









Working With Clients

Results Reproduction:
Computational reproduction of
research results to ensure
Reproducibility and Transparency.

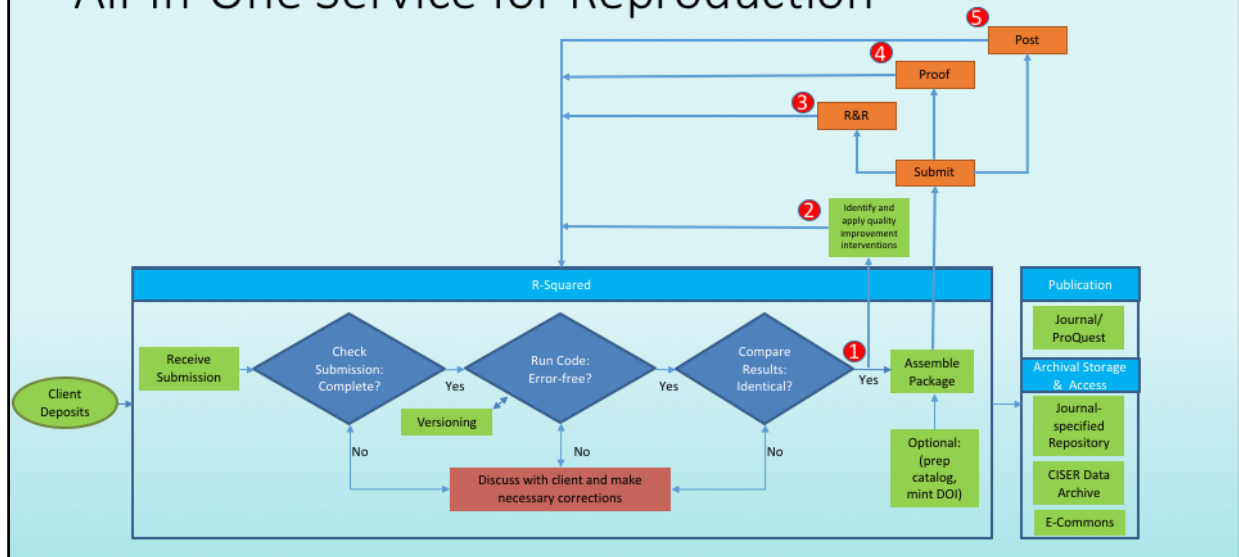


Outline

-  Current Workflow
-  Differences in Projects
-  Most Common Script Errors
-  Most Common Manuscript Errors
-  Good Practices
-  Future Workflow

Here is an outline of what I will speak about today. Starting out I will share the current status of our workflow and discuss some of the elements of it. This will include a brief discussion of the various elements of the workflow. Then I will discuss differences between projects and the inherent difficulty that brings with it. From there I will move to the most common errors that we see, first discussing the most common computer script errors and then the most common manuscript errors. After that, I will discuss good practices to avoid some of these common errors and solutions to them and finally I will discuss our future plans for the workflow.

All-In-One Service for Reproduction



- ❖ This is our current workflow and I will quickly walk you through it. First we receive the submission and check it for completeness. Then we run it and begin to version the code and submission. Once we manage to get it to run all the way through we compare the results with the manuscript. If we find any inconsistencies or differences, we will discuss with the client and make the necessary corrections. During this process we will also look at the code and suggest improvements for either curation or for future usability such as better commenting. Finally we assemble the final package and check that against the manuscript once again to ensure that any changes were made didn't affect the results. As you can see from the diagram, we make multiple checks if our clients want us to, either at Revise and Resubmit, after they have a Proof, and then even post publication to ensure everything lines up.

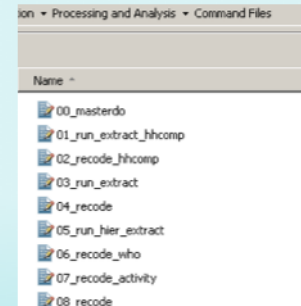
Step 1. Check the submission

Create a project folder following this convention:

R2-<year>-<lastname>-[#]

(eg. R2-2018-KIM-1)

1. Presence of a Readme (often non existant)
2. File names of code contain the order of execution, i.e., prefixed with a step ##.
3. Presence of master program file – often non-existent
4. Identify software used and use base version



First we check the files submitted, with regards to the code:

To do this we create a project folder with the naming convention of “R2-Year-Lastname-Number”

1. We check presence of a readme that would provide instructions on the requirements to reproduce the results. Often there is no readme file. This is most important once the project is finished or in final stages but getting one started early isn't bad.
2. We check the program file names that they indicate the order of execution and/or what it does, so we can easily identify which code to run first or its function.
3. We check for the presence of master program file that executes all the program files in sequence – to simplify the experience of reproducing results
4. We identify the software used by the researchers – once identified we use a freshly installed version of the software, just the base, so we can detect all the packages/modules/libraries/ados that installed by the researchers that were not part of the base package. The code will break and produce an error if the command or library being called by the code is not part of the base.



ARTICLE

- ➔ Highlight all sections (e.g., paragraphs, sentences, tables, charts) that reference output derived from your data.
- ➔ Highlight sections that reference methodological assumptions have been met.

...ion of sampled activities in which the paren.
is alone with children—solo parenting. Mothers
engaged in activities with their children are
significantly more likely to be solo parenting:
49 percent of mothers' activities with children
do not include another adult present compared
to 32 percent of fathers' activities. This is true,
in part, because women are much more likely
than men to be parenting without a spouse or
partner in the household. In some work...

```
acts6|.0393|.0265|.0493|.034|.0291|
acts7|.1149|.1869|.1131|.232|.1168|.
acts8|.1145|.1537|.1081|.1697|.121|.
acts9|.0296|.0346|.0194|.0363|.0397|.
acts10|.0792|.1537|.0748|.1666|.0835|
wykidonly|0|.4299|0|.3249|0|.4874|.
.
} //end T2
```

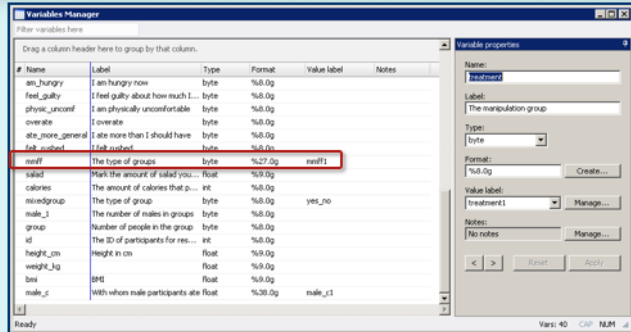
```
acts7|.1149|.1869|.1131|.232|.1168|.
acts8|.1145|.1537|.1081|.1697|.121|.1
acts9|.0296|.0346|.0194|.0363|.0397|.
acts10|.0792|.1537|.0748|.1666|.0835|
wykidonly|0|.4299|0|.3249|0|.4874|.
.
} //end T2
```

To do any further checks (or even to check the code itself) we require something to check it against which is where the article comes in. We ask the researcher to highlight the sections of their article that contain data to be checked, to insure we don't miss any in line numbers. To that end, we also ask that they print out an output file and highlight where the numbers came from. This serves as a sort of check in itself and helps to clarify which numbers we should be looking at, though authors are not great about doing this for us.



DATA

- ➔ All variables and values labeled
- ➔ Free of errors and inconsistencies



When we check the data itself to ensure that it is free of errors and inconsistencies if possible and that variable values are well labeled and named.



Step 2 – Inspect the contents of code

1. Change absolute paths to relative paths – relative to the location of the code
2. Run the code
3. Errors?
 - Correct the error in order to proceed (and not waste time waiting for authors)
 - If unclear how to correct and can't proceed, communicate with researchers





CODE

- ➔ Add comments that map sections of code to results in the manuscript. Make sure commands that generate results are preceded by comments that indicate which result the command generates. For example:

```
24 /* Run the following command to generate the results for Table 1 */
25 use data1.dta
26
27 /* Table 1: Descriptive Statistics of the Sample (N = 100) */
28
29 /* The following command generates the mean age of the sample */
30
31 /* The following command generates the standard deviation of the sample */
32
33 /* The following command generates the 95% confidence interval for the mean age */
34
35 /* The following command generates the 95% confidence interval for the standard deviation of the sample */
36
37 /* The following command generates the 95% confidence interval for the 95th percentile of the sample */
38
39 /* The following command generates the 95% confidence interval for the 5th percentile of the sample */
40
41 /* The following command generates the 95% confidence interval for the 1st percentile of the sample */
42
43 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
44
45 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
46
47 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
48
49 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
50
51 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
52
53 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
54
55 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
56
57 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
58
59 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
60
61 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
62
63 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
64
65 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
66
67 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
68
69 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
70
71 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
72
73 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
74
75 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
76
77 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
78
79 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
80
81 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
82
83 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
84
85 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
86
87 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
88
89 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
90
91 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
92
93 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
94
95 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
96
97 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
98
99 /* The following command generates the 95% confidence interval for the 99th percentile of the sample */
100
```

*The following command generates column 1 of Table 1

*The following command generates the mean age mentioned on page 3, paragraph 3

When we check the code we make sure that it is prefixed with a number to en



```

*****
*****CODE TITLE*****
*****
clear
log using "..\Documents\demoCodeLog.smcl", replace text
import delimited "..\Original Data\HelloWorld.txt"

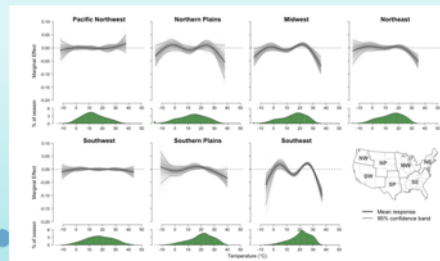
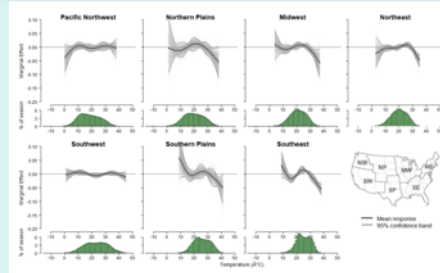
//Labeling the variables and values
label variable treatment "The manipulation group"
label define treatment1 1 "Well Documented" 2 "Poorly Documented"
label value treatment treatment1
label variable pieces "Good labels are descriptive but not too long"
label variable gender "Gender"
label define gender1 1 "Male" 2 "Female"

***** Table 1 - Descriptive statistics of the sample
tab mmff
ttest age if mmff ==1 | mmff == 2, by(mmff) unequal
ttest age if mmff ==3 | mmff == 4, by(mmff) unequal
    
```

- ✓ New variables and values are labeled
- ✓ Codes are commented to describe processes and mapped to paper sections
- ✓ Code produces outputs in the same order as they appear in the manuscript
- ✓ Remove or comment out unnecessary code

Lack of Seed

Any random generation requires a seed or else it can have different results and won't be reproduced. Example on the top is without the seed and bottom is with it.



Lack of a seed is a serious problem in reproduction, and is a tricky thing to handle. Findings should, in theory, be robust enough to be able to withstand using a random seed and get the same results, without the specific seed it is impossible to reproduce the exact findings. The example shown displays how findings can change between two seeds even though there were thousands of generated numbers and the findings were robust.

Insert header/metadata

1. Title of study
2. Author(s)
3. Article DOI
4. Article Citation
5. Primary author contact info and ORCID
6. Date of last code update
7. Code citation
8. License

```

analysis_final
1 version 14
2 log using "../Output/comparisonLog", replace
3
4 .....
5 ** Title: The Role of Education Interventions in Improving Economic Rationality
6 ** Authors: Hyuncheol Bryant Kim, Syngjoo Choi, Booyuel Kim, Cristian Pop-Eleches
7 ** August 16th, 2018
8 ** DOI: https://doi.org/10.1126/science.aar6987
9
10 ** Citation: Hyuncheol Bryant Kim, Syngjoo Choi, Booyuel Kim, Cristian Pop-Eleches.
11 ** "The role of education interventions in improving economic rationality,"
12 ** Science 362(6410), pp. 83-86.
13
14 ** Primary Author: Hyuncheol Bryant Kim
15 ** Primary Author ORCID: https://orcid.org/0000-0001-5304-0274
16 ** Primary Author Contact: hk788@cornell.edu
17
18 ** Software and version: Stata 14 for Windows
19 ** Code last updated: 5 October 2018
20 .....
21
22
    
```

** Code Citation (i.e., how to cite this code):

<author(s) names> (<date>) <title of program/source code>
(<code version>) [<type>]. Web address or publisher.

** License: CC-BY

Usually there is no header in the code.
We put a header.

It contains:

1. Software and version used
2. Command to save the output file. Reusers can compare this output file to the one we include in the package so they know they've executed the code correctly.
3. The title of the study
4. The DOI which we mint
5. Authors, their contact info, and ORCID number so the authors can be contacted should re-users have questions.

The path is anchored on the location of the code, i.e. we make the location of the code as the active working directory, so that reproduction materials package is portable and installable in any folder location without having to recreate or mimic the directory path of the authors.

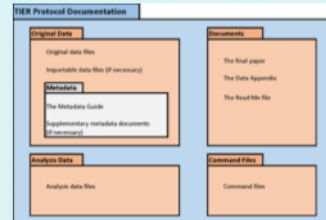
The code should save the output file as these results are the once the researchers will

compare to the manuscript.

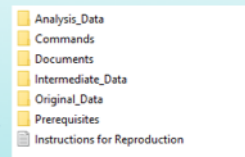
ORCID is important so that in case the author moves to a different institution, researchers will still be able to get in touch with them, assuming they update their ORCID profile.

Step 3. Organize and Package the Reproduction Materials

1. If files are unorganized within the project folder, we use **TIER Protocol** (projectier.org).
2. If files are organized within project folder, retain researcher's existing structure
3. Add a **Prerequisites folder** to store the version of the libraries, packages, ados that researchers used to generate their results



Example of our expanded TIER structure



We use a modified version of the TIER protocol to organize our materials, but really any sort of organization is acceptable, as long as there is a structure and it is organized.

Step 4. Validate

1. Trained CISER Helpdesk Consultants will run the entire package and verify no errors are encountered and results are produced exactly.
2. Author(s) also will be tasked to do the same and confirm.

Trust nothing. Validate everything!

Multiple checks are better, as we are all only human and make mistakes. Multiple checks by multiple people are better still!



Differences in Projects



- Levels of work to be done
- Levels of complexity
- Timeline requirements
- Skill and programming experience of the client
- Level of comfort with reproduction tools

Now onto the difficulties associated with projects. First are the differences between projects which include the level of work that needs to be done, level of complexity, the timeline requirements, the skill and experience of the client as well as their level of comfort with reproduction tools.



Different levels of work need to be done

- The more work needed the harder each part is to do
- Finding all the various scripts that were used and removing unneeded parts
- Creating Master Files
- Re-ordering code to match order of output





Levels of Complexity

- Number of different programming languages
- Length of execute time
 - Can lead to lots of dead time and add complexity
- Amount of outputs to be checked





Timeline requirements

- Still planning out the paper/data not collected
- Pre-Prints
- Not submitted yet
- Revise and Resubmit
- Already Published





Differences in Programming Experience/Skill

- Some disciplines focus more on coding skills, while some don't place as much of an emphasis on it
- Something that may be obvious to someone in one discipline may not be inherently obvious to everyone, most prevalently in unneeded acronyms that can be confusing.





Levels of Comfort With Reproducibility Tools

- Version control such as GitHub
- Dynamic Documentation
- Relative paths in coding





Most Common Error Types present in Scripts



- Not calling the data initially
- Hard Coded Variables
- Missing Files
- Poor or No Commenting
- Uncertain Order of Execution
- Doesn't follow the order of the Output
- Prerequisites
- Extraneous output/lines of code

For the next two sections, I will talk about the most common types of errors we see, and how often I've seen them in the past 10 papers that I have reviewed to give an idea of how common they are.



Not calling the data initially

- Many just open the data with the GUI instead of via code
- Can cause confusion when there are multiple data files present (including raw data, cleaned data, intermediate data, data from different sources, etc.)

5 of last 10 had this problem



Hard Coded Variables

- Generally done either by hand calculator or excel
- Can cause errors when data is further refined, or if mistake was made when calculating it
- Causes confusion as to where the number came from and what it represents

3 of last 10



Missing Files

- Impossible to fully reproduce results with needed files missing
- Often done because researchers will work on multiple computers or utilize poor version control and delete things
- Can also be caused by hard coding a variable and then discarding the code

2 of last 10



Poor or No Commenting and Poor Variable Names

- People will write the code for themselves, and thus commenting can tend to fall by the wayside instead of commenting as they go
- What makes sense to someone who has worked with the code/data for 4 years may not make sense to someone else and need further explanation

3 of last 10, but much more common than that



Uncertain or No Order Specified in Title of Scripts

- Researcher knows the steps so doesn't order them
 - Alternatively they do, but with confusing names such as "first-first..." vs "actually-first..." vs "initialize..."
- Results often generated in the final step of the code, and so if things have been run in the right order previously at that location then the intermediate data will be available so the order doesn't matter on their computer

7 of last 10 missing Master File, but not always uncertain order (some of those no order)



Doesn't Follow the Order of the Manuscript

- Things are rarely coded in the exact order they appear in the manuscript, and can be impossible to do so. However having the results generally follow manuscript dramatically increases transparency
- Often calculated on other computers and then the code thrown together to put everything in the same script file

Hard to say, as code and manuscript evolve over time.



Prerequisites (ado, library, extension, etc.)

- The idea of reproducibility is that it will work on any computer and not having the list of prerequisites will make that not work
- Very easy to forget as they don't need to be re-installed from project to project
- Prone to breaking things when they or the base software gets updated, so having unneeded ones can be just as bad

3 of last 10 were missing a prerequisite library when it was needed/suggestable



Extraneous outputs/lines of code

- Makes the code and output files longer and more confusing unnecessarily
- Not very transparent to have outputs that aren't used and not needed
- Can cause mis-identification or other issues detailed in the common errors in manuscript.

Similar to following the order of the manuscript, as calculations enter or leave the final manuscript, code will become unnecessary or new code needed.



Most Common Error Types present in Manuscripts



- Rounding
- Typos
- Misidentification
- Out of Date Tables
- Missing from Package
- Incorrectly Marking p-values



Rounding

- Extremely common and easy to miss
- Different software/people round in different ways, compounding if it is an error or difference of style
 - However inconsistent rounding is always a problem
- Double rounding can cause errors as well

This is present in virtually all manuscripts, 9 of last 10 had an issue where rounding was most likely to blame.



Typos

- Sometimes obvious, such as missing a decimal place when all the others have it
- Can also just be a mistyped number or two

Can be hard differentiating this from rounding errors or the other type. It was certainly present in at least 2 of last 10.



Misidentification (looking at wrong cell in output file)

- Very easy when there aren't output management software/practices
- Difficult to detect from repeated checks, as the number is present on the output page, but simply is not the result that should be included

This wasn't present in any of the last 10, however we have had issues with this in the past, and it can be hard differentiating rounding/typos from this as well.



Out of Date Tables

- Descriptive Tables are most common culprit
 - Researchers will generate their descriptive tables and then further clean data
 - Can be calculated inconsistently as well, such as valid vs total percentage
- Plenty of other tables can have this happen as well, as researchers often don't re-run their entire code (especially for long execute times)

4 of the last 10 had issues that appear be related to out of date tables



Missing From Package of Scripts

- Can be similar to “hard coded variables” in that the calculation just done by hand/excel and put into the manuscript
- Very common with regards to figures as they are often edited or done in non or different statistical package programs (such as excel) and then edited to look better

3/10 of last 10. This occurs to differing degrees in many papers. Often this is done in excel or just by hand/hand calculator.



Incorrectly Marking P-Values

- Often happens with regards to rounding
 - 0.045 rounds to 0.05 but should be marked as ≤ 0.05 , but can mistakenly be marked as $=0.05$
 - 0.101 shouldn't be marked as ≤ 0.10 , as it is >0.10
- Can come from inconsistency of having multiple people doing the tables
- Easy to make one of the other errors

This happens fairly often when output management systems aren't used which we will talk about in the next section.



Good Practices



- Shared Workspace
- Versioning
- Output Management
- Relative Paths
- Folder Organization
- Scripts



Shared Workspace

- Shared both between the team working on reproduction and the client
- Helps to organize everything and ensures that changes are reflected between both
- Facilitates better back and forth between client and reproducer





Versioning

- Often will go hand in hand with the shared workspace
 - Such as github or even box
- Allows progress to be made on multiple ends, and for mistakes to happen but easily fixed all while retaining a clean work space





Output Management

- Actual output management code/prerequisites
 - Ado/library prerequisite such as Outreg2 in Stata or Huxreg/Huxtables in R
 - Isolating and saving just the values used in the paper into their own table
- Dynamic Documentation
 - Markdown type
 - R Markdown or Stata Markdown and pandoc
 - Notebook type
 - Jupyter Notebook or Mathematica Notebook





Relative Paths

- Allows for the script to be run on other computers without editing them
- Helps to ensure that everything is present, as if file is missing and only relative paths are used the script will fail, whereas if there is an absolute path it will still work on that computer but not any others





Folder Organization

- Facilitates relative paths and leads to a much cleaner workspace
- TIER protocol
- Personal preference of organization with named folders and not everything in same place



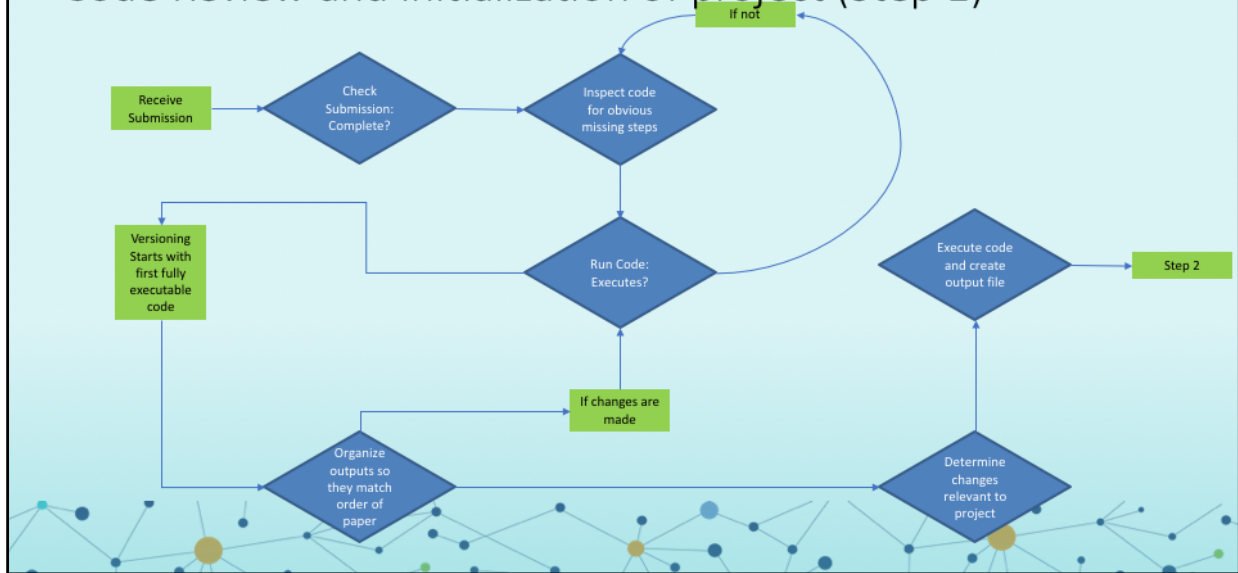


Scripts

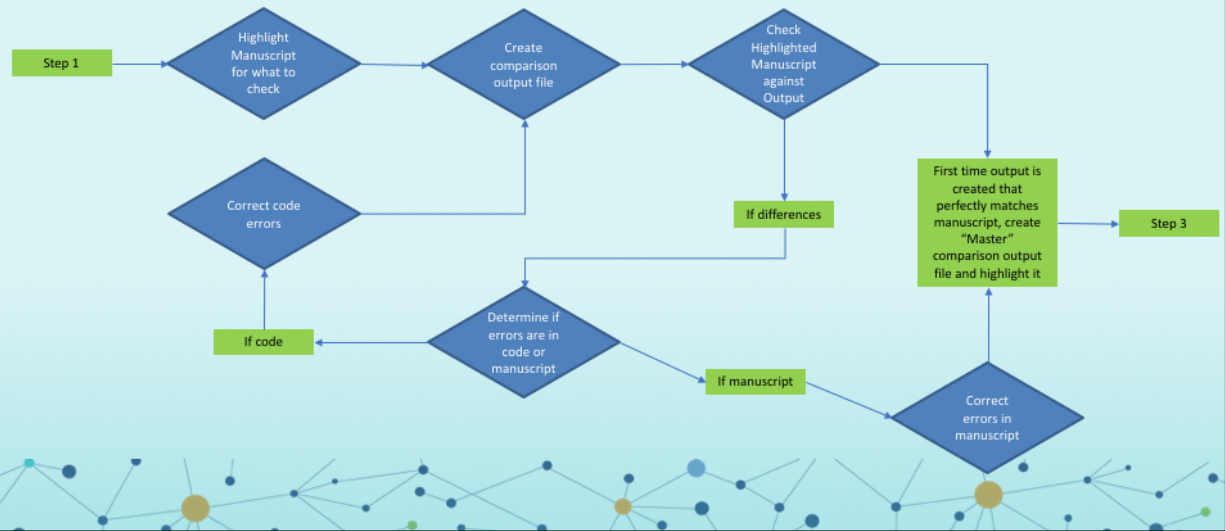
- Well named and numbered to show sequence order
- Master script so only one script needs to be run
- Prerequisites listed out or automatically installed
 - Best to have them locally installed to specific folder to ensure none are missing
 - Prerequisites folder should be part of folder structure
- Code header present
- **Execute beginning to end fresh installation of software**



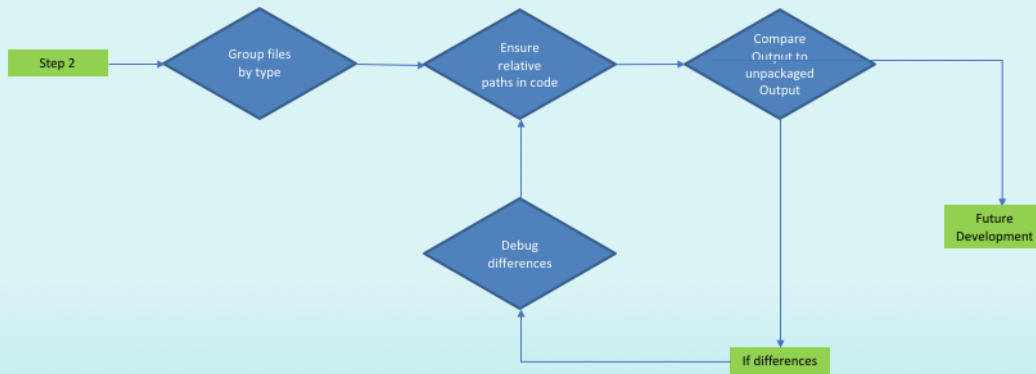
Code Review and initialization of project (Step 1)



Comparison Output (Step 2)



Package Creation (Step 3)





Conclusion

Many of the common errors are related and can be solved in similar ways such as output management.

Each project is fairly unique as clients come with very different starting places and very different needs





THANK YOU

ciser@cornell.edu
Contact information

