

# Licensing in the Computational Sciences

Victoria Stodden

Cyberscholars Research Group

Yale Law School

April 22, 2008

# The “Third Branch” of the Scientific Method

- Scientific study is no longer only *deductive* or *empirical*, it is also *computational*:
  - Simulations
  - Massive data research
- Across many fields:
  - Genetics, Biostatistics
  - Machine Learning / Network Analysis
  - Computer Science, Statistics
  - Geophysics/Earth Sciences
  - Law ...

# Computational Research and the Ubiquity of Error

- Standards are established in deductive and empirical fields (proofs, error detection)
- Need **standards** for computational research -> reproducibility
- Publicly funded work is not made publicly available, as often mandated
- Scientific research remains closed and difficult to access; fields fragment

# Solution: Standards for Computational Research

- Need clear explanations of the process to find and eliminate errors
- Need standardized mechanisms for verifying work and validating conclusions

# Definition: Research *Compendium*

*(Gentleman & Lang 2004)*

The entire set of materials required to reproduce the results:

- Research paper and source files
- Data and its documentation, methodology, code (meta-data)
- Code, parameters, instructions from the experiment
- Results and documentation
- Auxiliary materials

# Definition: Research *Compendium*

*(Gentleman & Lang 2004)*

The entire set of materials required to reproduce the results:

- **Research paper** and source files
- Data and its documentation, methodology, code (meta-data)
- Code, parameters, instructions from the experiment
- Results and documentation
- Auxiliary materials

# The Solution: Reproducible Research

- Align incentives to promote reproducible research:
  - Allay attribution concerns
  - Educate about copyright and compendium release
  - Peer review of open research, verification of results

# The Computational Sciences Public License

- Ensures attribution
  - Viral attribution attached to all derivative works that use the original research product, limited to the researcher's contribution in derivative products (not Share Alike)
- Ensures entire compendium released
  - includes meta-data but excepting the data itself (“selection and arrangement”)



# Why Reproducible Research?

- Computational research is becoming more pervasive
- Reproducible papers cited more frequently
- Publicity creates an incentive for better quality work (“sunshine principle”)
- Accountability and oversight
- Knowledge extends outside the immediate field to other fields, regular citizens

# Why? The Scientist

- Incentive for the scientist to release work with attribution
- Easier to use CSPL than a GPL/CC hybrid
- Publicity of the licensed release concept

# Why a New License?

- GPL: focus is code, other elements of the research compendium eschewed
- CC: focus is media
- Attention to reproducible research and computational standards
- Promotion of knowledge, science
- Vehicle to tie funding to reproducibility