

# Reproducibility and Literate Programming (CMP595 PPGC/INF/UFRGS)

Lucas Mello Schnorr, Jean-Marc Vincent

INF/UFRGS

Porto Alegre, Brazil – October 20th, 2017



# Frustration as an Author

- ▶ I thought I used the same parameters but I'm getting different results!
- ▶ The new student wants to compare with the method I proposed last year
- ▶ My advisor asked me whether I took care of setting this or this but I can't remember
- ▶ The damned fourth reviewer asked for a major revision and wants me to change figure 3 :(
- ▶ Which code and which data set did I use to generate this figure?
- ▶ It worked yesterday!
- ▶ 6 months later: why did I do that?

# Frustration as a Reviewer

This may be an interesting contribution but:

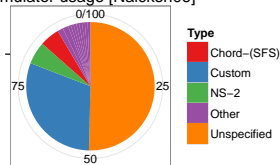
- ▶ This **average value** must hide something
- ▶ As usual, there is no **confidence interval**, I wonder about the variability and whether the difference is **significant** or not
- ▶ That can't be true, I'm sure they **removed some points**
- ▶ Why is this graph in **logscale**? How would it look like otherwise?
- ▶ The authors decided to show only a **subset of the data**. I wonder what the rest looks like
- ▶ There is no label/legend/... What is the **meaning of this graph**? If only I could access the generation script

# A Few Edifying Examples

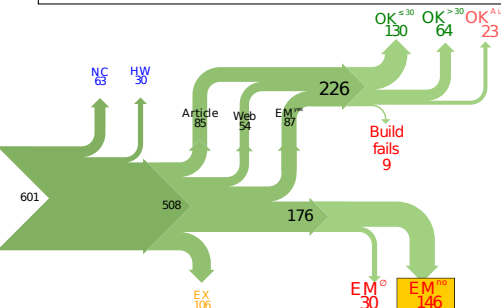
Naicken, Stephen *et Al.*, *Towards Yet Another Peer-to-Peer Simulator*, HET-NETs'06.

From 141 P2P sim.papers, 30% use a custom tool,  
50% don't report used tool

Simulator usage [Naicken06]



Collberg, Christian *et Al.*, *Measuring Reproducibility in Computer Systems Research*, <http://reproducibility.cs.arizona.edu/>



- ▶ 8 ACM conferences (ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12) and 5 journals
- ▶ EM<sup>no</sup> = the code cannot be provided

# The Dog Ate my Homework !!!

## ► Versioning Problems

*Thanks for your interest in the implementation of our paper. The good news is that I was able to find some code. I am just **hoping** that **it** is a stable working version of the code, and **matches the implementation we finally used for the paper**. Unfortunately, I have **lost some data** when **my laptop was stolen** last year. The bad news is that the code is not commented and/or clean.*

*Attached is the  $\langle$ system $\rangle$  source code of our algorithm. I'm **not** very **sure whether it is the final version of the code used in our paper**, but it should be at least 99% close. Hope it will help.*

# The Dog Ate my Homework !!!

- ▶ Versioning Problems
- ▶ Bad Backup Practices

*Unfortunately, the server in which my implementation was stored had a **disk crash in April and three disks crashed simultaneously**. While the help desk made significant effort to save the data, my entire implementation for this paper was not found.*

# The Dog Ate my Homework !!!

- ▶ Versioning Problems
- ▶ Bad Backup Practices
- ▶ Code Will be Available Soon

*Unfortunately the current system is **not mature enough at the moment**, so it's not yet publicly available. We are actively working on a number of extensions and **things are somewhat volatile**. However, once things stabilize we plan to release it to outside users. At that point, we would be happy to send you a copy.*

# The Dog Ate my Homework !!!

- ▶ Versioning Problems
- ▶ Bad Backup Practices
- ▶ Code Will be Available Soon
- ▶ No Intention to Release

*I am afraid that the source code was never released. The code was never intended to be released so is not in any shape for general use.*



# The Dog Ate my Homework !!!

- ▶ Versioning Problems
- ▶ Bad Backup Practices
- ▶ Code Will be Available Soon
- ▶ No Intention to Release
- ▶ Programmer Left

*⟨STUDENT⟩ was a graduate student in our program but **he left a while back** so I am responding instead. For the paper we used a prototype that included many moving pieces that only ⟨STUDENT⟩ knew how to operate and we did not have the time to integrate them in a ready-to-share implementation before he left. Still, I hope you can build on the ideas/technique of the paper.*

*Unfortunately, the author who has done most of the coding for this paper has **passed away** and the code is no longer maintained.*

# The Dog Ate my Homework !!!

- ▶ Versioning Problems
- ▶ Bad Backup Practices
- ▶ Code Will be Available Soon
- ▶ No Intention to Release
- ▶ Programmer Left
- ▶ Commercial Code

*Since this work has been done at  $\langle$ COMPANY $\rangle$  we don't open-source code unless there is a compelling business reason to do so. So unfortunately I don't think we'll be able to share it with you.*

*The code owned by  $\langle$ COMPANY $\rangle$ , and AFAIK the code is not open-source. Your best bet is to reimplement :( Sorry.*

# The Dog Ate my Homework !!!

- ▶ Versioning Problems
- ▶ Bad Backup Practices
- ▶ Code Will be Available Soon
- ▶ No Intention to Release
- ▶ Programmer Left
- ▶ Commercial Code
- ▶ Proprietary Academic Code

*Unfortunately, the  $\langle \text{SYSTEM} \rangle$  sources are **not meant to be opensource** (the code is partially **property of  $\langle \text{UNIVERSITY 1} \rangle$ ,  $\langle \text{UNIVERSITY 2} \rangle$  and  $\langle \text{UNIVERSITY 3} \rangle$ .**)*

*If this will change I will let you know, albeit I do not think there is an intention to make the  $\langle \text{SYSTEM} \rangle$  sources opensource in the near future.*

*If you're interested in obtaining the code, **we only ask for a description of the research project** that the code will be used in (**which may lead to some joint research**), and we also have a software license agreement that the University would need to sign.*

# The Dog Ate my Homework !!!

- ▶ Versioning Problems
- ▶ Bad Backup Practices
- ▶ Code Will be Available Soon
- ▶ No Intention to Release
- ▶ Programmer Left
- ▶ Commercial Code
- ▶ Proprietary Academic Code
- ▶ Research vs. Sharing
- ▶ ...
- ▶ ...

*In the past when we attempted to share it, we found ourselves spending more time getting outsiders up to speed than on our own research. So I finally had to establish the policy that we will not provide the source code outside the group.*

# Structure

Research articles are often structured in this basic order:

**Introduction** Why was the study undertaken? What was the research question, the tested hypothesis or the purpose of the research?

**Methods** When, where, and how was the study done? What materials/hardware were used? How was it configured?

**Results** What answer was found to the research question; what did the study find? Was the tested hypothesis true? **Present useful results in a synthetic way with a logical order.**

**Discussion** What might the answer imply and why does it matter? How does it fit in with what other researchers have found? What are the possible bias and points to improve? What are the perspectives for future research?

Such structure **facilitates literature review** and is a very effective way to convey information.

If the report is a few pages long then **an abstract is required.**

# Step 0: Taking Notes

Document your:

- ▶ **Hypotheses**: keep track of your ideas/line of thoughts
- ▶ **Experiments**: details on how and why an experiment was run, including failed or ambiguous attempts.
- ▶ **Initial analysis or interpretation** of these experiments: was the outcome conform to the expectation or not? does it (in)validate the hypothesis?
- ▶ **Organization**: keep track of things to do/fix/test/improve

Structure:

1. General information about the document and organization **conventions** (e.g., directory structure, notebook structure, experimental result storing mechanism, ...)
2. Documentation of **commonly used commands** and of how to set up experiments (e.g., git cloning, environment deployment, connection to machines, compiling scripts)
3. Experiment results can be either structured **by dates** (↪ add tags) or **by experiment campaigns** (↪ add date/time)

# Which format should I use ?

- ▶ Wikis are encouraged to favor collaboration but I do not find them really effective
- ▶ Blogging systems are also a way of managing such notebook but they should rather be considered as an effective way to share information with others
- ▶ I recommend to use basic plain-text format and to structure it hierarchically

Here is a link to an excerpt of the journal of one of my PhD student, managed with git/org-mode. More detailed links are given in slide ??.

Last but not least:

Provide links to Raw Data!!!

# When/How Often Should I Use it?

I have a very intense usage (demo to **general journal** and specific **BOINC journal**) and I tend to capture a lot of information but you do not have to be as extreme as I am. Here are a few advices:

- ▶ Spending **more than an hour without** at least **writing** what you're working on **is not right**. . .
  - ▶ **Take a 5 minutes** break and ask yourself what you're doing, what is keeping you busy and where all this is leading you
- ▶ While working on something, you will often notice/think about something you should fix/improve but you just don't want to do it now. Take 20 seconds to write a **TODO** entry.
- ▶ There are moments where you have to **wait for something** (compiling, deployment, . . . ). It is generally the perfect time for improving your notes (e.g., detail the steps to accomplish a TODO entry).
- ▶ **By the end of the day**: daily (and weekly) **review!**
  - ▶ Update your lists, write what the next steps are
  - ▶ **Summarize in a 2-4 lines** (for your advisor) what you did, what was difficult, what you learnt.



# Step 1: Sharing Code and Data

What kinds of systems are available?

- ▶ "Good- The cloud (Dropbox, Google Drive, Figshare)
- ▶ Better - Version control systems (SVN, Git and Mercurial)
- ▶ "Best- Version control systems on the cloud (GitHub, Bitbucket)

Depends on the level of privacy you expect but you probably already know these tools.

**Few handle GB files...**

Is this enough?

1. Use a workflow that documents both data and process
2. Use the machine readable CSV format
3. Provide raw data and meta data, not just statistical outputs
4. Never do data manipulation and statistical tests by hand
5. Use R, Python or another free software to read and process raw data (ideally to produce complete reports with code, results and prose)

Courtesy of Adam J. Richards

## Step 2: Literate Programming

**Donald Knuth**: explanation of the program logic in a **natural language interspersed with snippets of** macros and traditional **source code**.

I'm way too stupid to program this way but that's  
**exactly what we need for writing a reproducible article/analysis!**

### Org-mode (requires emacs)

My favorite tool.

- ▶ plain text, very smooth, works both for html, pdf, ...
- ▶ allows to combine all my favorite languages even with sessions

### Ipython notebook

If you are a python user, go for it! Web app, easy to use/setup...

### KnitR (a.k.a. Sweave)

For non-emacs users and as a first step toward *reproducible papers*:

- ▶ Click and play with a modern IDE (e.g., Rstudio)