10 represents access to all the materials and being able to reproduce the target outcome from the raw data.

- **Level 1 (L1):** No data or code are available. Possible improvements include adding: raw data (+AD), analysis data (+RD), cleaning code (+CC), and analysis code (+AC).

You will have detected papers that are reproducible at Level 1 as part of the Scoping stage (unsuccessful candidate papers). Make sure to take record them in Survey 1.

- **Level 2 (L2):** Code scripts are available (partial or complete), but no data are available. Possible improvements include adding: raw data (+AD) and analysis data (+RD).
- **Level 3 (L3):** Analytic data and code are partially available, but raw data and cleaning code are not. Possible improvements include: completing analysis data and/or code, adding raw data (+RD), and adding analysis code (+AC).
- **Level 4 (L4):** All analytic data sets and analysis code are available, but code does not run or produces results different than those in the paper (not CRA). Possible improvements include: debugging the analysis code (DAC) or obtaining raw data (+RD).
- **Level 5 (L5):** Analytic data sets and analysis code are available. They produce the same results as presented in the paper (CRA). The reproducibility package may be improved by obtaining the original raw data sets.

This is the highest level that most published research papers can attain currently. Computational reproducibility *from raw data* is required for papers that are reproducible at Level 6 and above.

- **Level 6 (L6):** Cleaning code is partially available, but raw data is not. Possible improvements include: completing cleaning code (+CC) and/or raw data (+RD).
- **Level 7 (L7):** Cleaning code is available and complete, but raw data is not. Possible improvements include: adding raw data (+RD).
- **Level 8 (L8):** Cleaning code is available and complete, and raw data is partially available. Possible improvements include: adding raw data (+RD).
- **Level 9 (L9):** All the materials (raw data, analytic data, cleaning code, and analysis code) are available. The analysis code produces the same output as presented in the paper (CRA). However, the cleaning code does not run or produces different results than those presented in the paper (not CRR). Possible improvements include: debugging the cleaning code (DCC).
- **Level 10 (L10):** All the materials are available and produce the same results as presented in the paper with minimal effort, starting from the analytic data (yes CRA) or the raw data (yes CRR). Note that Level 10 is aspirational and may be very difficult to attain for most research published today.

Table 3.5: Levels of Computational Reproducibility (P denotes "partial", C denotes "complete")

	Availability of materials, and reproducibility									
	Analysis Code		Analysis Data		CRA	Cleaning Code		Raw Data		
	P	C	P	C		P	C	P	C	
L1: No materials	–	–	–	–	–	–	–	–	–	–
L2: Only code			–	–	–	–	–	–	–	–
L3: Partial analysis data & code				–	–	–	–	–	–	–
L4: All analysis data & code					–	–	–	–	–	–
L5: Reproducible from analysis						–	–	–	–	–
L6: Some cleaning code							–	–	–	–
L7: All cleaning code								–	–	–
L8: Some raw data									–	–
L9: All raw data										–
L10: Reproducible from raw data										

^a **Computationally Reproducible from Analytic data (CRA):** The output can be reproduced with minimal effort

^b **Computationally Reproducible from Raw data (CRR):** The output can be reproduced with minimal effort

The following figure summarizes the different levels of computational reproducibility (for any given output). For each level, there will be improvements that have been made () or can be made to move up one level of reproducibility (-).

You may disagree with some of the levels outlined above, particularly wherever subjective judgment may be required. If so, you are welcome to interpret the levels as unordered categories (independent from their sequence) and suggest improvements using the "Edit" button above (top left corner if you are reading this document in your browser).

Adjusting Levels To Account for Confidential/Proprietary Data

A large portion of published research in economics uses confidential or proprietary data, most often government data from tax records or service provision and what is generally referred to as *administrative data*. Since administrative and proprietary data are rarely publicly accessible, some of the reproducibility levels presented above only apply once modified. The underlying theme of these modifications is that when data cannot be provided, you can assign a reproducibility score based on the level of detail in the instructions for accessing the data. Similarly, when reproducibility cannot be verified based on publicly available materials, the reproduction materials should demonstrate that a competent and unbiased third party (not involved in the original research team) has been able to reproduce the results.

- **Levels 1 and 2** can be applied as described above.

- **Adjusted Level 3 (L3*):** All analysis code is provided, but only partial instructions on how to access the *analysis data* are available. This means that the authors have provided some, but not all, of the following information:
 - a. *Contact information*, including name of the organization(s) that provides access to the data and contact information of at least one individual.
 - b. *Terms of use*, including licenses and eligibility criteria for accessing the data, if any.
 - c. *Information on data files (meta-data)*, including the name(s) and number of files, file size(s), relevant file version(s), and number of variables and observations in each file. Though not required, other relevant information may be included, including a description dataset dictionary, summary statistics, and synthetic data (fake data with the same statistical properties as the original data)
 - d. *Estimated costs for access*, including monetary costs such as fees and licences required to access the data, and non-monetary costs such as wait times and specific geographical locations from where researchers need to access the data.
- **Adjusted Level 4 (L4*):** All analysis code is provided, and complete and detailed instructions on how to access the *analysis data* are available.
- **Adjusted Level 5 (L5*):** All requirements for Level 4* are met, and the authors provide a certification that the output can be reproduced from the analysis data (CRA) by a third party. Examples include a signed letter by a disinterested reproducer or an official reproducibility certificate from a certification agency for data and code (e.g., see *cascad*).
- **Levels 6 and 7** can be applied as described above.
- **Adjusted Level 8 (L8*):** All requirements for Level 7* are met, but instructions for accessing the *raw data* are incomplete. Use the instructions described in Level 3 above to assess the instructions' completeness.
- **Adjusted Level 9 (L9*):** All requirements for Level 8* are met, and instructions for accessing the *raw data* are complete.
- **Adjusted Level 10 (L10*):** All requirements for Level 9* are met, and a certification that the output can be reproduced from the raw data is provided.

3.3.2 Reproducibility dimensions at the paper level

In addition to the output-specific assessment and improvement of computational reproducibility, several practices can facilitate reproducibility at the level of the overall paper. You can read about such practices in greater detail in the next chapter, dedicated to Stage 3: *Improvements*. In this Assessment section, you

Table 3.6: Levels of Computational Reproducibility with Proprietary/Confidential Data (P denotes "partial", C denotes "complete")

	Availability of materials, and reproducibility								
	Analysis Code		Instr. Analysis Data		CRA	Cleaning Code		Instr. Raw Data	
	P	C	P	C		P	C	P	C
L1: No materials	—	—	—	—	—	—	—	—	—
L2: Only code			—	—	—	—	—	—	—
L3: Partial analysis data & code				—	—	—	—	—	—
L4*: All analysis data & code					—	—	—	—	—
L5*: Proof of third party CRA						—	—	—	—
L6: Some cleaning code							—	—	—
L7: All cleaning code								—	—
L8*: Some instr. for raw data									—
L9*: All instr. for raw data									
L10*: Proof of third party CRR									

^a **Computationally Reproducible from Analytic data (CRA):** The output can be reproduced with minimal effort

^b **Computationally Reproducible from Raw data (CRR):** The output can be reproduced with minimal effort

should only verify whether the original reproduction package made use of any of the following:

- Master script that runs all steps
- Readme file
- Standardized file organization
- Version control
- Open source (statistical) software
- Dynamic document
- Computing capsule (e.g. CodeOcean, Binder, etc.)

Congratulations! You have now completed the *Assessment* stage of this exercise. You have provided a concrete building block of knowledge to improve understanding of the state of reproducibility in Economics.

Please continue to the next section where you can help improve it!

Chapter 4

Improvements

As you assess the paper’s reproducibility, you can start proposing ways to improve its reproducibility. These improvements can be specific to a display item or at the paper level. Improvements can be either implemented or in the form of specific suggestion for future reproducers to implement. Considering improvements is an opportunity to gain a deeper understanding of the paper’s methods, findings, and overall contribution. Each contribution can also be assessed and used by the wider Social Science Reproduction Platform (SSRP) community, including other students and researchers using the SSRP.

Part of the improvements might require you to engage with the original authors of the study you are reproducing. This stage will help you identify if the authors have already been contacted with a similar request, and if not, on how to approach them in order to have a constructive exchange.

As with the *Assessment* stage, we recommend that you first focus on one specific display item (e.g., “Table 1”). After making improvements to this first item, you will have a much easier time translating those improvements to other items.

4.1 Display item improvements

As part of your assessment of specific display items, you will identify potential issues with the original reproduction package (for any score lower than level 10). In addition to identifying these gaps, you are encouraged to implement specific improvements. In this section we suggest steps on how to add missing materials (data or code), or debug analysis or cleaning code. Record these improvements in the “Display item improvements” section.

4.1.1 Adding raw data: missing files or metadata

Reproduction packages often do not include all original raw datasets. To obtain any missing raw data or information about them, follow these steps:

1. Identify the missing file. During the Assessment stage, you identified all data sources from the paper’s body and appendices (column `data_source` in this standardized spreadsheet). However, some data sources (as collected by the original investigators) might be missing one or more files. You can sometimes find the specific name of those files by looking at the beginning of the cleaning code scripts. If you find the name of the file, record it in the `known_missing` field of the same spreadsheet as above. If not, record it as “Some/All” in the `known_missing` field of the for each specific data source.
2. Verify whether this file (or files) can be easily obtained from the web.
 - 2.1 - If yes: obtain the missing files and add them to the reproduction package. Make sure to obtain permission from the original author to publicly share this data. See tips for communication for more guidance.
 - 2.2 - If no: proceed to step 3.
3. Use the SSRP to verify whether there previous reproducers have contacted the authors regarding this paper and the specific missing files.
4. Contact the original authors and politely request the original materials. Be mindful of their time, and remember that the paper you are trying to reproduce was possibly published at a time when standards for computational reproducibility were different. See tips for communication for sample language on how to approach the authors for this specific scenario.
5. If the datasets are not available due to legal or ethical restrictions, you can still improve the reproduction package by providing detailed instructions for future researchers to follow, including contact information and possible costs of obtaining the raw data.

4.1.2 Adding missing analytic data files

Analytic data might be missing for two reasons: (1) raw data exists, but the procedures to transform it into analytic data are not fully reproducible, or (2) some or all raw data is missing, and some or all analytic data is not included in the original reproduction package. To obtain any missing analytic data, follow these steps:

1. Identify the specific name of the missing data set. Typically this information can be found in some of the analysis code that calls the data

to perform an analysis (e.g., `analysis_data_03.csv`).

2. Verify that the data cannot be obtained by running the data cleaning code over the raw data.
3. Use the SSRP to verify if previous attempts have been made to contact the authors about this data.
4. Contact the authors and request the specific data set.

4.1.3 Adding missing analysis code

Analysis code can be added when analytic data files are available, but some or all methodological steps are missing from the code. In this case, follow these steps:

1. Identify the specific line or paragraph in the paper that describes the analytic step that is missing from the code (e.g., “We impute missing values to...” or “We estimate this regression using a bandwidth of...”).
2. Identify the code file and the approximate line in the script where the analysis can be carried out. If you cannot find the relevant code file, identify its location relative to the main folder using the steps in the reproduction diagram.
3. Use the ACRE database to verify if previous attempts have been made to contact the authors about this issue.
4. Contact the authors and request the specific code files.
5. If step #4 does not work, we encourage you to attempt to recreate the analysis using your own interpretation of the paper, and making explicit your assumptions when filling in any gaps.

4.1.4 Adding missing data cleaning code

Data cleaning (processing) code might be added when steps are missing in the creation or re-coding of variables, merging, subsetting of the data sets, or other steps related to data cleaning and processing. You should follow the same steps you used when adding missing analysis code (1-5).

4.1.5 Debugging analysis code

Whenever code is available in the reproduction package, you should be able to debug those scripts. There are four types of debugging that can improve the reproduction package:

Table 4.1: Level-specific quality improvements: add data/code, debug code

output_name	imprv	description_of_added_files	lvl
table 1	+AD	ADD EXAMPLES	5
table 1	+RD	ADD EXAMPLES	5
table 1	DCC	ADD EXAMPLES	5
figure 1	+CC		6
figure 1	DAC		6
inline 1	DAC		8
...

- *Code cleaning*: Simplify the instructions (e.g., by wrapping repetitive steps in a function or a loop) or remove redundant code (i.e., old code that was commented out) while keeping the original output intact.
 - *Performance improvement*: Replace the original instructions with new ones that perform the same tasks but take less time (e.g., choose one numerical optimization algorithm over another while still obtaining the same results).
 - *Environment set up*: Modify the code to include correct paths to files, specific versions of software, and instructions to install missing packages or libraries.
 - *Correcting errors*: A coding error will occur when a section of the code in the reproduction package executes a procedure that is in direct contradiction with the intended procedure expressed in the documentation (i.e., paper or code comments). For example, an error will occur if the paper specifies that the analysis is performed on a population of males, but the code restricts the analysis to females only. Please follow the ACRE procedure to report coding errors.

4.1.6 Debugging cleaning code

Follow the same steps that you did to debug the analysis code, but report them separately.

4.1.7 Reporting results

Track all the different types of improvements you make and record in this standardized spreadsheet with the following structure:

4.2 Paper-level improvements

There are at least six additional improvements you can make to improve a paper’s overall reproducibility. These additional improvements can be applied across all reproducibility levels (including level 10). Record these improvements in the “Paper-level improvements” section.

1. Set up the reproduction package using version control software, such as Git.
2. Improve documentation by adding comments to the code.
3. Integrate the documentation with the code by adapting the paper into a literate programming environment (e.g., using Jupyter notebooks, RMarkdown, or a Stata Dynamic Doc).
4. If the code was written using a proprietary statistical software (e.g., Stata or Matlab), re-write some parts of it using an open source statistical software (e.g., R, Python, or Julia).
5. Re-organize the reproduction package into a set of folders and sub-folders that follow standardized best practices, and add a master script that executes all the code in order, with no further modifications. See AEA’s reproduction template.
6. Set up a computing capsule that executes the entire reproduction in a web browser without needing to install any software. For examples, see Binder and Code Ocean.

Chapter 5

Checking for Robustness

Once you have assessed, and potentially improved, the computational reproducibility of the display items for a claim within a paper, you can assess the robustness of these results by modifying some analytic choices and reporting their subsequent effects on the estimates of interest, i.e. conducting **robustness checks**. The universe of robustness checks can be very large (potentially infinite!) and they can pertain to both data analysis and data cleaning. In these guidelines, we will distinguish between **reasonable** and **feasible** robustness checks.

Reasonable robustness checks (Simonsohn et. al., 2018) are defined as (i) sensible tests of the research question, (ii) expected to be statistically valid, and (iii) not redundant with other specifications in the set. The set of **feasible robustness checks** is defined by all the specifications that can be computationally reproduced. We assume that the specifications already published in the paper are part of the reasonable set of specifications.

The size of the feasible set of robustness checks, and the likelihood that it contains reasonable specifications, will depend on the current level of reproducibility of the results supporting a claim. This is illustrated in Figure 5.1. At levels 1-2, it won't be possible to perform additional robustness checks because there is no data to work with. It may be possible to perform additional robustness checks for claims supported by display items reproducible at levels 3-4, but not using the specific estimates declared in the *Scoping Stage* since the display items are not computationally reproducible from analysis data (lacking CRA). It is possible to conduct additional robustness checks to validate the core conclusions of a claim based on a display item reproducible at level 5. Finally, claims associated with display items reproducible at level 6 or higher allow for robustness checks that involve variable definitions other types of analytical choices.

The size of feasible robustness checks grows exponentially as higher levels of computational reproducibility are achieved. For example, when checking the



Figure 5.1: Universe of robustness tests and its elements

robustness to a new variable definition, you will also be able to test the combination of how the main estimate changes under an alternative variable definition *and* an alternative variable definitions.

Robustness is assessed at the claim level (see our diagram representing a paper’s components 3). For a given claim, there will be several specifications presented, one of which will be identified by the authors (or yourself if the authors did not identify one) as the main or preferred specification. Identify which display item contains this specification and refer to the reproduction tree to identify the code files in which you can modify a computational choice. Using the example tree discussed in the *Assessment* stage, we can obtain the following (we removed the data files for simplicity). This simplified tree provides a list of potential files in which you can test different specifications:

```
table1.tex (contains preferred specification of a given claim)
  |___[code] analysis.R
    |___[code] final_merge.do
      |___[code] clean_merged_1_2.do
        |___[code] merge_1_2.do
          |___[code] clean_raw_1.py
          |___[code] clean_raw_2.py
        |___[code] clean_merged_3_4.do
          |___[code] merge_3_4.do
            |___[code] clean_raw_3.py
            |___[code] clean_raw_4.py
```

Here we suggest two types of contributions to robustness checks: (1) increasing the number of feasible robustness checks by identifying key analytical choices in code scripts, and (2) justifying and testing reasonable specifications within the set of feasible checks. Both contributions should be recorded on the SSRP Platform and refer to specific files in the reproduction package.

5.1 Feasible robustness checks: increasing the number of feasible specifications

Increasing the number of feasible robustness checks requires that you, as the reproducer identify the specific line(s) in the code scripts that execute an analytical choice. An advantage of this type of contribution is that you don’t need to have an in-depth knowledge of the paper and its methodology to contribute. This allows you to potentially map several code files, achieving a broader understanding of the paper, and also building on top of the work of others. The disadvantage is that you are not expected to test and justify the reasonableness of an alternative specifications.

Analytical choices can include those behind data cleaning and data analysis.

Below are some proposed types for each category.

Analytical choices in data cleaning code

- Variable definition
 - Data sub-setting
 - Data re-shaping (merge, append, long/gather, wide/spread)
 - Others (specify as “processing - other”)

Analytical choices in analysis code

- Regression function (link function)
 - Key parameters (tuning, tolerance parameters, etc.)
 - Controls
 - Adjustment of standard errors
 - Choice of weights
 - Treatment of missing values
 - Imputations
 - Other (specify as “methods - other”)

To record a specific analytical choice in the ACRE platform, please follow these steps:

1. Review a specific code file (e.g. `clean_merged_1_2.do`) and identify an analytical choice (e.g. `regress y x if gender == 1`).
2. Record the file name, line number, reproduction package (original or name of revised version), choice type, and choice value. For the `source` field, type “*original*” whenever the analytical choice is identified for the first time, and `line number` each time the same analytical choice is applied thereafter (for example, if an analytical choice is identified for the first time in line #103 and for the second time in line #122 their respective values for the `source` field should be `original` and 103). For each analytical choice recorded, add the specific choice used in the paper and, optionally, describe what alternatives could have been used. The resulting database should have the following structure:

entry	file_name	line_number	choice_type	choice_value	choice_range	Source
1	code_0173	10	data sub-setting	males	males, female,	original
2	code_01102	22	variable definition	income = wages + capital gains	wages, capital gains, gifts	“code_01.do-L103”
3	code_05113	13	controls	age, income, education	age, income, education, region	original
...

5.2 Reasonable robustness check: justifying and testing.

Justifying and testing a specific analytical choice requires that the reproducer identifies a feasible analytical choice, conducts a variation on it, and justifies its reasonableness. The advantage of this approach is that it allows for an in-depth inspection of a specific section of the paper. The main limitation is that justifying sensibility and validity (and non-redundancy, to an extent) requires a deeper understanding of the paper's topic and the methods, making it less feasible for undergraduate students or graduates with only a surface-level interest in the paper and limited time.

When performing a specific robustness check, follow these steps:

1. Search the database of feasible robustness checks (discussed above) and record the identifier(s) corresponding to the analytical choice to test (`entry_id`). If there is no entry corresponding for the specific lines, create one yourself.
2. Propose a specific variation to this analytical choice.
3. Discuss whether you think this variation is sensible, specifically in the context of the claim tested (e.g., does it make sense to include or exclude low-income Hispanic people from the sample when assessing the impact of a large wave of new immigrants?).
4. Discuss how this variation could affect the validity of the results (e.g., likely effects on omitted variable bias, measurement error, change in the Local Average Treatment Effects for the underlying population).
5. Confirm that test is not redundant with other tests in the paper or robustness exercise.
6. Report the results from the robustness check (new estimate, standard error, and units)[is what's in the brackets the suggested reporting format for this – if so, be explicit about it].

Chapter 6

Concluding the Reproduction

Once you have completed each of the reproduction stages for all of your claims of interest, you will be ready to submit your work. For a reproduction to be considered complete and ready for submission, you must have **assessed at least one display item**.

Before submitting a reproduction, you will be able to modify your answers to any entry in the ACRE platform. After you hit “Submit”, however, you will not be able to modify your reproduction attempt any further. If you wish to modify your reproduction after submitting it, you will have to record *a new reproduction attempt* on the platform and link to the previously completed reproduction.

6.1 Outputs

A completed reproduction will consist of three different types of outputs:

1. Revised reproduction package – Deposit your revised version of the original reproduction package in a trusted repository, such as Dataverse, openICPSR, Figshare, Dryad, Zenodo, or the Open Science Framework. You should submit a revised reproduction package any time that you perform any type of improvement to the original reproduction package. This revised reproduction package is expected to be self-contained as it might be used by future reproducers for assessment and improvement. If your new reproduction package is larger than 2Gb (?) [A: how is this determined, it sounds arbitrary?], or it contains data that you don’t have permission to share, remove the specific files from your reproduction package and add a reference to the original reproduction package [A: what does making a reference mean and look like in this context? I suggest providing an example of how to do that.].

Before submitting the reproduction, make sure that your revised reproduction package fulfills the following requirements:

- Has a digital object identifier (DOI). All of the trusted repositories referenced above generate DOIs. This is a stable, citable link that will allow others to easily find your work.
- Is titled following this convention: **revised reproduction package for - title of the paper - last name of reproducer - year when the reproduction was completed**
- If it is not self-contained, indicate the DOI of the last reproduction package used for the reproduction (the original or a previously revised reproduction package), and the relative file location of the missing files (e.g. `/data/raw/large_file.csv`).

2. Reproduction report(s) – At the end of each stage, the Social Science Reproduction Platform will give you the option to generate a report summarizing your work for that stage. Once you’ve submitted your final reproduction, the platform will also automatically generate a final report with its own DOI (different from the DOI used for the revised reproduction package). If you are conducting this reproduction as part of a supervised course or project, your instructor or advisor should have defined the structure of those reports (e.g., Stage 1 and 2 reports may be part of problem set 1, and Stage 3 and 4 reports may be a part of problem set 2). Each report will contain the following information:

- **Scoping Report:** basic information about the paper, descriptive statistics about the number of claims identified and the subset that was assessed, and a summary of the claims and associated specifications.
- **Assessment Report:** a summary of the display items assessed and their connections with the claims.
- **Improvement Report:** descriptions of improvements implemented and of improvements suggested for future reproductions, and an updated reproducibility score, if any.
- **Robustness Report:** summary of all of the analytical choices identified, and a description of the results of robustness tests to reasonable new specifications, if any.
- **Final reproduction report:** includes everything, plus final comments, and it cannot be edited after submission.

3. Your data – You will be able to export a .csv file containing all of your responses recorded as part of your reproduction attempt.

6.2 Anonymity and data sharing

- **Anonymity:** You may choose to post your reproduction anonymously for up to one year after submitting it. During this embargo period, unidentifiable data from your reproduction (i.e., numerical and categorical form responses) will be visible to other platform users, but your identity and identifiable data (i.e., text form responses) will not be revealed.
- **Using the data from your reproduction:** The Social Science Reproduction Platform will aggregate descriptive statistics from all submitted reproductions recorded to produce reproducibility metrics across disciplines, sub-

disciplines, journals, and topical bodies of literature.

Chapter 7

Guidance for a Constructive Exchange Between Reproducers and Original Authors

This chapter contains guidance for constructive and respectful communication between reproducers and original authors. Exchanges that contain charged or adversarial language can damage professional relationships and hamper scientific progress. Janz and Freese (2019) articulate two important steps reproducers can take to ensure that their interactions with original authors are constructive. We summarize and build on this approach below and encourage you to follow this guidance. Remember the **golden rule of reproductions** (and replications): *treat others and their work, as you would like others to treat you and your work!*

1. Carefully and transparently plan your study.

- a. Clearly state that you are conducting a reproduction of their original work.
- b. Explain why you have chosen this study.
- c. Explain how “far” your results must deviate from the original work before claiming that the study could not be reproduced. Engage deeply with the substantive literature to ensure that your interpretation of differences between the original and reproduction is thorough and acceptable to other authors in the field.

2. Use professional and sensitive language. Discuss potential discrepancies between your work and the original paper just like you might do for your own work.

- a. Avoid binary judgments and statements like “failed to reproduce.” Clearly

state which results reproduced and which did not (e.g., “we successfully reproduced X, but failed to reproduce Y”) unless you uncover apparent scientific misconduct (e.g., see Broockman, Kalla and Aronow, 2015).

- b. Talk about *the study*, not *the author*, to avoid making it personal. Make clear what the positive contribution of the original article is. Consider sending a copy of your reproduction report to the original authors.
- c. Discuss what your reproduction contributes to the literature, and refrain from claiming to give the final answer to the question.
- d. For papers published five or more years ago, be mindful that norms for reproducibility have evolved since then.
- e. Remember, *the goal is not to criticize previous work or hunt for errors, but to move the literature forward!*

To help you put these recommendations into practice, we’ve developed template language for common scenarios that reproducers and authors may encounter in their interactions.

While we hope that you find these useful, note that they are *only recommendations*, and you are welcome to modify them based on the context and needs of your specific project. Feel free to contact us if you need more guidance or would like to provide feedback on these materials.

7.1 For reproducers contacting the authors of the original study

Consider the following *before* you contact the original author:

1. Carefully read all footnotes, appendices, tables, captions, etc. to learn if, how, and where reproduction materials are provided. Follow this Data and Code Guidance to determine whether you have everything before you start. A few things to consider:
 - A *Readme* file, if available, would be a good place to start. For economics, all papers published in AEA journals after July 2019 should have such documentation.
 - Check whether there are any restrictions on accessing the data or code, and whether there are instructions on how to access these files for the purpose of reproduction.
2. If a reproduction package is not readily available in the location where the article is published (e.g., the journal website), check the authors’ websites,

7.1. FOR REPRODUCERS CONTACTING THE AUTHORS OF THE ORIGINAL STUDY⁵³

Dataverse profiles, or other relevant archives and/or data repositories like the ICPSR Publications Related Archive.

3. If steps 1 and 2 don't yield anything, contact the corresponding author (copying the co-authors, if any), consolidating your requests into as few emails as possible. In your email, make sure to include the following details:
 - Basic information about the paper being reproduced (include title, version, date, and a DOI (or just a URL));
 - Context for the reproduction (as part of a class exercise, thesis, personal project, etc.) and a note that the outcome will be recorded on the Social Science Reproduction Platform (SSRP);
 - Items from the reproduction package that are missing, as well as locations where you had (unsuccessfully) searched for them;
 - Your use plan: Will the materials be used exclusively for this project? Ask for permission to share the data publicly.
 - Right to consultation and results: Will you share the outcome of the reproduction with the original authors?
 - A deadline to respond (we suggest at least two weeks).
4. Follow up if you don't get a response within two weeks (or whatever deadline you set), and include any details or clarifications that were left out in your first email.
5. Record the outcome of your interaction with the original author on the SSRP. You can qualify the outcome as one of the following:
 - A *complete* reproduction package was provided
 - An *incomplete* reproduction package was provided. You can also select one of the following reasons:
 - Data is of sensitive, confidential, or proprietary character and cannot be shared;
 - Data is of sensitive, confidential, or proprietary character, but access instructions were provided.
 - The author *declined* to share the reproduction package
 - The author *did not respond* (including after a reminder was sent) within 4 weeks after the initial request.

7.1.1 Contacting the original author(s) when there is no reproduction package

Template email:

Subject: Reproduction package for ["Title of the paper"]

Dear [Title (e.g., "Dr.") Last name of Corresponding Author],

I am contacting you to request a reproduction package for your paper titled [Title] which was published in [Journal] in [year] (vol [volume], no. [no.]), [link]. A reproduction package contains (raw and/or analytic) data, code, and other documentation that makes it possible to reproduce the paper. Would you be able to share any of these items?

I am a [graduate student/postdoc/other position] at [Institution], and I would like to reproduce the results, tables, and other figures using the reproduction materials mentioned above. I have chosen this paper because [add context for why you want to reproduce this particular paper using neutral language (e.g., "This is a seminal paper in my field"), avoiding any statements that would put the respondent on the defensive]. Unfortunately, I was not able to locate any of these materials on the journal website, Dataverse [or other data and code repositories], or your website.

I will record the result of my reproduction attempt on the Social Science Reproduction Platform (SSRP), an open-source repository for the results of verifications of computational reproducibility of published work in the social sciences. SSRP is hosted by the Berkeley Initiative for Transparency in the Social Sciences). With your permission, I will also record the materials you share with me, which would allow access for other reproducers and avoid repeated requests directed to you. Please let me know if there are any legal or ethical restrictions that apply to any of the reproduction materials so that I can take that into consideration during this exercise.

In addition to your response above, would you be available to respond to future (non-repetitive) inquiries from me or other SSRP users? Though your cooperation with my and/or future requests would be extremely helpful, please note that you are *not required to respond*.

Since I am required to complete this project by [date], I would appreciate your response by [deadline].

Let me know if you have any questions. Please also feel free to contact my supervisor/instructor [Name (email)] for further details on this exercise. Thank you in advance for your help!

Best regards,
[Reproducer]

7.1.2 Contacting the original author(s) to request specific missing items of a reproduction package

Template email:

Subject: Reproduction materials for ["Title of the paper"]

Dear [Title (e.g., "Dr.") Last name of Corresponding Author],

I am contacting you regarding reproduction materials for your paper titled [Title] which was published in [Journal] in [year] (vol [volume], no. [no.]), [link].

I am a [graduate student/postdoc/other position] at [Institution], and I'm working to reproduce this paper as part of a class exercise. [Add context for why you want to reproduce this particular paper using neutral language (e.g., "This is a seminal paper in my field"), avoiding any statements that would put the respondent on the defensive].

To help me reproduce the paper in full, I hope that you can share the following items: [list items missing from reproduction package, preferably bulleted if more than one (e.g., raw/analytic data, code, protocols for conducting the experiment, etc.)]. I have already searched [locations where you searched for items, with links provided], but I was unable to locate the items. You can be assured that I will not share any of the materials without your permission, and I will use them exclusively for the purpose of this exercise. Let me know if there are any legal or ethical restrictions that apply to any of the reproduction materials so that I can take that into consideration during this exercise.

Note that I will record the outcome of my reproduction on the Social Science Reproduction Platform(SSRP), an open-source repository for the results of verifications of computational reproducibility of published work in the social sciences. SSRP is hosted by the Berkeley Initiative for Transparency in the Social Sciences). Let me know if you would like me to share the outcome of my reproduction with you, and whether you might be interested in providing a response.

Since I am required to complete this project by [date], I would appreciate your response by [deadline].

Let me know if you have any questions. Please also feel free to contact my supervisor/instructor [Name (email)] for further details on this exercise. Thank you in advance for your help!

Best regards,
[Reproducer]

7.1.3 Asking for additional guidance when some materials have been shared

Note: Even when a corresponding author has shared a reproduction package, you may still run into challenges in interpreting or executing the materials. That shouldn't discourage you from asking the corresponding author to provide clarifications or to share missing materials. As in the previous scenario described above, demonstrate that you've made an honest effort to reproduce the work using the available resources and try to consolidate your requests into as few emails as possible.

Template email:

Subject: Clarification for reproduction materials for ["Title of the paper"]

Dear [Title (e.g., "Dr.") Last name of Corresponding Author],

Thank you for sharing your materials. They have been immensely helpful for my work.

Unfortunately, I ran into a few issues as I delved into the reproduction, and I think your guidance would be helpful in resolving them. [Describe the issues and how you have tried to resolve them. Describe whatever files or parts of the data or code are missing. Refer to examples 1 and 2 below for more details].

Thank you in advance for your help.

Best regards,
[Reproducer]

1: An example of well-described issues:

Specifically, I am attempting to reproduce OUTPUT X (e.g., table 1, figure 3). I found that the following components are required to reproduce OUTPUT X:

OUTPUT X

```
[code] formatting_table1.R
      output1_part1.txt
|      [code] output_table1.do
|      [data] analysis_data01.csv
|      [code] data_cleaning01.R*
|      [data] UNKNOWN
```


7.1. FOR REPRODUCERS CONTACTING THE AUTHORS OF THE ORIGINAL STUDY⁵⁷

```
output1_part2.txt
[code] output_table2.do
[data] analysis_data02.csv
[code] data_cleaning02.R
[data] admin_01raw.csv*
```

I have marked with an asterisk (*) the items that I could not find in the reproduction materials: **data_cleaning01.R** and **admin_01raw.csv**. After accessing these files, I will also be able to identify the name of the raw data set required to obtain output1_part1.txt. This is to let you know that I may need to contact you again if I cannot find this file (labeled as UNKNOWN above) in the reproduction materials.

I understand that this request will require some work for you, but I want to assure you that I will add these missing files to the reproduction package for your paper on the Social Science Reproduction Platform. **Doing this will ensure that you will not be asked twice for the same missing file.**

2. An example of poorly described issues:

Your paper does not reproduce. I have tried for several hours now, and can't get the DO files to run. Could you please share all the missing reproduction materials? Data and code sharing are basic principles of open science, so I am confident that you will do the right thing.

7.1.4 Response when the original author has declined to share due to *undisclosed reasons*

Note: You can also use this template if a corresponding author has not submitted a response after two or more follow-up emails.

Template email:

Subject: Re: Reproduction materials for ["Title of the paper"]

Dear [Title (e.g., "Dr.") Last name of Corresponding Author],

Thank you for considering my request. I will try to reproduce the paper using the available materials and will record the missing items accordingly on the Social Science Reproduction Platform (SSRP). I will also post my assessment of the reproducibility of the paper in its current form based on the SSRP reproducibility scale.

Let me know if you have any questions.

Best regards,
[Reproducer]

7.1.5 Response when the original author has declined to share due to legal or ethical restrictions of the data

Template email:

Subject: Re: Reproduction materials for ["Title of the paper"]

Dear [Title (e.g., "Dr.") Last name of Corresponding Author],

Thank you for your response and for clarifying the terms of use for the reproduction materials.

Though I understand you are unable to share the raw data, there may be alternative steps you can take that would help me improve the reproducibility of your paper. These include:

1. Sharing the analytic version of the data (the version of the dataset that was used for analysis in the final version of your paper);
2. Providing a public description of the steps other researchers can follow to request access to the raw data or materials, including an estimate of the costs and the duration of the process. You can find examples of data availability statements for proprietary or restricted-access data [here](#); and
3. Providing access to all data and materials for which the constraints do not apply.

Based on my assessment, your paper would currently rank at [level X] on the SSRP reproducibility scale. However, *this score can be easily improved*. Being able to provide analytic data would elevate the reproducibility of your paper to [level Y]. Providing public instructions on how other parties can access the data would further elevate its reproducibility to [level Z].

I would be happy to help if you are interested in taking any of the steps I outlined above.

Thank you for your help!

Best regards,
[Reproducer]

7.1.6 Contacting the original author to share the results of your reproduction exercise

Note: Reporting the results of reproductions can be the most contentious part of the process, particularly in instances where the reproducer is not able to fully reproduce the paper or finds significant deviations from the original work. However, if the reproduction can correctly identify the sources of such deviations, it may be viewed as an improved version of the original work.

Regardless of the outcome of the reproduction exercise, the guidance from the introduction of this chapter still stands: *reproduce the work of others as you would like for others to reproduce yours*, and make sure that is reflected in how you discuss any discrepancies between your and the original work.

Template email:

Subject: Reproducibility Assessment of ["Title of the paper"]

Dear [Title (e.g., "Dr.") Last name of Corresponding Author],

Thank you for your support throughout my project as I worked to verify and advance the reproducibility of [Paper]. I'm writing now to share the results of my project and to invite your feedback.

The results of each step of my reproduction include i) Assessment, ii) Improvements, iii) Robustness Checks, (and iv) Extensions, if applicable).

[Include the following items in the body of your email:

- Briefly describe which parts of the paper you tried to reproduce (e.g., a specific estimate, a table, etc.).
- Within the scope of your reproduction, describe exactly which items you were able to reproduce.
- Discuss the differences you observed between the results of your reproduction and the original work, and demonstrate that you did your due diligence in trying to reproduce each item. Remember that it is more constructive to discuss discrepancies, differences or deviations, rather than errors, mistakes, or failures, and *always talk about the work – not the author!*
- Use sensitive language when presenting discrepancies, e.g., "Unfortunately, I found X, which differs from the Y result in the original paper...". Be cognizant of any potential limitations of your work, and explain how you have tried to address them – that way you will proactively address potential criticism!

- Describe how you tried to improve the reproducibility of the paper. If some of the improvements are based on discretionary judgment (e.g., file organization or code commenting), try to explain why you think they are an improvement over the original work. If you didn't make improvements, point out some concrete steps that the author(s) can take to improve the reproducibility of the section you reproduced.]'

I look forward to your questions, comments, and suggestions for my work. As discussed previously, I will record the outcomes of my reproduction, along with the improvements, on the Social Science Reproduction Platform.

Best regards,
[Reproducer]

7.1.7 Responding to hostile responses from original authors

Note: Planning your study carefully and transparently, and using professional and sensitive language are the best ways to ensure that the interaction will be beneficial to both you and the original author. However, unpleasant interactions may happen despite your best efforts, and can range anywhere from dismissive comments to bullying, discrimination, and harassment. Find guidance at the end of this chapter on how to deal with instances of bullying, harassment, or discrimination.

7.1.7.1 Dismissive comments

In cases of dismissive comments, the best course of action may be to simply thank the author for their response and continue with the exercise.

Template email:

Subject: Re: Reproduction materials for ["Title of the paper"]

Dear [Title (e.g., "Dr.") Last name of Corresponding Author],

Thank you for your response. I will work to reproduce your paper using the available materials and will record my results accordingly on the Social Science Reproduction Platform. I will also post my assessment of the reproducibility of the paper in its current form based on the SSRP reproducibility scale.

Let me know if you have any questions.

Best regards,
[Reproducer]

7.2 For original authors Responding to requests from reproducers

This section contains guidance for authors of papers involved in reproductions on the Social Science Reproduction Platform. We present language that may be helpful for various scenarios in which authors find themselves when interacting with reproducers. Though every interaction between authors and reproducers takes place in a distinct context and may carry its own unique challenges, the guiding principle of this chapter always applies: “Treat others and their work as you would like others to treat you and your work!” We hope that these resources will facilitate more efficient and constructive exchanges between the parties involved. Let us know if you need guidance in other scenarios!

7.2.1 Responding to a repeated request that has been addressed in an earlier interaction

Dear [Reproducer],

Thank you for your interest in my work. Note that I have been contacted about this issue by another SSRP reproducer before and provided a response, which I suspect may be already recorded on the SSRP. I’m copying my original response below for your reference. You may find further guidance in the readme file in the reproduction package.

If there are no prior records of these issues on the SSRP, please record the enclosed response [and materials]. This will also help avoid the duplication of effort on the part of others who may be interested in reproducing this work. Good luck with the remainder of your project and thank you in advance for your cooperation!

Best regards, [Author]

7.2.2 Acknowledging that the author no longer has access to certain part(s) of the reproduction package

Dear [Reproducer],

Thank you for reviewing my work closely. I wish I could be of more help, but unfortunately I no longer have access to the requested materials due to [briefly describe the circumstances that prevent you from providing the materials].

While I recognize that the current standards in the discipline suggest that work should be reproducible using materials that are readily available for this purpose, note that this paper was written at a time when different standards of reproducibility were mandated.

Please feel free to evaluate the paper as is and propose any improvements wherever possible.

I look forward to working with you to address this and improve the overall reproducibility of the paper.

Best regards, [Author]

7.2.3 Acknowledging that some material is still embargoed for future research

Dear [Reproducer],

Thank you for your interest in my work. The data/materials/program that you reference are not publicly accessible for the time being because they are embargoed until [embargo period].

[Depending on the restrictions that apply to the reproduction package, consider alternatives to sharing the reproduction materials in full. These include: 1. Sharing the analytic version of the data (the version of the dataset that was used for analysis in the final version of your paper); 2. Providing a public description of the steps other researchers can follow to request access to the raw data or materials, including an estimate of the costs and the duration of the process. Find examples of data availability statements for proprietary or restricted-access data [here] (<https://social-science-data-editors.github.io/guidance>) and 3. Providing access to all data and materials for which the constraints do not apply.]

I hope you find this helpful for reproducing the paper. Please feel free to contact me if you have any further questions.

Best regards, [Author]

7.2.4 Responding to incomplete/unclear requests

Dear [Reproducer],

Thank you for your interest in my work. I would be happy to assist you and other reproducers to assess and improve the reproducibility of this paper.

To help me give more concrete guidance on this issue, I'd appreciate if you could provide a more specific description of the items that you need from me. You can find helpful information and resources in Chapter 6 of the Guide, specifically here [based on the context, you may need to point the reproducer to a different scenario and/or provide further information].

Feel free to contact me if you have any further questions. Thank you for your cooperation.

Best regards, [Author]

7.3 Harassment and/or discrimination

The American Economic Association and other economic societies have strict policies against harassment and discrimination. Here are some of the behaviors that the AEA Policy on Harassment and Discrimination has listed as unacceptable and could emerge in a hostile exchange regarding a reproduction:

- Intentionally intimidating, threatening, harassing, or abusive actions or remarks (both spoken and in other media)
 - Prejudicial actions or comments that undermine the principles of equal opportunity, fair treatment, or free academic exchange
 - Deliberate intimidation, stalking, or following
 - Real or implied threat of physical harm.

Here are a some steps you can take if you believe you have experienced bullying, discrimination or harassment:

- **File a complaint with the AEA Ombudsperson.** Any AEA member can file a complaint. (You can also join the AEA solely for the purpose of filing a report.) The person about whom you are making the complaint need not be an AEA member. A non-AEA member can also file a report if the act of harassment or discrimination was committed by an AEA member or in the context of an AEA-sponsored activity. Learn more about the process [here](#).
 - **File a report with your institution’s office for the prevention of harassment & discrimination.** US-based institutions have internal mechanisms that allow students and faculty to seek support in cases of discrimination and harassment on the basis of race, color, national origin, gender, age, or sexual orientation/identity, including allegations of sexual harassment and sexual violence. Formal titles of this office vary across institutions, but common names include “Office for the Prevention of Harassment and Discrimination” (in institutions that are part of the University of California system), “Office of Equity and Title IX”, etc.
 - **Contact your institution’s Ombudsperson/Ombuds Office.** If you believe that you have experienced academic bullying or other forms of disrespectful behavior that fall outside the scope of harassment and/or discrimination as described above, you should know that university ombuds officers are a confidential, impartial resource to discuss your concerns and learn about potential next steps available in your case.
 - **Access mental health services at your institution.** While no

amount of bullying, discrimination, or harassment is acceptable or the fault of the victim, these unfortunately still occur and can take a toll on victims' mental health. Many universities offer short-term Counseling & Psychological Services (CAPS) for academic, career, and personal issues.

- **Ask for support from your academic supervisor.** If you are unsure on how to proceed, consult your academic supervisor on whether continuing the reproduction is appropriate.

Chapter 8

Examples of Reproduction Trees

A diagram generated by the ACRE Diagram Builder which represents all the available data and code on behind a specific display item. The tree is meant to represent the entire computational workflow behind a result from the paper. It allows reproducers to trace a display item to its primary sources. It can also be used to guide users of the reproduction package and/or to identify missing components for a complete reproduction.

A reproduction tree is complete when it is possible to connect its output (a given display item) with all of its inputs down to the raw data. A reproduction is incomplete when it is not possible to connect all the inputs to the resulting display item. Paraphrasing the author Leo Tolstoy, *complete workflows are all alike; every incomplete computational workflow is incomplete in its own way.* This chapter presents a few examples of reproduction trees, focusing particularly on the many possible ways in which a tree could be incomplete. If you have a reproduction tree that contains an instructive example please contribute to this chapter (via a pull request or emailing your reproduction tree to ACRE@berkeley.edu.)

8.1 Stylized examples

8.1.1 Complete reproduction tree

Below is an example of output from the ACRE Diagram builder for a display item that can be fully constructed using the files contained in the reproduction package. The diagram displays files as outputs and inputs to code scripts all the way down to the raw data.

table 1

```

[code] formatting_table1.R
      output1_part1.txt
      |   [code] output_table1.do
      |       [data] analysis_data01.csv
      |           [code] data_cleaning01.R
      |               [data] survey_01raw.csv
      output1_part2.txt
      |   [code] output_table2.do
      |       [data] analysis_data02.csv
      |           [code] data_cleaning02.R
      |               [data] admin_01raw.csv

```

8.1.2 Incomplete reproduction tree

8.1.2.1 Raw data and analytic data are available, but cleaning code is missing.

Below is an example of output from the ACRE Diagram Builder for a display item that is missing some of the code needed to generate it from the raw data. There are two reasons to suspect that this workflow is incomplete: (i) there is no clear data cleaning step (only analysis that generates output and formatting), and (ii) there are unused files that are likely to be raw data. None of these reasons can confirm unequivocally that the tree is incomplete, but a reproducer familiar with the paper and its data sources could use the tree to certify its (in)completeness and request missing files.

```

table 1
[code] formatting_table1.R
      output1_part1.txt
      |   [code] output_table1.do
      |       [data] analysis_data01.csv
      output1_part2.txt
      |   [code] output_table2.do
      |       [data] analysis_data02.csv

```

Unused files:

- survey_01raw.csv
- admin_01raw.csv

Reproducers are asked to speculate on where the missing files might go, and hence propose how a complete tree might look like (where possible). For this example, we have assumed there are missing code scripts that at some point take in `survey_01raw.csv` and `admin_01raw.csv`, and eventually output `analysis_data01.csv` and `analysis_data02.csv`, though this requires the reproducer's discretion.

```

table 1
[code] formatting_table1.R

```

```

output1_part1.txt
| [code] output_table1.do
| [data] analysis_data01.csv
| [code] MISSING FILE(S)
| [data] survey_01raw.csv
output1_part2.txt
| [code] output_table2.do
| [data] analysis_data02.csv
| [code] MISSING FILE(S)
| [data] admin_01raw.csv

```

8.1.3 Unused Data Sources

It is possible that not all data included in a replication package are actually used in code scripts in the reproduction package. This would be the case if, for example, the raw data and analysis data are included, but not the script that generates the analysis data. As a concrete example, consider what the original diagram above would look like if the only code included in the reproduction package were analysis.R:

```

table1.tex
|___[code] analysis.R
|___analysis_data.dta

```

Unused data sources:

```

raw_1.dta
raw_2.dta
raw_3.dta
raw_4.dta

```

Unused analysis data:

```

cleaned_1.dta
cleaned_2.dta
cleaned_3.dta
cleaned_4.dta
merged_1_2.dta
merged_3_4.dta
cleaned_1_2.dta
cleaned_3_4.dta

```

In this case, there are many data files that were listed in the raw data and analytic data spreadsheets that are not used by any code script in the replication package.

8.2 Examples from real reproduction attempts

8.2.1 Possibly missing code for producing a display item¹.

This reproduction diagram fragment likely shows a missing piece of code. In a complete reproduction package there are no unused files and all final outputs are display items. `cps2018.dta` and `cps_march2017.dta` are included as inputs in the reproduction kit, but never used to make a display item, i.e., there is no piece of code that is listed as using them as inputs. Likewise, `PublicSalary.dta` is listed as a final output, meaning it, too, is not used to make a display item since it is not listed as an input for any code script. Perhaps `Paper2_PoExitDataset.dta` is made by a missing code script that takes `PublicSalary.dta`, `cps2018.dta`, and `cps_march2017.dta` as inputs? It may be the case that the only way to find out for sure is to contact the study author(s), in which case, such a diagram can be used to help them identify any missing files.

```
PublicSalary.dta
|___MakeData2.do
|   |___Paper2_ProviderDataset.dta
|   |___Paper2_SSPDataset_wIRT_final.dta
|   |___PublicFacilitySurvey_clean.dta

Table1.xml
|___Table1.do
|   |___ProviderData.dta
|   |   |___MakeData4.do
|   |   |   |___Paper2_PoExitDataset.dta
|   |   |   |___Paper2_ProviderDataset.dta
|   |   |   |___Paper2_SSPDataset_wIRT_final.dta
|   |___VillageDataset.dta
|   |   |___MakeData8.do
|   |   |   |___Paper2_HouseholdDataset1.dta
|   |   |   |___Paper2_VillageDataset.dta
|   |___HouseholdDataset.dta
|   |   |___MakeData8.do
|   |   |   |___Paper2_HouseholdDataset1.dta
|   |   |   |___Paper2_VillageDataset.dta

Unusued data sources:
cps2018.dta
cps_march2017.dta
```

¹This is from a reproduction attempt conducted as part of a UC Berkeley Development Economics course

8.2.2 Long, complicated tree².

This diagram shows that the production of a given display item may be very complicated, highlighting the usefulness of the ACRE Diagram Builder as a visualization tool. In reproducing such a complicated display item, it can be useful to have such a diagram to determine, for example, the order in which code scripts should be run, what files might depend on a faulty code script, or which files are necessary to keep if the goal is to only produce the specific display item.

```

results_days.dta
|___lightspaper_replication.do
|   |___isocvout.dbf
|   |   |___v4cvggini_bycountry_prep.aml
|   |   |   |___ctryliggrid
|   |   |   |   |___v4unzip.bat
|   |   |   |   |   |___rad_cal.tar
|   |   |   |   |   |___phys_geo.zip
|   |   |   |   |   |___world_avg_fixed.tfw
|   |   |   |   |   |___af_grumpv1_ppoints_csv.zip
|   |   |   |   |   |___gl_grumpv1_area_ascii_30.zip
|   |   |   |   |   |___malaria_ecology.zip
|   |   |   |   |   |___gl_gpww3_pcount_00_wrk_25.zip
|   |   |   |   |   |___gasflares.zip
|   |   |   |   |   |___F%sy%.v4.tar
|   |   |   |   |   |___boundaries.zip
|   |   |   |___l%sy%nm
|   |   |   |___v4lightsprep.aml
|   |   |   |   |___F%sy%.v4b_web.stable_lights.avg_vis.tfw
|   |   |   |   |   |___v4unzip.bat
|   |   |   |   |   |   |___rad_cal.tar
|   |   |   |   |   |   |___phys_geo.zip
|   |   |   |   |   |   |___world_avg_fixed.tfw
|   |   |   |   |   |   |___af_grumpv1_ppoints_csv.zip
|   |   |   |   |   |   |___gl_grumpv1_area_ascii_30.zip
|   |   |   |   |   |   |___malaria_ecology.zip
|   |   |   |   |   |   |___gl_gpww3_pcount_00_wrk_25.zip
|   |   |   |   |   |   |___gasflares.zip
|   |   |   |   |   |   |___F%sy%.v4.tar
|   |   |   |   |   |   |___boundaries.zip
|   |   |   |   |___F%sy%.v4b_web.stable_lights.avg_vis.tif
|   |   |   |   |   |___v4unzip.bat
|   |   |   |   |   |   |___rad_cal.tar
|   |   |   |   |   |   |___phys_geo.zip

```

²This is from a reproduction attempt conducted as part of a UC Berkeley Development Economics course

```

| | | | | world_avg_fixed.tfw
| | | | | ___af_grumpv1_ppoints_csv.zip
| | | | | ___gl_grumpv1_area_ascii_30.zip
| | | | | ___malaria_ecology.zip
| | | | | ___gl_gpww3_pcount_00_wrk_25.zip
| | | | | ___gasflares.zip
| | | | | ___F%sy%.v4.tar
| | | | | ___boundaries.zip
| | | | | ___world_avg.tfw
| | | | | |___v4unzip.bat
| | | | | |___rad_cal.tar
| | | | | |___phys_geo.zip
| | | | | |___world_avg_fixed.tfw
| | | | | |___af_grumpv1_ppoints_csv.zip
| | | | | |___gl_grumpv1_area_ascii_30.zip
| | | | | |___malaria_ecology.zip
| | | | | |___gl_gpww3_pcount_00_wrk_25.zip
| | | | | |___gasflares.zip
| | | | | |___F%sy%.v4.tar
| | | | | |___boundaries.zip
| | | | | ___gluareag_alpha1.asc
| | | | | |___v4unzip.bat
| | | | | |___rad_cal.tar
| | | | | |___phys_geo.zip
| | | | | |___world_avg_fixed.tfw
| | | | | |___af_grumpv1_ppoints_csv.zip
| | | | | |___gl_grumpv1_area_ascii_30.zip
| | | | | |___malaria_ecology.zip
| | | | | |___gl_gpww3_pcount_00_wrk_25.zip
| | | | | |___gasflares.zip
| | | | | |___F%sy%.v4.tar
| | | | | |___boundaries.zip
| | | | | ___F%sy%.v4b_web.cf_cvg.tif
| | | | | |___v4unzip.bat
| | | | | |___rad_cal.tar
| | | | | |___phys_geo.zip
| | | | | |___world_avg_fixed.tfw
| | | | | |___af_grumpv1_ppoints_csv.zip
| | | | | |___gl_grumpv1_area_ascii_30.zip
| | | | | |___malaria_ecology.zip
| | | | | |___gl_gpww3_pcount_00_wrk_25.zip
| | | | | |___gasflares.zip
| | | | | |___F%sy%.v4.tar
| | | | | |___boundaries.zip
| | | | | ___F%sy%.v4b_web.cf_cvg.tfw
| | | | | |___v4unzip.bat

```

```
| | | | |__rad_cal.tar  
| | | | |__phys_geo.zip  
| | | | |__world_avg_fixed.tfw  
| | | | |__af_grumpv1_ppoints_csv.zip  
| | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | |__malaria_ecology.zip  
| | | | |__gl_gpwv3_pcount_00_wrk_25.zip  
| | | | |__gasflares.zip  
| | | | |__F%sy%.v4.tar  
| | | | |__boundaries.zip  
| | | | |__world_avg.tif  
| | | | | |__v4unzip.bat  
| | | | | |__rad_cal.tar  
| | | | | |__phys_geo.zip  
| | | | | |__world_avg_fixed.tfw  
| | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | |__malaria_ecology.zip  
| | | | | |__gl_gpwv3_pcount_00_wrk_25.zip  
| | | | | |__gasflares.zip  
| | | | | |__F%sy%.v4.tar  
| | | | | |__boundaries.zip  
| | | | |__cvgsy%  
| | | | | |__v4lightsprep.aml  
| | | | | |__F%sy%.v4b_web.stable_lights.avg_vis.tfw  
| | | | | | |__v4unzip.bat  
| | | | | | |__rad_cal.tar  
| | | | | | |__phys_geo.zip  
| | | | | | |__world_avg_fixed.tfw  
| | | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | | |__malaria_ecology.zip  
| | | | | | |__gl_gpwv3_pcount_00_wrk_25.zip  
| | | | | | |__gasflares.zip  
| | | | | | |__F%sy%.v4.tar  
| | | | | | |__boundaries.zip  
| | | | | |__F%sy%.v4b_web.stable_lights.avg_vis.tif  
| | | | | | |__v4unzip.bat  
| | | | | | |__rad_cal.tar  
| | | | | | |__phys_geo.zip  
| | | | | | |__world_avg_fixed.tfw  
| | | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | | |__malaria_ecology.zip  
| | | | | | |__gl_gpwv3_pcount_00_wrk_25.zip  
| | | | | | |__gasflares.zip
```



```

|___ginioutu.dbf
|   |___v4cvggini_bycountry_prep.aml
|       |___ctryliggrid
|           |___v4unzip.bat
|               |___rad_cal.tar
|               |___phys_geo.zip
|               |___world_avg_fixed.tfw
|               |___af_grumpv1_ppoints_csv.zip
|               |___gl_grumpv1_area_ascii_30.zip
|               |___malaria_ecology.zip
|               |___gl_gpwv3_pcount_00_wrk_25.zip
|               |___gasflares.zip
|               |___F%sy%.v4.tar
|               |___boundaries.zip
|   |___l%sy%nm
|       |___v4lightsprep.aml
|           |___F%sy%.v4b_web.stable_lights.avg_vis.tfw
|               |___v4unzip.bat
|                   |___rad_cal.tar
|                   |___phys_geo.zip
|                   |___world_avg_fixed.tfw
|                   |___af_grumpv1_ppoints_csv.zip
|                   |___gl_grumpv1_area_ascii_30.zip
|                   |___malaria_ecology.zip
|                   |___gl_gpwv3_pcount_00_wrk_25.zip
|                   |___gasflares.zip
|                   |___F%sy%.v4.tar
|                   |___boundaries.zip
|           |___F%sy%.v4b_web.stable_lights.avg_vis.tif
|               |___v4unzip.bat
|                   |___rad_cal.tar
|                   |___phys_geo.zip
|                   |___world_avg_fixed.tfw
|                   |___af_grumpv1_ppoints_csv.zip
|                   |___gl_grumpv1_area_ascii_30.zip
|                   |___malaria_ecology.zip
|                   |___gl_gpwv3_pcount_00_wrk_25.zip
|                   |___gasflares.zip
|                   |___F%sy%.v4.tar
|                   |___boundaries.zip
|           |___world_avg.tfw
|               |___v4unzip.bat
|                   |___rad_cal.tar
|                   |___phys_geo.zip
|                   |___world_avg_fixed.tfw
|                   |___af_grumpv1_ppoints_csv.zip

```



```
|
|       |___world_avg_fixed.tfw
|       |___af_grumpv1_ppoints_csv.zip
|       |___gl_grumpv1_area_ascii_30.zip
|       |___malaria_ecology.zip
|       |___gl_gpww3_pcount_00_wrk_25.zip
|       |___gasflares.zip
|       |___F%sy%.v4.tar
|       |___boundaries.zip
|___cvg%sy%
|   |___v4lightsprep.aml
|   |___F%sy%.v4b_web.stable_lights.avg_vis.tfw
|   |   |___v4unzip.bat
|   |   |___rad_cal.tar
|   |   |___phys_geo.zip
|   |   |___world_avg_fixed.tfw
|   |   |___af_grumpv1_ppoints_csv.zip
|   |   |___gl_grumpv1_area_ascii_30.zip
|   |   |___malaria_ecology.zip
|   |   |___gl_gpww3_pcount_00_wrk_25.zip
|   |   |___gasflares.zip
|   |   |___F%sy%.v4.tar
|   |   |___boundaries.zip
|   |___F%sy%.v4b_web.stable_lights.avg_vis.tif
|   |   |___v4unzip.bat
|   |   |___rad_cal.tar
|   |   |___phys_geo.zip
|   |   |___world_avg_fixed.tfw
|   |   |___af_grumpv1_ppoints_csv.zip
|   |   |___gl_grumpv1_area_ascii_30.zip
|   |   |___malaria_ecology.zip
|   |   |___gl_gpww3_pcount_00_wrk_25.zip
|   |   |___gasflares.zip
|   |   |___F%sy%.v4.tar
|   |   |___boundaries.zip
|   |___world_avg.tfw
|   |   |___v4unzip.bat
|   |   |___rad_cal.tar
|   |   |___phys_geo.zip
|   |   |___world_avg_fixed.tfw
|   |   |___af_grumpv1_ppoints_csv.zip
|   |   |___gl_grumpv1_area_ascii_30.zip
|   |   |___malaria_ecology.zip
|   |   |___gl_gpww3_pcount_00_wrk_25.zip
|   |   |___gasflares.zip
|   |   |___F%sy%.v4.tar
|   |   |___boundaries.zip
```

```
| | | ---glureag_alpha1.asc  
| | |   |--v4unzip.bat  
| | |     |--rad_cal.tar  
| | |       |--phys_geo.zip  
| | |         |--world_avg.fixed.tfw  
| | |           |--af_grumpv1_ppoints_csv.zip  
| | |             |--gl_grumpv1_area_ascii_30.zip  
| | |               |--malaria_ecology.zip  
| | |                 |--gl_gpww3_pcount_00_wrk_25.zip  
| | |                   |--gasflares.zip  
| | |                     --F%sy%.v4.tar  
| | |                       |--boundaries.zip  
| | | ---F%sy%.v4b_web.cf_cvg.tif  
| | |   |--v4unzip.bat  
| | |     |--rad_cal.tar  
| | |       |--phys_geo.zip  
| | |         |--world_avg.fixed.tfw  
| | |           |--af_grumpv1_ppoints_csv.zip  
| | |             |--gl_grumpv1_area_ascii_30.zip  
| | |               |--malaria_ecology.zip  
| | |                 |--gl_gpww3_pcount_00_wrk_25.zip  
| | |                   |--gasflares.zip  
| | |                     --F%sy%.v4.tar  
| | |                       |--boundaries.zip  
| | | ---F%sy%.v4b_web.cf_cvg.tif  
| | |   |--v4unzip.bat  
| | |     |--rad_cal.tar  
| | |       |--phys_geo.zip  
| | |         |--world_avg.fixed.tfw  
| | |           |--af_grumpv1_ppoints_csv.zip  
| | |             |--gl_grumpv1_area_ascii_30.zip  
| | |               |--malaria_ecology.zip  
| | |                 |--gl_gpww3_pcount_00_wrk_25.zip  
| | |                   |--gasflares.zip  
| | |                     --F%sy%.v4.tar  
| | |                       |--boundaries.zip  
| | | ---world_avg.tif  
| | |   |--v4unzip.bat  
| | |     |--rad_cal.tar  
| | |       |--phys_geo.zip  
| | |         |--world_avg.fixed.tfw  
| | |           |--af_grumpv1_ppoints_csv.zip  
| | |             |--gl_grumpv1_area_ascii_30.zip  
| | |               |--malaria_ecology.zip  
| | |                 |--gl_gpww3_pcount_00_wrk_25.zip  
| | |                   |--gasflares.zip
```



```

|_F%sy%.v4.tar
|_boundaries.zip
|_l%sy%nm
|_v4lightsprep.aml
|_F%sy%.v4b_web.stable_lights.avg_vis.tfw
|_v4unzip.bat
|_rad_cal.tar
|_phys_geo.zip
|_world_avg_fixed.tfw
|_af_grumpv1_ppoints_csv.zip
|_gl_grumpv1_area_ascii_30.zip
|_malaria_ecology.zip
|_gl_gpww3_pcount_00_wrk_25.zip
|_gasflares.zip
|_F%sy%.v4.tar
|_boundaries.zip
|_F%sy%.v4b_web.stable_lights.avg_vis.tif
|_v4unzip.bat
|_rad_cal.tar
|_phys_geo.zip
|_world_avg_fixed.tfw
|_af_grumpv1_ppoints_csv.zip
|_gl_grumpv1_area_ascii_30.zip
|_malaria_ecology.zip
|_gl_gpww3_pcount_00_wrk_25.zip
|_gasflares.zip
|_F%sy%.v4.tar
|_boundaries.zip
|_world_avg.tfw
|_v4unzip.bat
|_rad_cal.tar
|_phys_geo.zip
|_world_avg_fixed.tfw
|_af_grumpv1_ppoints_csv.zip
|_gl_grumpv1_area_ascii_30.zip
|_malaria_ecology.zip
|_gl_gpww3_pcount_00_wrk_25.zip
|_gasflares.zip
|_F%sy%.v4.tar
|_boundaries.zip
|_gluareag_alpha1.asc
|_v4unzip.bat
|_rad_cal.tar
|_phys_geo.zip
|_world_avg_fixed.tfw
|_af_grumpv1_ppoints_csv.zip

```

```

|___gl_grumpv1_area_ascii_30.zip
|___malaria_ecology.zip
|___gl_gpww3_pcount_00_wrk_25.zip
|___gasflares.zip
|___F%sy%.v4.tar
|___boundaries.zip
|___F%sy%.v4b_web.cf_cvg.tif
|___v4unzip.bat
|___rad_cal.tar
|___phys_geo.zip
|___world_avg_fixed.tfw
|___af_grumpv1_ppoints_csv.zip
|___gl_grumpv1_area_ascii_30.zip
|___malaria_ecology.zip
|___gl_gpww3_pcount_00_wrk_25.zip
|___gasflares.zip
|___F%sy%.v4.tar
|___boundaries.zip
|___F%sy%.v4b_web.cf_cvg.tfw
|___v4unzip.bat
|___rad_cal.tar
|___phys_geo.zip
|___world_avg_fixed.tfw
|___af_grumpv1_ppoints_csv.zip
|___gl_grumpv1_area_ascii_30.zip
|___malaria_ecology.zip
|___gl_gpww3_pcount_00_wrk_25.zip
|___gasflares.zip
|___F%sy%.v4.tar
|___boundaries.zip
|___world_avg.tif
|___v4unzip.bat
|___rad_cal.tar
|___phys_geo.zip
|___world_avg_fixed.tfw
|___af_grumpv1_ppoints_csv.zip
|___gl_grumpv1_area_ascii_30.zip
|___malaria_ecology.zip
|___gl_gpww3_pcount_00_wrk_25.zip
|___gasflares.zip
|___F%sy%.v4.tar
|___boundaries.zip
|___cvg%sy%
|___v4lightsprep.aml
|___F%sy%.v4b_web.stable_lights.avg_vis.tfw
|___v4unzip.bat

```



```
|
|
|
|
|rad_cal.tar
|phys_geo.zip
|world_avg_fixed.tfw
|af_grumpv1_ppoints_csv.zip
|gl_grumpv1_area_ascii_30.zip
|malaria_ecology.zip
|gl_gpww3_pcount_00_wrk_25.zip
|gasflares.zip
|F%sy%.v4.tar
|boundaries.zip
|F%sy%.v4b_web.stable_lights.avg_vis.tif
|__v4unzip.bat
|__rad_cal.tar
|__phys_geo.zip
|__world_avg_fixed.tfw
|__af_grumpv1_ppoints_csv.zip
|__gl_grumpv1_area_ascii_30.zip
|__malaria_ecology.zip
|__gl_gpww3_pcount_00_wrk_25.zip
|__gasflares.zip
|__F%sy%.v4.tar
|__boundaries.zip
|__world_avg.tfw
|__v4unzip.bat
|__rad_cal.tar
|__phys_geo.zip
|__world_avg_fixed.tfw
|__af_grumpv1_ppoints_csv.zip
|__gl_grumpv1_area_ascii_30.zip
|__malaria_ecology.zip
|__gl_gpww3_pcount_00_wrk_25.zip
|__gasflares.zip
|__F%sy%.v4.tar
|__boundaries.zip
|__gluareag_alpha1.asc
|__v4unzip.bat
|__rad_cal.tar
|__phys_geo.zip
|__world_avg_fixed.tfw
|__af_grumpv1_ppoints_csv.zip
|__gl_grumpv1_area_ascii_30.zip
|__malaria_ecology.zip
|__gl_gpww3_pcount_00_wrk_25.zip
|__gasflares.zip
|__F%sy%.v4.tar
|__boundaries.zip
```

```

|___F%sy%.v4b_web.cf_cvg.tif
|   |___v4unzip.bat
|       |___rad_cal.tar
|       |___phys_geo.zip
|       |___world_avg_fixed.tfw
|       |___af_grumpv1_ppoints_csv.zip
|       |___gl_grumpv1_area_ascii_30.zip
|       |___malaria_ecology.zip
|       |___gl_gpww3_pcount_00_wrk_25.zip
|       |___gasflares.zip
|       |___F%sy%.v4.tar
|       |___boundaries.zip
|___F%sy%.v4b_web.cf_cvg.tfw
|   |___v4unzip.bat
|       |___rad_cal.tar
|       |___phys_geo.zip
|       |___world_avg_fixed.tfw
|       |___af_grumpv1_ppoints_csv.zip
|       |___gl_grumpv1_area_ascii_30.zip
|       |___malaria_ecology.zip
|       |___gl_gpww3_pcount_00_wrk_25.zip
|       |___gasflares.zip
|       |___F%sy%.v4.tar
|       |___boundaries.zip
|___world_avg.tif
|   |___v4unzip.bat
|       |___rad_cal.tar
|       |___phys_geo.zip
|       |___world_avg_fixed.tfw
|       |___af_grumpv1_ppoints_csv.zip
|       |___gl_grumpv1_area_ascii_30.zip
|       |___malaria_ecology.zip
|       |___gl_gpww3_pcount_00_wrk_25.zip
|       |___gasflares.zip
|       |___F%sy%.v4.tar
|       |___boundaries.zip
|___gluareag
|   |___v4lightsprep.aml
|       |___F%sy%.v4b_web.stable_lights.avg_vis.tfw
|           |___v4unzip.bat
|               |___rad_cal.tar
|               |___phys_geo.zip
|               |___world_avg_fixed.tfw
|               |___af_grumpv1_ppoints_csv.zip
|               |___gl_grumpv1_area_ascii_30.zip
|               |___malaria_ecology.zip

```

```
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar  
| | | | |___boundaries.zip  
| | | ___F%sy%.v4b_web.stable_lights.avg_vis.tif  
| | | |___v4unzip.bat  
| | | | |___rad_cal.tar  
| | | | |___phys_geo.zip  
| | | | |___world_avg_fixed.tfw  
| | | | |___af_grumpv1_ppoints_csv.zip  
| | | | |___gl_grumpv1_area_ascii_30.zip  
| | | | |___malaria_ecology.zip  
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar  
| | | | |___boundaries.zip  
| | | ___world_avg.tfw  
| | | |___v4unzip.bat  
| | | | |___rad_cal.tar  
| | | | |___phys_geo.zip  
| | | | |___world_avg_fixed.tfw  
| | | | |___af_grumpv1_ppoints_csv.zip  
| | | | |___gl_grumpv1_area_ascii_30.zip  
| | | | |___malaria_ecology.zip  
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar  
| | | | |___boundaries.zip  
| | | ___gluareag_alpha1.asc  
| | | |___v4unzip.bat  
| | | | |___rad_cal.tar  
| | | | |___phys_geo.zip  
| | | | |___world_avg_fixed.tfw  
| | | | |___af_grumpv1_ppoints_csv.zip  
| | | | |___gl_grumpv1_area_ascii_30.zip  
| | | | |___malaria_ecology.zip  
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar  
| | | | |___boundaries.zip  
| | | ___F%sy%.v4b_web.cf_cvg.tif  
| | | |___v4unzip.bat  
| | | | |___rad_cal.tar  
| | | | |___phys_geo.zip  
| | | | |___world_avg_fixed.tfw  
| | | | |___af_grumpv1_ppoints_csv.zip
```

```
| | | | |___gl_grumpv1_area_ascii_30.zip  
| | | | |___malaria_ecology.zip  
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar  
| | | | |___boundaries.zip  
| | | ___F%sy%.v4b_web.cf_cvg.tfw  
| | | |___v4unzip.bat  
| | | | |___rad_cal.tar  
| | | | |___phys_geo.zip  
| | | | |___world_avg_fixed.tfw  
| | | | |___af_grumpv1_ppoints_csv.zip  
| | | | |___gl_grumpv1_area_ascii_30.zip  
| | | | |___malaria_ecology.zip  
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar  
| | | | |___boundaries.zip  
| | | ___world_avg.tif  
| | | |___v4unzip.bat  
| | | | |___rad_cal.tar  
| | | | |___phys_geo.zip  
| | | | |___world_avg_fixed.tfw  
| | | | |___af_grumpv1_ppoints_csv.zip  
| | | | |___gl_grumpv1_area_ascii_30.zip  
| | | | |___malaria_ecology.zip  
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar  
| | | | |___boundaries.zip  
| | ___ctrynum.dbf  
| ___gasflares.shp  
| | ___merge_gf.pyw  
| | | ___afprimates.csv  
| | | |___pre_aml.do  
| | | |___afpv1.csv  
| | | | |___v4unzip.bat  
| | | | |___rad_cal.tar  
| | | | |___phys_geo.zip  
| | | | |___world_avg_fixed.tfw  
| | | | |___af_grumpv1_ppoints_csv.zip  
| | | | |___gl_grumpv1_area_ascii_30.zip  
| | | | |___malaria_ecology.zip  
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar
```



```
| | | | |rad_cal.tar  
| | | | |__phys_geo.zip  
| | | | |__world_avg_fixed.tfw  
| | | | |__af_grumpv1_ppoints_csv.zip  
| | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | |__malaria_ecology.zip  
| | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | |__gasflares.zip  
| | | | |F%sy%.v4.tar  
| | | | |__boundaries.zip  
| | | | |__F%sy%.v4b_web.stable_lights.avg_vis.tif  
| | | | | |__v4unzip.bat  
| | | | | |__rad_cal.tar  
| | | | | |__phys_geo.zip  
| | | | | |__world_avg_fixed.tfw  
| | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | |__malaria_ecology.zip  
| | | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | | |__gasflares.zip  
| | | | | |F%sy%.v4.tar  
| | | | | |__boundaries.zip  
| | | | |__world_avg.tfw  
| | | | | |__v4unzip.bat  
| | | | | |__rad_cal.tar  
| | | | | |__phys_geo.zip  
| | | | | |__world_avg_fixed.tfw  
| | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | |__malaria_ecology.zip  
| | | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | | |__gasflares.zip  
| | | | | |F%sy%.v4.tar  
| | | | | |__boundaries.zip  
| | | | |__gluareag_alpha1.asc  
| | | | | |__v4unzip.bat  
| | | | | |__rad_cal.tar  
| | | | | |__phys_geo.zip  
| | | | | |__world_avg_fixed.tfw  
| | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | |__malaria_ecology.zip  
| | | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | | |__gasflares.zip  
| | | | | |F%sy%.v4.tar  
| | | | | |__boundaries.zip
```



```
| | | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | | |___gasflares.zip  
| | | | |___F%sy%.v4.tar  
| | | | |___boundaries.zip  
| | | ___F%sy%.v4b_web.stable_lights.avg_vis.tif  
| | | |___v4unzip.bat  
| | | |___rad_cal.tar  
| | | |___phys_geo.zip  
| | | |___world_avg_fixed.tfw  
| | | |___af_grumpv1_ppoints_csv.zip  
| | | |___gl_grumpv1_area_ascii_30.zip  
| | | |___malaria_ecology.zip  
| | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | |___gasflares.zip  
| | | |___F%sy%.v4.tar  
| | | |___boundaries.zip  
| | | ___world_avg.tfw  
| | | |___v4unzip.bat  
| | | |___rad_cal.tar  
| | | |___phys_geo.zip  
| | | |___world_avg_fixed.tfw  
| | | |___af_grumpv1_ppoints_csv.zip  
| | | |___gl_grumpv1_area_ascii_30.zip  
| | | |___malaria_ecology.zip  
| | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | |___gasflares.zip  
| | | |___F%sy%.v4.tar  
| | | |___boundaries.zip  
| | | ___gluareag_alpha1.asc  
| | | |___v4unzip.bat  
| | | |___rad_cal.tar  
| | | |___phys_geo.zip  
| | | |___world_avg_fixed.tfw  
| | | |___af_grumpv1_ppoints_csv.zip  
| | | |___gl_grumpv1_area_ascii_30.zip  
| | | |___malaria_ecology.zip  
| | | |___gl_gpww3_pcount_00_wrk_25.zip  
| | | |___gasflares.zip  
| | | |___F%sy%.v4.tar  
| | | |___boundaries.zip  
| | | ___F%sy%.v4b_web.cf_cvg.tif  
| | | |___v4unzip.bat  
| | | |___rad_cal.tar  
| | | |___phys_geo.zip  
| | | |___world_avg_fixed.tfw  
| | | |___af_grumpv1_ppoints_csv.zip
```



```
| | | | |__rad_cal.tar  
| | | | |__phys_geo.zip  
| | | | |__world_avg_fixed.tfw  
| | | | |__af_grumpv1_ppoints_csv.zip  
| | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | |__malaria_ecology.zip  
| | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | |__gasflares.zip  
| | | | |__F%sy%.v4.tar  
| | | | |__boundaries.zip  
| | | | |__world_avg.tfw  
| | | | | |__v4unzip.bat  
| | | | | |__rad_cal.tar  
| | | | | |__phys_geo.zip  
| | | | | |__world_avg_fixed.tfw  
| | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | |__malaria_ecology.zip  
| | | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | | |__gasflares.zip  
| | | | | |__F%sy%.v4.tar  
| | | | | |__boundaries.zip  
| | | | |__gluareag_alpha1.asc  
| | | | | |__v4unzip.bat  
| | | | | |__rad_cal.tar  
| | | | | |__phys_geo.zip  
| | | | | |__world_avg_fixed.tfw  
| | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | |__malaria_ecology.zip  
| | | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | | |__gasflares.zip  
| | | | | |__F%sy%.v4.tar  
| | | | | |__boundaries.zip  
| | | | |__F%sy%.v4b_web.cf_cvg.tif  
| | | | | |__v4unzip.bat  
| | | | | |__rad_cal.tar  
| | | | | |__phys_geo.zip  
| | | | | |__world_avg_fixed.tfw  
| | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | |__malaria_ecology.zip  
| | | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | | |__gasflares.zip  
| | | | | |__F%sy%.v4.tar  
| | | | | |__boundaries.zip
```



```

|___world_avg_fixed.tfw
|___af_grumpv1_ppoints_csv.zip
|___gl_grumpv1_area_ascii_30.zip
|___malaria_ecology.zip
|___gl_gpww3_pcount_00_wrk_25.zip
|___gasflares.zip
|___F%sy%.v4.tar
|___boundaries.zip
|___glp00ag
|___v4unzip.bat
|___rad_cal.tar
|___phys_geo.zip
|___world_avg_fixed.tfw
|___af_grumpv1_ppoints_csv.zip
|___gl_grumpv1_area_ascii_30.zip
|___malaria_ecology.zip
|___gl_gpww3_pcount_00_wrk_25.zip
|___gasflares.zip
|___F%sy%.v4.tar
|___boundaries.zip
|___ginioutu.dbf
|___v4cvggini_bycountry_prep.aml
|___ctryliggrid
|___v4unzip.bat
|___rad_cal.tar
|___phys_geo.zip
|___world_avg_fixed.tfw
|___af_grumpv1_ppoints_csv.zip
|___gl_grumpv1_area_ascii_30.zip
|___malaria_ecology.zip
|___gl_gpww3_pcount_00_wrk_25.zip
|___gasflares.zip
|___F%sy%.v4.tar
|___boundaries.zip
|___l%sy%nm
|___v4lightsprep.aml
|___F%sy%.v4b_web.stable_lights.avg_vis.tfw
|___v4unzip.bat
|___rad_cal.tar
|___phys_geo.zip
|___world_avg_fixed.tfw
|___af_grumpv1_ppoints_csv.zip
|___gl_grumpv1_area_ascii_30.zip
|___malaria_ecology.zip
|___gl_gpww3_pcount_00_wrk_25.zip
|___gasflares.zip

```

```
| | | | F%sy%.v4.tar  
| | | | ___boundaries.zip  
|___F%sy%.v4b_web.stable_lights.avg_vis.tif  
|   |___v4unzip.bat  
|     |___rad_cal.tar  
|     |___phys_geo.zip  
|     |___world_avg_fixed.tfw  
|     |___af_grumpv1_ppoints_csv.zip  
|     |___gl_grumpv1_area_ascii_30.zip  
|     |___malaria_ecology.zip  
|     |___gl_gpww3_pcount_00_wrk_25.zip  
|     |___gasflares.zip  
|     |___F%sy%.v4.tar  
|     |___boundaries.zip  
|___world_avg.tfw  
|   |___v4unzip.bat  
|     |___rad_cal.tar  
|     |___phys_geo.zip  
|     |___world_avg_fixed.tfw  
|     |___af_grumpv1_ppoints_csv.zip  
|     |___gl_grumpv1_area_ascii_30.zip  
|     |___malaria_ecology.zip  
|     |___gl_gpww3_pcount_00_wrk_25.zip  
|     |___gasflares.zip  
|     |___F%sy%.v4.tar  
|     |___boundaries.zip  
|___gluareag_alpha1.asc  
|   |___v4unzip.bat  
|     |___rad_cal.tar  
|     |___phys_geo.zip  
|     |___world_avg_fixed.tfw  
|     |___af_grumpv1_ppoints_csv.zip  
|     |___gl_grumpv1_area_ascii_30.zip  
|     |___malaria_ecology.zip  
|     |___gl_gpww3_pcount_00_wrk_25.zip  
|     |___gasflares.zip  
|     |___F%sy%.v4.tar  
|     |___boundaries.zip  
|___F%sy%.v4b_web.cf_cvg.tif  
|   |___v4unzip.bat  
|     |___rad_cal.tar  
|     |___phys_geo.zip  
|     |___world_avg_fixed.tfw  
|     |___af_grumpv1_ppoints_csv.zip  
|     |___gl_grumpv1_area_ascii_30.zip  
|     |___malaria_ecology.zip
```

```

|__gl_gpwv3_pcount_00_wrk_25.zip
|__gasflares.zip
|__F%sy%.v4.tar
|__boundaries.zip
|__F%sy%.v4b_web.cf_cvg.tfw
|__v4unzip.bat
|__rad_cal.tar
|__phys_geo.zip
|__world_avg_fixed.tfw
|__af_grumpv1_ppoints_csv.zip
|__gl_grumpv1_area_ascii_30.zip
|__malaria_ecology.zip
|__gl_gpwv3_pcount_00_wrk_25.zip
|__gasflares.zip
|__F%sy%.v4.tar
|__boundaries.zip
|__world_avg.tif
|__v4unzip.bat
|__rad_cal.tar
|__phys_geo.zip
|__world_avg_fixed.tfw
|__af_grumpv1_ppoints_csv.zip
|__gl_grumpv1_area_ascii_30.zip
|__malaria_ecology.zip
|__gl_gpwv3_pcount_00_wrk_25.zip
|__gasflares.zip
|__F%sy%.v4.tar
|__boundaries.zip
|__cvg%sy%
|__v4lightsprep.aml
|__F%sy%.v4b_web.stable_lights.avg_vis.tfw
|__v4unzip.bat
|__rad_cal.tar
|__phys_geo.zip
|__world_avg_fixed.tfw
|__af_grumpv1_ppoints_csv.zip
|__gl_grumpv1_area_ascii_30.zip
|__malaria_ecology.zip
|__gl_gpwv3_pcount_00_wrk_25.zip
|__gasflares.zip
|__F%sy%.v4.tar
|__boundaries.zip
|__F%sy%.v4b_web.stable_lights.avg_vis.tif
|__v4unzip.bat
|__rad_cal.tar
|__phys_geo.zip

```



```
| | | | |__rad_cal.tar  
| | | | |__phys_geo.zip  
| | | | |__world_avg_fixed.tfw  
| | | | |__af_grumpv1_ppoints_csv.zip  
| | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | |__malaria_ecology.zip  
| | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | |__gasflares.zip  
| | | | |__F%sy%.v4.tar  
| | | | |__boundaries.zip  
| | | | |__world_avg.tif  
| | | | | |__v4unzip.bat  
| | | | | |__rad_cal.tar  
| | | | | |__phys_geo.zip  
| | | | | |__world_avg_fixed.tfw  
| | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | |__malaria_ecology.zip  
| | | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | | |__gasflares.zip  
| | | | | |__F%sy%.v4.tar  
| | | | | |__boundaries.zip  
| | | | |__isocvgst.dbf  
| | | | | |__pre_aml.do  
| | | | | | |__afpv1.csv  
| | | | | | | |__v4unzip.bat  
| | | | | | |__rad_cal.tar  
| | | | | | |__phys_geo.zip  
| | | | | | |__world_avg_fixed.tfw  
| | | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | | |__gl_grumpv1_area_ascii_30.zip  
| | | | | | |__malaria_ecology.zip  
| | | | | | |__gl_gpww3_pcount_00_wrk_25.zip  
| | | | | | |__gasflares.zip  
| | | | | | |__F%sy%.v4.tar  
| | | | | | |__boundaries.zip  
| | | | | |__ctrynum.dbf  
| | | | |__ginistrt.dbf  
| | | | | |__pre_aml.do  
| | | | | | |__afpv1.csv  
| | | | | | | |__v4unzip.bat  
| | | | | | |__rad_cal.tar  
| | | | | | |__phys_geo.zip  
| | | | | | |__world_avg_fixed.tfw  
| | | | | | |__af_grumpv1_ppoints_csv.zip  
| | | | | | |__gl_grumpv1_area_ascii_30.zip
```



```
| | | | | rad_cal.tar  
| | | | | phys_geo.zip  
| | | | | world_avg_fixed.tfw  
| | | | | af_grumpv1_ppoints_csv.zip  
| | | | | gl_grumpv1_area_ascii_30.zip  
| | | | | malaria_ecology.zip  
| | | | | gl_gpww3_pcount_00_wrk_25.zip  
| | | | | gasflares.zip  
| | | | | F%sy%.v4.tar  
| | | | | boundaries.zip  
| | | | | gluareag_alpha1.asc  
| | | | | v4unzip.bat  
| | | | | rad_cal.tar  
| | | | | phys_geo.zip  
| | | | | world_avg_fixed.tfw  
| | | | | af_grumpv1_ppoints_csv.zip  
| | | | | gl_grumpv1_area_ascii_30.zip  
| | | | | malaria_ecology.zip  
| | | | | gl_gpww3_pcount_00_wrk_25.zip  
| | | | | gasflares.zip  
| | | | | F%sy%.v4.tar  
| | | | | boundaries.zip  
| | | | | F%sy%.v4b_web.cf_cvg.tif  
| | | | | v4unzip.bat  
| | | | | rad_cal.tar  
| | | | | phys_geo.zip  
| | | | | world_avg_fixed.tfw  
| | | | | af_grumpv1_ppoints_csv.zip  
| | | | | gl_grumpv1_area_ascii_30.zip  
| | | | | malaria_ecology.zip  
| | | | | gl_gpww3_pcount_00_wrk_25.zip  
| | | | | gasflares.zip  
| | | | | F%sy%.v4.tar  
| | | | | boundaries.zip  
| | | | | F%sy%.v4b_web.cf_cvg.tfw  
| | | | | v4unzip.bat  
| | | | | rad_cal.tar  
| | | | | phys_geo.zip  
| | | | | world_avg_fixed.tfw  
| | | | | af_grumpv1_ppoints_csv.zip  
| | | | | gl_grumpv1_area_ascii_30.zip  
| | | | | malaria_ecology.zip  
| | | | | gl_gpww3_pcount_00_wrk_25.zip  
| | | | | gasflares.zip  
| | | | | F%sy%.v4.tar  
| | | | | boundaries.zip
```

[illegible]


```

|      |___F%sy%.v4b_web.cf_cvg.tif
|      |___v4unzip.bat
|      |___rad_cal.tar
|      |___phys_geo.zip
|      |___world_avg_fixed.tfw
|      |___af_grumpv1_ppoints_csv.zip
|      |___gl_grumpv1_area_ascii_30.zip
|      |___malaria_ecology.zip
|      |___gl_gpwv3_pcount_00_wrk_25.zip
|      |___gasflares.zip
|      |___F%sy%.v4.tar
|      |___boundaries.zip
|      |___F%sy%.v4b_web.cf_cvg.tfw
|      |___v4unzip.bat
|      |___rad_cal.tar
|      |___phys_geo.zip
|      |___world_avg_fixed.tfw
|      |___af_grumpv1_ppoints_csv.zip
|      |___gl_grumpv1_area_ascii_30.zip
|      |___malaria_ecology.zip
|      |___gl_gpwv3_pcount_00_wrk_25.zip
|      |___gasflares.zip
|      |___F%sy%.v4.tar
|      |___boundaries.zip
|      |___world_avg.tif
|      |___v4unzip.bat
|      |___rad_cal.tar
|      |___phys_geo.zip
|      |___world_avg_fixed.tfw
|      |___af_grumpv1_ppoints_csv.zip
|      |___gl_grumpv1_area_ascii_30.zip
|      |___malaria_ecology.zip
|      |___gl_gpwv3_pcount_00_wrk_25.zip
|      |___gasflares.zip
|      |___F%sy%.v4.tar
|      |___boundaries.zip
|___global_total_dn_uncal.dta
|___v4lights_stataprep_uncal.do
|___ctryoutu.dbf
|   |___v4ctrytables_uncal.aml
|   |___ctryliggrid
|   |   |___v4unzip.bat
|   |   |___rad_cal.tar
|   |   |___phys_geo.zip
|   |   |___world_avg_fixed.tfw
|   |   |___af_grumpv1_ppoints_csv.zip

```



```
| | | | | world_avg_fixed.tfw  
| | | | | ___af_grumpv1_ppoints_csv.zip  
| | | | | ___gl_grumpv1_area_ascii_30.zip  
| | | | | ___malaria_ecology.zip  
| | | | | ___gl_gpwv3_pcount_00_wrk_25.zip  
| | | | | ___gasflares.zip  
| | | | | ___F%sy%.v4.tar  
| | | | | ___boundaries.zip  
| | | ___F%sy%.v4b_web.stable_lights.avg_vis.tif  
| | | | ___v4unzip.bat  
| | | | | ___rad_cal.tar  
| | | | | ___phys_geo.zip  
| | | | | ___world_avg_fixed.tfw  
| | | | | ___af_grumpv1_ppoints_csv.zip  
| | | | | ___gl_grumpv1_area_ascii_30.zip  
| | | | | ___malaria_ecology.zip  
| | | | | ___gl_gpwv3_pcount_00_wrk_25.zip  
| | | | | ___gasflares.zip  
| | | | | ___F%sy%.v4.tar  
| | | | | ___boundaries.zip  
| | | ___world_avg.tfw  
| | | | ___v4unzip.bat  
| | | | | ___rad_cal.tar  
| | | | | ___phys_geo.zip  
| | | | | ___world_avg_fixed.tfw  
| | | | | ___af_grumpv1_ppoints_csv.zip  
| | | | | ___gl_grumpv1_area_ascii_30.zip  
| | | | | ___malaria_ecology.zip  
| | | | | ___gl_gpwv3_pcount_00_wrk_25.zip  
| | | | | ___gasflares.zip  
| | | | | ___F%sy%.v4.tar  
| | | | | ___boundaries.zip  
| | | ___gluareag_alpha1.asc  
| | | | ___v4unzip.bat  
| | | | | ___rad_cal.tar  
| | | | | ___phys_geo.zip  
| | | | | ___world_avg_fixed.tfw  
| | | | | ___af_grumpv1_ppoints_csv.zip  
| | | | | ___gl_grumpv1_area_ascii_30.zip  
| | | | | ___malaria_ecology.zip  
| | | | | ___gl_gpwv3_pcount_00_wrk_25.zip  
| | | | | ___gasflares.zip  
| | | | | ___F%sy%.v4.tar  
| | | | | ___boundaries.zip  
| | | ___F%sy%.v4b_web.cf_cvg.tif  
| | | | ___v4unzip.bat
```

```

|_|_|_|rad_cal.tar
|_|_|_|phys_geo.zip
|_|_|_|world_avg_fixed.tfw
|_|_|_|af_grumpv1_ppoints_csv.zip
|_|_|_|gl_grumpv1_area_ascii_30.zip
|_|_|_|malaria_ecology.zip
|_|_|_|gl_gpwv3_pcount_00_wrk_25.zip
|_|_|_|gasflares.zip
|_|_|_|F%sy%.v4.tar
|_|_|_|boundaries.zip
|_|_|F%sy%.v4b_web.cf_cvg.tfw
|_|_|_|v4unzip.bat
|_|_|_|_|rad_cal.tar
|_|_|_|_|phys_geo.zip
|_|_|_|_|world_avg_fixed.tfw
|_|_|_|_|af_grumpv1_ppoints_csv.zip
|_|_|_|_|gl_grumpv1_area_ascii_30.zip
|_|_|_|_|malaria_ecology.zip
|_|_|_|_|gl_gpwv3_pcount_00_wrk_25.zip
|_|_|_|_|gasflares.zip
|_|_|_|_|F%sy%.v4.tar
|_|_|_|_|boundaries.zip
|_|_|world_avg.tif
|_|_|_|v4unzip.bat
|_|_|_|_|rad_cal.tar
|_|_|_|_|phys_geo.zip
|_|_|_|_|world_avg_fixed.tfw
|_|_|_|_|af_grumpv1_ppoints_csv.zip
|_|_|_|_|gl_grumpv1_area_ascii_30.zip
|_|_|_|_|malaria_ecology.zip
|_|_|_|_|gl_gpwv3_pcount_00_wrk_25.zip
|_|_|_|_|gasflares.zip
|_|_|_|_|F%sy%.v4.tar
|_|_|_|_|boundaries.zip
|_|_|ctrynum.dbf
|_|_|gasflares.shp
|_|_|_|merge_gf.pyw
|_|_|_|_|afprimates.csv
|_|_|_|_|_|pre_aml.do
|_|_|_|_|_|afpv1.csv
|_|_|_|_|_|_|v4unzip.bat
|_|_|_|_|_|_|_|rad_cal.tar
|_|_|_|_|_|_|_|phys_geo.zip
|_|_|_|_|_|_|_|world_avg_fixed.tfw
|_|_|_|_|_|_|_|af_grumpv1_ppoints_csv.zip
|_|_|_|_|_|_|_|gl_grumpv1_area_ascii_30.zip

```

```

|      |      |      |      |__malaria_ecology.zip
|      |      |      |      |__gl_gpww3_pcount_00_wrk_25.zip
|      |      |      |      |__gasflares.zip
|      |      |      |      |__F%sy%.v4.tar
|      |      |      |      |__boundaries.zip
|      |      |      |      |__ctrynum.dbf
|      |      |      |      |__gas flares by country
|      |      |      |      |__v4unzip.bat
|      |      |      |      |__rad_cal.tar
|      |      |      |      |__phys_geo.zip
|      |      |      |      |__world_avg_fixed.tfw
|      |      |      |      |__af_grumpv1_ppoints_csv.zip
|      |      |      |      |__gl_grumpv1_area_ascii_30.zip
|      |      |      |      |__malaria_ecology.zip
|      |      |      |      |__gl_gpww3_pcount_00_wrk_25.zip
|      |      |      |      |__gasflares.zip
|      |      |      |      |__F%sy%.v4.tar
|      |      |      |      |__boundaries.zip
|      |      |      |      |__glp00ag
|      |      |      |      |__v4unzip.bat
|      |      |      |      |__rad_cal.tar
|      |      |      |      |__phys_geo.zip
|      |      |      |      |__world_avg_fixed.tfw
|      |      |      |      |__af_grumpv1_ppoints_csv.zip
|      |      |      |      |__gl_grumpv1_area_ascii_30.zip
|      |      |      |      |__malaria_ecology.zip
|      |      |      |      |__gl_gpww3_pcount_00_wrk_25.zip
|      |      |      |      |__gasflares.zip
|      |      |      |      |__F%sy%.v4.tar
|      |      |      |      |__boundaries.zip
|      |      |      |      |__dhselect.xls
|      |      |      |      |__wb_dq.xls
|      |      |      |      |__ctryout2.dbf
|      |      |      |      |__v4ctrytables_uncal.aml
|      |      |      |      |__ctryliggrid
|      |      |      |      |__v4unzip.bat
|      |      |      |      |__rad_cal.tar
|      |      |      |      |__phys_geo.zip
|      |      |      |      |__world_avg_fixed.tfw
|      |      |      |      |__af_grumpv1_ppoints_csv.zip
|      |      |      |      |__gl_grumpv1_area_ascii_30.zip
|      |      |      |      |__malaria_ecology.zip
|      |      |      |      |__gl_gpww3_pcount_00_wrk_25.zip
|      |      |      |      |__gasflares.zip
|      |      |      |      |__F%sy%.v4.tar
|      |      |      |      |__boundaries.zip

```

```

|_---1%sy%nm
|   |___v4lightsprep.aml
|       |___F%sy%.v4b_web.stable_lights.avg_vis.tfw
|           |___v4unzip.bat
|               |___rad_cal.tar
|               |___phys_geo.zip
|               |___world_avg.fixed.tfw
|               |___af_grumpv1_ppoints_csv.zip
|               |___gl_grumpv1_area_ascii_30.zip
|               |___malaria_ecology.zip
|               |___gl_gpww3_pcount_00_wrk_25.zip
|               |___gasflares.zip
|               |___F%sy%.v4.tar
|               |___boundaries.zip
|       |___F%sy%.v4b_web.stable_lights.avg_vis.tif
|           |___v4unzip.bat
|               |___rad_cal.tar
|               |___phys_geo.zip
|               |___world_avg.fixed.tfw
|               |___af_grumpv1_ppoints_csv.zip
|               |___gl_grumpv1_area_ascii_30.zip
|               |___malaria_ecology.zip
|               |___gl_gpww3_pcount_00_wrk_25.zip
|               |___gasflares.zip
|               |___F%sy%.v4.tar
|               |___boundaries.zip
|       |___world_avg.tfw
|           |___v4unzip.bat
|               |___rad_cal.tar
|               |___phys_geo.zip
|               |___world_avg.fixed.tfw
|               |___af_grumpv1_ppoints_csv.zip
|               |___gl_grumpv1_area_ascii_30.zip
|               |___malaria_ecology.zip
|               |___gl_gpww3_pcount_00_wrk_25.zip
|               |___gasflares.zip
|               |___F%sy%.v4.tar
|               |___boundaries.zip
|       |___gluareag_alpha1.asc
|           |___v4unzip.bat
|               |___rad_cal.tar
|               |___phys_geo.zip
|               |___world_avg.fixed.tfw
|               |___af_grumpv1_ppoints_csv.zip
|               |___gl_grumpv1_area_ascii_30.zip
|               |___malaria_ecology.zip

```



```
| | | | |gl_grumpv1_area_ascii_30.zip  
| | | | |__malaria_ecology.zip  
| | | | |__gl_gpwv3_pcount_00_wrk_25.zip  
| | | | |__gasflares.zip  
| | | | |__F%sy%.v4.tar  
| | | | |__boundaries.zip  
|___F%sy%.v4b_web.stable_lights.avg_vis.tif  
|   |___v4unzip.bat  
|     |___rad_cal.tar  
|     |___phys_geo.zip  
|     |___world_avg_fixed.tfw  
|     |___af_grumpv1_ppoints_csv.zip  
|     |___gl_grumpv1_area_ascii_30.zip  
|     |___malaria_ecology.zip  
|     |___gl_gpwv3_pcount_00_wrk_25.zip  
|     |___gasflares.zip  
|     |___F%sy%.v4.tar  
|     |___boundaries.zip  
|___world_avg.tfw  
|   |___v4unzip.bat  
|     |___rad_cal.tar  
|     |___phys_geo.zip  
|     |___world_avg_fixed.tfw  
|     |___af_grumpv1_ppoints_csv.zip  
|     |___gl_grumpv1_area_ascii_30.zip  
|     |___malaria_ecology.zip  
|     |___gl_gpwv3_pcount_00_wrk_25.zip  
|     |___gasflares.zip  
|     |___F%sy%.v4.tar  
|     |___boundaries.zip  
|___gluareag_alpha1.asc  
|   |___v4unzip.bat  
|     |___rad_cal.tar  
|     |___phys_geo.zip  
|     |___world_avg_fixed.tfw  
|     |___af_grumpv1_ppoints_csv.zip  
|     |___gl_grumpv1_area_ascii_30.zip  
|     |___malaria_ecology.zip  
|     |___gl_gpwv3_pcount_00_wrk_25.zip  
|     |___gasflares.zip  
|     |___F%sy%.v4.tar  
|     |___boundaries.zip  
|___F%sy%.v4b_web.cf_cvg.tif  
|   |___v4unzip.bat  
|     |___rad_cal.tar  
|     |___phys_geo.zip
```



```
|
|
|      |___F%sy%.v4.tar
|      |___boundaries.zip
|    ___F%sy%.v4b_web.cf_cvg.tfw
|      |___v4unzip.bat
|      |___rad_cal.tar
|      |___phys_geo.zip
|      |___world_avg_fixed.tfw
|      |___af_grumpv1_ppoints_csv.zip
|      |___gl_grumpv1_area_ascii_30.zip
|      |___malaria_ecology.zip
|      |___gl_gpww3_pcount_00_wrk_25.zip
|      |___gasflares.zip
|      |___F%sy%.v4.tar
|      |___boundaries.zip
|    ___world_avg.tif
|      |___v4unzip.bat
|      |___rad_cal.tar
|      |___phys_geo.zip
|      |___world_avg_fixed.tfw
|      |___af_grumpv1_ppoints_csv.zip
|      |___gl_grumpv1_area_ascii_30.zip
|      |___malaria_ecology.zip
|      |___gl_gpww3_pcount_00_wrk_25.zip
|      |___gasflares.zip
|      |___F%sy%.v4.tar
|      |___boundaries.zip
|  ___isocvgst.dbf
|    |___pre_aml.do
|      |___afpv1.csv
|        |___v4unzip.bat
|          |___rad_cal.tar
|          |___phys_geo.zip
|          |___world_avg_fixed.tfw
|          |___af_grumpv1_ppoints_csv.zip
|          |___gl_grumpv1_area_ascii_30.zip
|          |___malaria_ecology.zip
|          |___gl_gpww3_pcount_00_wrk_25.zip
|          |___gasflares.zip
|          |___F%sy%.v4.tar
|          |___boundaries.zip
|        |___ctrynum.dbf
|  ___ginistrt.dbf
|    |___pre_aml.do
|      |___afpv1.csv
|        |___v4unzip.bat
|          |___rad_cal.tar
```

```

|                                     |__phys_geo.zip
|                                     |__world_avg_fixed.tfw
|                                     |__af_grumpv1_ppoints_csv.zip
|                                     |__gl_grumpv1_area_ascii_30.zip
|                                     |__malaria_ecology.zip
|                                     |__gl_gpwv3_pcount_00_wrk_25.zip
|                                     |__gasflares.zip
|                                     |__F%sy%.v4.tar
|                                     |__boundaries.zip
|                                     |__ctrynum.dbf
|__wdi_limited.dta
|__imf_dds.xls
|__ginistata_uncal.dta
|  |__v4ginicalc_uncal.do
|  |  |__ginioutu.dbf
|  |    |__v4cvggini_bycountry_prep.aml
|  |      |__ctryliggrid
|  |        |__v4unzip.bat
|  |          |__rad_cal.tar
|  |            |__phys_geo.zip
|  |              |__world_avg_fixed.tfw
|  |                |__af_grumpv1_ppoints_csv.zip
|  |                  |__gl_grumpv1_area_ascii_30.zip
|  |                    |__malaria_ecology.zip
|  |                      |__gl_gpwv3_pcount_00_wrk_25.zip
|  |                        |__gasflares.zip
|  |                          |__F%sy%.v4.tar
|  |                            |__boundaries.zip
|  |__l%sy%nm
|  |  |__v4lightsprep.aml
|  |    |__F%sy%.v4b_web.stable_lights.avg_vis.tif
|  |      |__v4unzip.bat
|  |        |__rad_cal.tar
|  |          |__phys_geo.zip
|  |            |__world_avg_fixed.tfw
|  |              |__af_grumpv1_ppoints_csv.zip
|  |                |__gl_grumpv1_area_ascii_30.zip
|  |                  |__malaria_ecology.zip
|  |                    |__gl_gpwv3_pcount_00_wrk_25.zip
|  |                      |__gasflares.zip
|  |                        |__F%sy%.v4.tar
|  |                          |__boundaries.zip
|  |__F%sy%.v4b_web.stable_lights.avg_vis.tif
|  |  |__v4unzip.bat
|  |    |__rad_cal.tar
|  |      |__phys_geo.zip

```

```

|         |         |___world_avg_fixed.tfw
|         |         |___af_grumpv1_ppoints_csv.zip
|         |         |___gl_grumpv1_area_ascii_30.zip
|         |         |___malaria_ecology.zip
|         |         |___gl_gpww3_pcount_00_wrk_25.zip
|         |         |___gasflares.zip
|         |         |___F%sy%.v4.tar
|         |         |___boundaries.zip
|___world_avg.tfw
|   |___v4unzip.bat
|         |         |___rad_cal.tar
|         |         |___phys_geo.zip
|         |         |___world_avg_fixed.tfw
|         |         |___af_grumpv1_ppoints_csv.zip
|         |         |___gl_grumpv1_area_ascii_30.zip
|         |         |___malaria_ecology.zip
|         |         |___gl_gpww3_pcount_00_wrk_25.zip
|         |         |___gasflares.zip
|         |         |___F%sy%.v4.tar
|         |         |___boundaries.zip
|___gluareag_alpha1.asc
|   |___v4unzip.bat
|         |         |___rad_cal.tar
|         |         |___phys_geo.zip
|         |         |___world_avg_fixed.tfw
|         |         |___af_grumpv1_ppoints_csv.zip
|         |         |___gl_grumpv1_area_ascii_30.zip
|         |         |___malaria_ecology.zip
|         |         |___gl_gpww3_pcount_00_wrk_25.zip
|         |         |___gasflares.zip
|         |         |___F%sy%.v4.tar
|         |         |___boundaries.zip
|___F%sy%.v4b_web.cf_cvg.tif
|   |___v4unzip.bat
|         |         |___rad_cal.tar
|         |         |___phys_geo.zip
|         |         |___world_avg_fixed.tfw
|         |         |___af_grumpv1_ppoints_csv.zip
|         |         |___gl_grumpv1_area_ascii_30.zip
|         |         |___malaria_ecology.zip
|         |         |___gl_gpww3_pcount_00_wrk_25.zip
|         |         |___gasflares.zip
|         |         |___F%sy%.v4.tar
|         |         |___boundaries.zip
|___F%sy%.v4b_web.cf_cvg.tfw
|   |___v4unzip.bat

```

```

|      |      |___rad_cal.tar
|      |      |___phys_geo.zip
|      |      |___world_avg_fixed.tfw
|      |      |___af_grumpv1_ppoints_csv.zip
|      |      |___gl_grumpv1_area_ascii_30.zip
|      |      |___malaria_ecology.zip
|      |      |___gl_gpwv3_pcount_00_wrk_25.zip
|      |      |___gasflares.zip
|      |      |___F%sy%.v4.tar
|      |      |___boundaries.zip
|      |___world_avg.tif
|      |___v4unzip.bat
|      |      |___rad_cal.tar
|      |      |___phys_geo.zip
|      |      |___world_avg_fixed.tfw
|      |      |___af_grumpv1_ppoints_csv.zip
|      |      |___gl_grumpv1_area_ascii_30.zip
|      |      |___malaria_ecology.zip
|      |      |___gl_gpwv3_pcount_00_wrk_25.zip
|      |      |___gasflares.zip
|      |      |___F%sy%.v4.tar
|      |      |___boundaries.zip
|___cvg%sy%
|   |___v4lightsprep.aml
|   |   |___F%sy%.v4b_web.stable_lights.avg_vis.tif
|   |   |   |___v4unzip.bat
|   |   |   |___rad_cal.tar
|   |   |   |___phys_geo.zip
|   |   |   |___world_avg_fixed.tfw
|   |   |   |___af_grumpv1_ppoints_csv.zip
|   |   |   |___gl_grumpv1_area_ascii_30.zip
|   |   |   |___malaria_ecology.zip
|   |   |   |___gl_gpwv3_pcount_00_wrk_25.zip
|   |   |   |___gasflares.zip
|   |   |   |___F%sy%.v4.tar
|   |   |   |___boundaries.zip
|   |   |___F%sy%.v4b_web.stable_lights.avg_vis.tif
|   |   |___v4unzip.bat
|   |   |   |___rad_cal.tar
|   |   |   |___phys_geo.zip
|   |   |   |___world_avg_fixed.tfw
|   |   |   |___af_grumpv1_ppoints_csv.zip
|   |   |   |___gl_grumpv1_area_ascii_30.zip
|   |   |   |___malaria_ecology.zip
|   |   |   |___gl_gpwv3_pcount_00_wrk_25.zip
|   |   |   |___gasflares.zip

```

```

|           |___F%sy%.v4.tar
|           |___boundaries.zip
|___world_avg.tfw
|   |___v4unzip.bat
|       |___rad_cal.tar
|       |___phys_geo.zip
|       |___world_avg_fixed.tfw
|       |___af_grumpv1_ppoints_csv.zip
|       |___gl_grumpv1_area_ascii_30.zip
|       |___malaria_ecology.zip
|       |___gl_gpww3_pcount_00_wrk_25.zip
|       |___gasflares.zip
|       |___F%sy%.v4.tar
|       |___boundaries.zip
|___gluareag_alpha1.asc
|   |___v4unzip.bat
|       |___rad_cal.tar
|       |___phys_geo.zip
|       |___world_avg_fixed.tfw
|       |___af_grumpv1_ppoints_csv.zip
|       |___gl_grumpv1_area_ascii_30.zip
|       |___malaria_ecology.zip
|       |___gl_gpww3_pcount_00_wrk_25.zip
|       |___gasflares.zip
|       |___F%sy%.v4.tar
|       |___boundaries.zip
|___F%sy%.v4b_web.cf_cvg.tif
|   |___v4unzip.bat
|       |___rad_cal.tar
|       |___phys_geo.zip
|       |___world_avg_fixed.tfw
|       |___af_grumpv1_ppoints_csv.zip
|       |___gl_grumpv1_area_ascii_30.zip
|       |___malaria_ecology.zip
|       |___gl_gpww3_pcount_00_wrk_25.zip
|       |___gasflares.zip
|       |___F%sy%.v4.tar
|       |___boundaries.zip
|___F%sy%.v4b_web.cf_cvg.tif
|   |___v4unzip.bat
|       |___rad_cal.tar
|       |___phys_geo.zip
|       |___world_avg_fixed.tfw
|       |___af_grumpv1_ppoints_csv.zip
|       |___gl_grumpv1_area_ascii_30.zip
|       |___malaria_ecology.zip

```

```

|      |      |___gl_gpww3_pcount_00_wrk_25.zip
|      |      |___gasflares.zip
|      |      |___F%sy%.v4.tar
|      |      |___boundaries.zip
|      |___world_avg.tif
|      |___v4unzip.bat
|      |___rad_cal.tar
|      |___phys_geo.zip
|      |___world_avg_fixed.tfw
|      |___af_grumpv1_ppoints_csv.zip
|      |___gl_grumpv1_area_ascii_30.zip
|      |___malaria_ecology.zip
|      |___gl_gpww3_pcount_00_wrk_25.zip
|      |___gasflares.zip
|      |___F%sy%.v4.tar
|      |___boundaries.zip
|___isocvgst.dbf
|  |___pre_aml.do
|  |  |___afpv1.csv
|  |  |  |___v4unzip.bat
|  |  |  |___rad_cal.tar
|  |  |  |___phys_geo.zip
|  |  |  |___world_avg_fixed.tfw
|  |  |  |___af_grumpv1_ppoints_csv.zip
|  |  |  |___gl_grumpv1_area_ascii_30.zip
|  |  |  |___malaria_ecology.zip
|  |  |  |___gl_gpww3_pcount_00_wrk_25.zip
|  |  |  |___gasflares.zip
|  |  |  |___F%sy%.v4.tar
|  |  |  |___boundaries.zip
|  |  |___ctrynum.dbf
|___ginistrt.dbf
|  |___pre_aml.do
|  |  |___afpv1.csv
|  |  |  |___v4unzip.bat
|  |  |  |___rad_cal.tar
|  |  |  |___phys_geo.zip
|  |  |  |___world_avg_fixed.tfw
|  |  |  |___af_grumpv1_ppoints_csv.zip
|  |  |  |___gl_grumpv1_area_ascii_30.zip
|  |  |  |___malaria_ecology.zip
|  |  |  |___gl_gpww3_pcount_00_wrk_25.zip
|  |  |  |___gasflares.zip
|  |  |  |___F%sy%.v4.tar
|  |  |  |___boundaries.zip
|  |___ctrynum.dbf

```

8.2.3 Example of completing reproduction tree

In many cases, some of the components of the workflow will not be easily identifiable (or missing) in the reproduction package. Here we present a more complex example than the one presented in the Assessment chapter. The Diagram Builder will return a partial reproduction tree diagram. For example, if the files `merge_1_2.do`, `merge_3_4.do`, and `final_merge.do` are missing from the previous diagram, the ACRE Diagram Builder will produce the following diagram:

```
cleaned_3.dta
  [code] clean_raw_3.py
      raw_3.dta

table1.tex
  [code] analysis.R
      analysis_data.dta

cleaned_3_4.dta
  [code] clean_merged_3_4.do
      merged_3_4.dta

cleaned_1.dta
  [code] clean_raw_1.py
      raw_1.dta

cleaned_2.dta
  [code] clean_raw_2.py
      raw_2.dta

cleaned_4.dta
  [code] clean_raw_4.py
      raw_4.dta

cleaned_1_2.dta
  [code] clean_merged_1_2.do
      merged_1_2.dta
Unused data sources: None.
```

In this case, you can still manually combine this partial information with your knowledge from the paper and own judgement to produce a “candidate” tree diagram (which might lead to different reproducers recreating different diagrams). This may look like the following:

```
table1.tex
  [code] analysis.R
      analysis_data.dta
      MISSING_CODE_FILE_3
```

Table 8.1: Adding rows to code spreadsheet

file_name	location	inputs	outputs	description	primary_type
...
missing_file1	unknown	cleaned_1.dta	merged_1_2.dta	missing code	unknown
missing_file2	unknown	cleaned_3.dta	merged_3_4.dta	missing code	unknown
missing_file3	unknown	merged_3_4.dta	analysis_data.dta	missing code	unknown

```

cleaned_3_4.dta
|
|   [code] clean_merged_3_4.do
|   merged_3_4.dta
|   MISSSING_CODE_FILE_2
|   cleaned_3.dta
|   |   [code] clean_raw_3.py
|   |   raw_3.dta
|   cleaned_4.dta
|   |   [code] clean_raw_4.py
|   |   raw_4.dta
|
cleaned_1_2.dta
|   [code] clean_merged_1_2.do
|   merged_1_2.dta
|   MISSSING_CODE_FILE_1
|   cleaned_1.dta
|   |   [code] clean_raw_1.py
|   |   raw_1.dta
|   |
|   cleaned_2.dta
|   |   [code] clean_raw_2.py
|   |   raw_2.dta

```

To leave a record of the reconstructed diagrams, you will have to amend the input spreadsheets using placeholders for the missing components. In the example above, you should add the following entries to the code description spreadsheet:

Chapter 9

Tips and Resources for Reproducible Workflow

9.1 Reproducible workflows:

Below is a summary from Chapter 11 of Christensen et al. (2019). If there is a book/chapter that you find particularly helpful for, please write a brief summary and submit a contribution.

Folder organization

Basic file organization is a critical component of a reproducible workflow. The following structure is recommended, but can be adapted to accommodate different reproducers or types of research. The name of the master folder should be easy to read and meaningful to all collaborators on the the project.

- Create a master folder with a descriptive name for the project. It should contain:
 - Separate folders for programming script files, raw data, edited data, output, and final paper or article text
 - A README file: description of contents of each folder, as well as installation and operating instructions for reproducers
- Keep raw data intact: Any edits or datasets generated using raw data should be stored in a “data” folder separate from the “raw data” folder.
- When naming a directory or file, stick to lowercase letters with underscores (instead of spaces) to avoid cross-operating-system issues.

Efficient and readable programming

The core of programming for reproducibility is to write code wherever possible. Writing scripts leaves a record of any changes to data, which allows other researchers to reproduce work exactly. It is also helpful to leave comments in your

code to explain the reasoning for changes or any gaps left if using point-and-click methods is necessary.

- Leave a record of any changes to the data: Write code in the programming environment, instead of modifying data by hand in a spreadsheet or relying on point-and-click options.
- Include comments in code to explain changes, and save intermediate datasets used in analysis.
- Give variables names that will be informative to reproducers.
- Use relative directory paths, not absolute paths, so the work can be more easily reproduced from different computers.

Version control

Version control software is used to keep a record of changes to project files. Although it is possible to manually track changes in a central research log or as notes in individual script files, many social scientists recognize the benefits of a distributed version control system. Because each collaborator is able to have a local copy of the project's entire work history, these systems are particularly suited to collaborative projects. Below are methods to manually track changes and a brief explanation of Git, a popular distributed version control system.

- Maintain a written record of work.
 - In a central research log: Log activities in a single central file as often as work on the project is being done (keep track of “which team member writes what code, produces what output, edits which files, and when”).
 - In individual script files: Record “who edited which part of which file when, and why.”
 - With a version control system, such as Git: Git records changes made to files, by whom, and when.
- A brief explanation of Git: Users add changed files to the staging area, then commit those changes to the project folder, or repository. Git keeps the filename and records the new version of each file from the staging area.

9.2 Links to resources, organizations and people for reproducible work

Below we point to an ever growing list of resources, organizations and specific researcher doing empirical work with a strong orientation towards reproducibility. The list are alphabetical order and contributions are welcome!

9.2.1 Resources

- Dynamic documents in R, Python and Stata
- Git resources:

9.2. LINKS TO RESOURCES, ORGANIZATIONS AND PEOPLE FOR REPRODUCIBLE WORK131

- Jenny Bryan’s book and video
- Github learning lab
- Udacity’s intro
- Git for poets
- Combining GitHub and Dropbox
- Atlassian intro to Git
- Software Carpentry tutorial from the command line
- IDB’s cheatsheet for transparency, reproducibility and ethics
- Lars Vilhuber LDI’s Wiki for Reproducibility. Particularly this section.
- Open Science Framework (OSF)
- Project TIER
- R for Stata users
- World Bank DIME’s Wiki for transparent and reproducible research.

9.2.2 Organizations

- Congressional Budget Office
- Gentzkow & Shapiro Lab
- LOST
- Opportunity Lab
- Policy Simulation Library
- Urban Institute

9.2.3 People

(by last name)

- Luiza Andrade
- Alvaro Carril
- Lachlan Deer
- Rebekah Din
- Richard Evans
- Andrew Heiss
- John Horton
- Nick Huntingon
- Matt Jensen
- Max Kasy
- Cora Kingdon

- Grant McDermott
 - Tyler Ransom
 - Lisa Rennels
 - Ed Rubin
 - Michael Stepner
 - Shoshana Vasserman
-
- Lars Vilhuber

Chapter 10

Contributions

10.1 Contributing feedback on these guidelines

We welcome feedback from participants and the wider social science community. If you wish to provide feedback on specific chapters or sections, click the “edit” icon at the top of this page (this will prompt you to sign into or create a GitHub account), after which you’ll be able to suggest changes directly to the text. Please submit your suggestions using the “create a new branch and start a pull request” option and provide a summary of the changes you’ve proposed in the description of the pull request. The ACRE project team will review all suggested changes and decide whether to “push” them to this Guide document or not. For more general feedback, please contact ACRE@berkeley.edu.

Major contributions to this Guide will be acknowledged below. This project employs the Contributor Roles Taxonomy (CRediT). Major contributions are defined as any pushed revisions to the Guide language or source code beyond corrections of spelling and grammar.

10.2 List of Contributors: Guidelines content and source code:

(in alphabetical order)

- Aleksandar Bogdanoski – Funding acquisition, Project administration, Writing (original draft), Writing (reviewing and editing)
- Carson Christiano – Funding acquisition, Project administration, Writing (reviewing and editing)
- Joel Ferguson – Writing (original draft), Writing (reviewing and editing)
- Fernando Hoces de la Guardia – Conceptualization, Funding acquisition, Writing (original draft), Writing (reviewing and editing)

- Katherine Hoerberling – Funding acquisition, Project administration, Writing (original draft), Writing (reviewing and editing)
- Edward Miguel – Conceptualization, Funding acquisition, Supervision
- Emma Ng – Visualization, Writing (original draft), Writing (reviewing and editing)
- Lars Vilhuber – Conceptualization, Funding acquisition, Supervision

The individuals below have contributed to the ACRE GitHub repository:

@abogdanoski, @albertchae, Emma Ng (@em-ng21), Fernando Hoces de la Guardia (@fhoces), Joel Ferguson (@joelferg), Katie Hoerberling (@khoeberling), Michael Weiss (@mweiss)

10.3 Suggested citation format

These follow the Chicago style citation format. Brackets indicate where you'll need to input more specific information. - When citing reproductions accessed on the Social Science Reproduction Platform: [Reproducer Last name], [Reproducer First name], [First Last], and [First Last]. [Year]. "Reproduction of [Title of original paper]." *Social Science Reproduction Platform*. doi: [doi]. [DOI or link to original paper]. - When citing these Guidelines: Berkeley Initiative for Transparency in the Social Sciences. 2020. "Guide for Advancing Computational Reproducibility in the Social Sciences." [Date accessed (Day, Month, Year)]. <https://bitss.github.io/ACRE/>.

10.4 Acknowledgments

Support for the development of this Guide was provided by Arnold Ventures.

Chapter 11

Definitions

11.1 Concepts in reproducibility

- **Analytic data** – Data used as the final input in a workflow in order to produce a statistic displayed in the paper (including appendices).
- **Claim (concept)** – A major hypothesis in a paper, whose results are presented in one or more **display items**. [ALEKS/FERNANDO]
 - **Causal claim** – An assertion that invokes causal relationships between variables. A paper may estimate the effect of X on Y for population P , using method F . Example: “This paper investigates the impact of bicycle provision on secondary school enrollment among young women in Bihar/India, using a Difference in Difference approach.”
 - **Descriptive/predictive claim** – A paper with such kind of a claim estimates the value of Y (estimated or predicted) for population P under dimensions X using method M . Example: “Drawing on a unique Swiss data set (population P) and exploiting systematic anomalies in countries’ portfolio investment positions (method M), I find that around 8% of the global financial wealth of households is held in tax havens (value of Y).”
- **Coding error** – A coding error will occur when a section of the code, of the reproduction package, executes a procedure that is in direct contradiction with the intended procedure expressed in the documentation (paper or comments of the code). For example an error happens if the paper specify that the analysis is perform on the population of males, but the code restricts the analysis to females only. Please follow the ACRE procedure to report coding errors. [ALEKS/FERNANDO]

- **Data availability statement** – A description, normally included in the paper, of the terms of use for data used in the paper, as well as the procedure to obtain the data (especially important for restricted-access data). Data availability statements expand on and complement data citations. Find guidance on data availability statements for reproducibility [here](#).
- **Data citation** – The practice of citing a dataset, rather than just the paper in which a dataset was used. This helps other researchers find data, and rewards researchers who share data. Find guidance on data citation [here](#).
- **Data sharing** – Making the data used in an analysis widely available to others, ideally through a trusted public repository/archive.
- **Disclosure** – In addition to publicly declaring all potential conflicts of interest, researchers should detail all the ways in which they test a hypothesis, e.g., by including the outcomes of all regression specifications tested. This can be presented in appendices or supplementary material if room is limited in the body of the text.
- **Intermediate data** – Data not directly used as final input for analyses presented in the final paper (including appendices). Intermediate data should not contain direct identifiers.
- **Literate programming** – Writing code to be read and easily understood by a human. This best practice can make a researcher’s code more easily reproducible.
- **Pre-specification** – The act of detailing the method of analysis before actually beginning data analysis.
- **Processed data** – Raw data that have gone through any transformation other than the removal of PII.
- **Raw data** – Unmodified data files obtained by the authors from the sources cited in the paper. Data from which personally identifiable information (PII) has been removed are *still considered raw*. All other modifications to raw data make it *processed*.
- **(Trial) registry** – A database of registered studies or trials, for example the AEA RCT Registry or clinicaltrials.gov. Some of the largest registries only accept randomized trials, hence the frequent discussion of ‘trial registries’. *Registration* is the act of publicly declaring that a hypothesis is being, has been, or will be tested, regardless of publication status. Registrations are time-stamped.
- **Replication** – Conducting an existing research project again. A subtle taxonomy exists and there is disagreement, as explained in Hamermesh, 2007 and Clemens, 2015. *Pure Replication, Reproduction, or Verification* entails re-running existing code, with error-checking, on the original dataset to check if the published results are obtained. *Scientific Repli-*

cation entails attempting to reproduce the published results with a new sample, either with the same code or with slight variations on the original analysis.

- **Reproducibility** – A research paper or a specific display item (an estimate, a table, or a graph) included in a research paper is reproducible if it is possible to reproduce within a reasonable margin of error (generally 10%) using the data, code, and materials made available by the author. Computational reproducibility is assessed through the process of **reproduction**.
- **Reproduction package** – A collection of all the materials associated with the reproduction of a paper. A reproduction package may contain data, code and documentation. When the materials are provided in the original publication they will be labeled as ‘original reproduction package’, when they provided by a previous reproducer they will be referred as ‘reproducer X’s reproduction package’. At this point you are only assessing the existence of one (or more) reproduction packages, you will not be assessing the quality of its content at this stage.
- **Researcher degrees of freedom** – The flexibility a researcher has in data analysis, whether consciously abused or not. This can take a number of forms, including specification searching, covariate adjustment, or selective reporting.
- **Robustness check:** – Any possible change in a computational choice, both in data analysis and data cleaning, and its subsequent effect on the main estimates of interest. In the context of ACRE, the focus should be on the set of **reasonable specifications** (Simonsohn et. al., 2018), defined as (1) sensible tests of the research question, (2) expected to be statistically valid, and (3) not redundant with other specifications in the dataset.
- **Reasonable specification** – [ALEKS/FERNANDO]
- **Specification** – [ALEKS/FERNANDO]
- **Specification searching** – Searching blindly or repeatedly through data to find statistically significant relationships. While not necessarily inherently wrong, if done without a plan or without adjusting for multiple hypothesis testing, test statistics and results no longer hold their traditional meaning, can result in false positives, and thus impede replicability.
- **Trusted digital repository** – An online platform where data can be stored such that it is not easily manipulated, and will be available into the foreseeable future. Storing data here is superior to simply posting on a personal website since it is more easily accessed, less easily altered, and more permanent.
- **Version control** – The act of tracking every change made to a computer file. This is quite useful for empirical researchers who may edit their programming code often.

11.2 Concepts in the ACRE exercise and the platform

- **Analysis code** – A script associated primarily with analysis. Most of its content is dedicated to actions like running regressions, running hypothesis tests, computing standard errors, and imputing missing values.
- **Candidate paper** – A paper that has been considered for reproduction, but the reproducer decided not to move forward with the analysis due to failure to locate a reproduction package. [Learn more here](#).
- **Cleaning code** – A script associated primarily with data cleaning. Most of its content is dedicated to actions like deleting variables or observations, merging data sets, removing outliers, or reshaping the structure of the data (from long to wide, or vice versa).
- **Declared paper** – The paper that the reproducer analyzes throughout the exercise.
- **Display item** – A display item is a figure or table that presents results described in the paper. Each display item contains several **specifications**. [ALEKS/FERNANDO]
- **Reproduction tree/ diagram** – A diagram generated by the ACRE Diagram Builder which represents all the available data and code on behind a specific display item. The tree is meant to represent the entire computational workflow behind a result from the paper. It allows reproducers to trace a display item to its primary sources. It can also be used to guide users of the reproduction package and/or to identify missing components for a complete reproduction.
- **Revised reproduction package** – [ALEKS/FERNANDO]

Bibliography

- Chang, A. and Li, P. (2015). Is economics research replicable? sixty published papers from thirteen journals say 'usually not'. *Available at SSRN 2669564*.
- Christensen, G., Freese, J., and Miguel, E. (2019). *Transparent and reproducible social science research: How to do open science*. University of California Press.
- Galiani, S., Gertler, P., and Romero, M. (2018). How to make replication the norm. *Nature*, 554(7693):417–419.
- King, G. (1995). Replication, replication. *PS: Political Science and Politics*, 28:444–452. See updates to this paper for how I use this paper as a class assignment now.
- King, H., Vilhuber, L., Herbert, S., and Stanchi, F. (2018). The reproducibility of economics research: A case study. Presented at the BITSS Annual Meeting 2018 and available at the Open Science
- of Sciences, N. A. (2019). *Reproducibility and replicability in science*. National Academies Press.