# Sharing Data with Dataverse

Mercè Crosas, Ph.D.
Director of Product Development
Institute for Quantitative Social Science (IQSS)
Harvard University

**IQSS**

The Institute for Quantitative Social Science
HARVARD UNIVERSITY

# 1. Data Sharing and Replication

# 2. A Solution with the Dataverse Network

SHARE your data

it's good for you, and for the world.

Come. Eat lunch. Accelerate the pace of science.

CfA, PHILLIPS Auditorium, 11:45 MONDAY 4/2/12

theastrodata.org

POWERED BY THE **Dataverse Network™** PROJECT & SEAMLESS ASTRONOMY Linking scientific data, publications, and communities

Slide acknowledgment: Alyssa Goodman

# From Data Sharing to Replication

"The replication standard holds that **sufficient information** exists with which to understand, evaluate, and build upon a prior work if a third party can replicate the results without any additional information from the author."
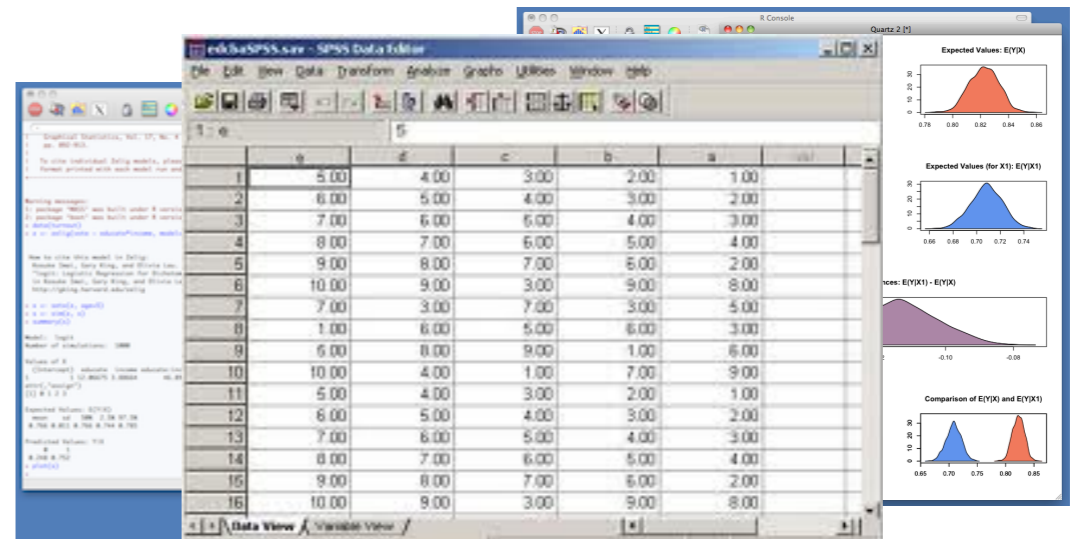
King, G. 1995 "Replication, Replication"
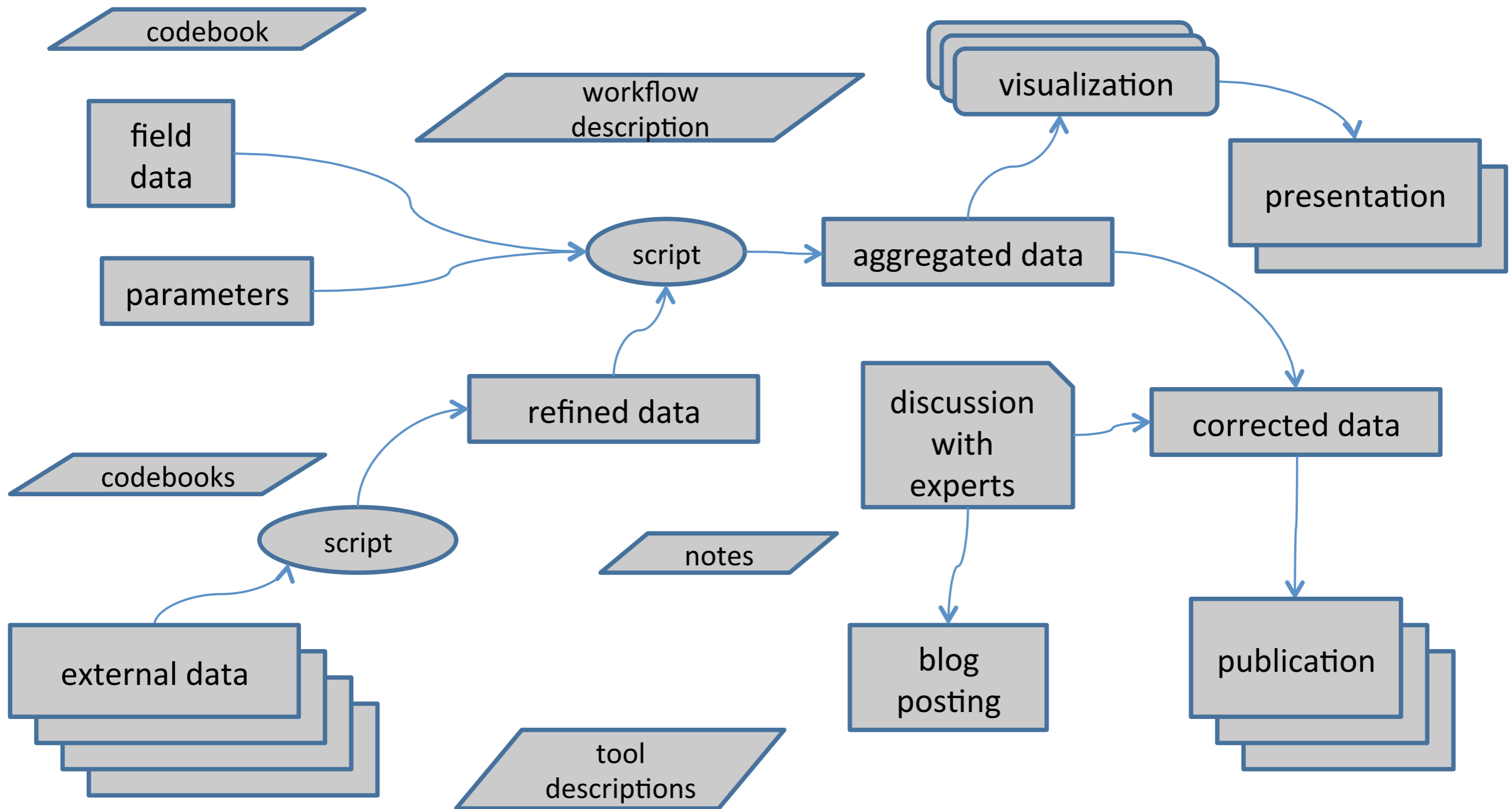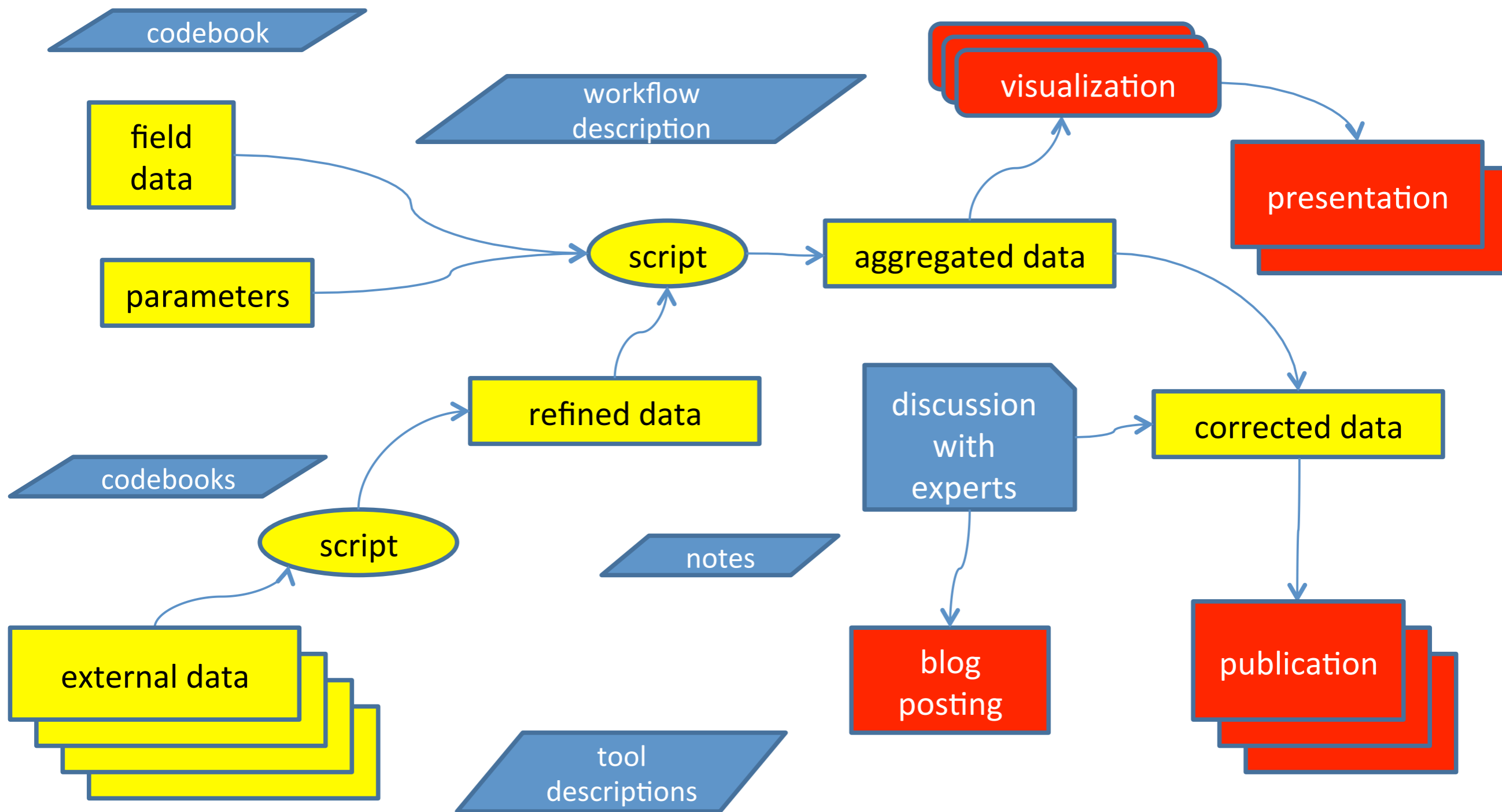
Published Work

Data + Metadata + Supporting Files



+

# "Sufficient Information"?



codebook

field data

parameters

workflow description

script

aggregated data

visualization

presentation

refined data

discussion with experts

corrected data

codebooks

script

notes

blog posting

publication

external data

tool descriptions

# "Sufficient Information"?



codebook

field data

workflow description

parameters

script → aggregated data

visualization

presentation

refined data

script

codebooks

external data

notes

discussion with experts → corrected data

blog posting

publication

tool descriptions

= data    = context (documentation)    = research products

Slide acknowledgment: Andrea Goethals

# Likelihood to Replicate

All the information about environment, data, models, agents, …

Data/ Information Shared

Very Possible

Moderately Possible

None or little information shared (need to contact author?)

Not Likely

Highly Unlikely

Now

50 yr (or ∞)

Time

At the very least, need sufficient documentation to judge the veracity and usefulness of the data

# What can happen if we don't share

**The New York Times**

**Research**

**Search Health** 3,000+ Topics

[ Go ]

In a survey of more than 2,000 American psychologists scheduled to be published this year, Leslie John of Harvard Business School and two colleagues found that 70 percent had acknowledged, anonymously, to cutting some corners in reporting data.

## Fraud Case Seen as a Red Flag for Psychology Research

By BENEDICT CAREY
Published: November 2, 2011

A well-known psychologist in the Netherlands whose work has been published widely in professional journals falsified data and made up entire experiments, an investigating committee has found. Experts say the case exposes deep flaws in the way science is done in a field, psychology, that has only recently earned a fragile respectability.

… an analysis of 49 studies appearing Wednesday in the journal PLoS One, by Dr. Wicherts, Dr. Bakker and Dylan Molenaar, found that the more reluctant that scientists were to share their data, the more likely that evidence contradicted their reported findings.

[+] SHARE

"We have the technology to share data and publish our initial hypotheses, and now's the time,"

The psychologist, Diederik Stapel, of Tilburg University, committed academic fraud in "several dozen" published papers, many accepted in respected journals and reported in the news media, according to a report released on Monday by the three Dutch institutions where he has worked: the University of Groningen, the University of Amsterdam, and Tilburg. The journal Science, which published one of Dr. Stapel's papers in April, posted an "editorial expression of concern" about the research online on Tuesday.

Joris Buijs/Pve

The psychologist Diederik Stapel in an undated photograph. "I have failed as a scientist and researcher," he said in a statement after a committee found problems in dozens of his papers.

# What can happen if we share

## Sharing Detailed Research Data Is Associated with Increased Citation Rate

| Article | Metrics | Related Content | Comments: 5 |

**Heather A. Piwowar[*], Roger S. Day, Douglas B. Fridsma**

Department of Biomedical Informatics, University of Pittsburgh School of Medicine, Pittsburgh, Pennsylvania, United States of America

## Abstract  Top

### Background

Sharing research data provides benefit to the general scientific community, but the benefit is less obvious for the investigator who makes his or her data available.

### Principal Findings

We examined the citation history of 85 cancer microarray clinical trial publications with respect to the availability of their data. The 48% of trials with publicly available microarray data received 85% of the aggregate data was significantly (p = 0.006) associated with a 69% increase in cita impact factor, date of publication, and author country of origin using linear regression.

### Significance

This correlation between publicly availab investigators to share their detailed res

"We found that cancer clinical trials which share their microarray data were cited about 70% more frequently than clinical trials which do not."

1. Data Sharing and Replication

2. A Solution with the Dataverse Network

**open source**

The Dataverse Network is a repository for research data that takes care of long term preservation and good archival practices, while researchers keep control of and get recognition for their data.

# Dataverse Network for Social Science Data

# Dataverse Network for Astronomy Data

# In NeuroInformatics, sharing data for long term access has similar challenges:

✓ **Incentives to share**

✓ Unit of data citation

✓ Sufficient metadata

✓ Obsolescence of formats and software

✓ Discoverability

✓ Increasing size of data

✓ ...

# Recognition and Credit for Author



Your own site or project site

Your Dataverse

# Formal Data Citation



Weisiger, Alex, "Replication data for: Logics of War: Explanations for Limited and Unlimited Conflicts", http://hdl.handle.net/1902.1/18738 UNF:5:OJCPMDOPJ96QO9V7fhXJMA== V1 [Version]

**Persistent URL (Handle)**
**for permanent reference**

**Universal Numerical Fingerprint**
**(UNF) for verification**

Altman, M., King, G., 2007 "A Proposed Standard for the Scholarly Citation of Quantitative Data"

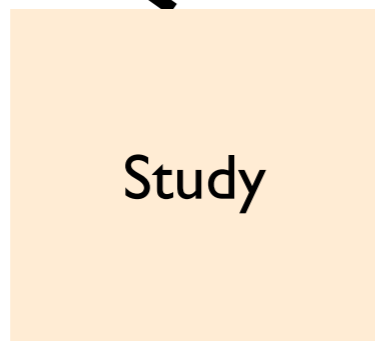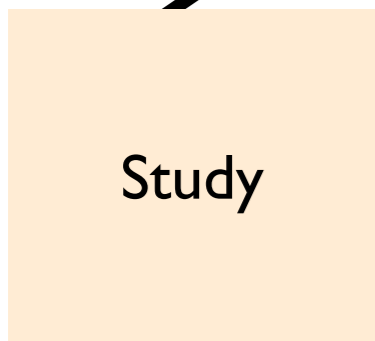**DATAVERSE NETWORK** — A **centralized** software installation and data repository

**Dataverse** — Individual virtual data archive with its own branding

**Collections** / **Study** — A study describes and holds the data (self-curated)

**Metadata** / **Data Files** — Metadata is searchable

# Data Versioning and Data Management

## Contributors, curators, admins view

## End user view

New Study → In Review → **Released version 1**

Metadata
title:
author:
abstract:
year:
methods:
...

Data File 1

Edit Study + Add New File

Draft → In Review → **Released version 2**

Data File 2

Christopher Casillas; Peter K. Enns; Patrick C. Wohlfarth, "How Public Opinion Constrains the Supreme Court", http://hdl.handle.net/1902.1/14568 UNF:5:YBYuXzOp6STpakyTEIoScQ== V1 [Version]

Christopher Casillas; Peter K. Enns; Patrick C. Wohlfarth, "How Public Opinion Constrains the Supreme Court", http://hdl.handle.net/1902.1/14568 UNF:5:dI8qi49P0uIB9pLfXA3RCw== V2 [Version]

# Connecting Publications to Data

# The Centralized Data Repository Provides:

✓ Backups and replication of data in different locations (LOCKSS)

✓ Re-formatting for preservation (e.g. from SPSS, STATA to archival format)

✓ Extraction of Metadata from data sets

✓ Metadata standards (DDI, Dublin Core)

✓ Inter-operability (OAI-PMH, APIs)

It handles good archival practices for you

# Open Source, Java EE 6 Architecture

**Glassfish**

JSF

OAICat Harverster

EJB 3.0 (Stateful Beans, JMS, etc.)

Lucene

**RServe**

Analysis Component (Zelig/R)

**PostgreSQL**

**File System**

Data Files

- JSF, Javascript (jquery): user interface

- EJB: middle-tier, business logic

- OAI-PMH client and server: Harvest Metadata

- Lucene: Indexes metadata

- PostgresSQL: Persistence storage of metadata

- Files system: Store data and complementary files.

- RServe: Analysis component for quantitative data files

# concerns and suggestions from the neuroinformatics community?

The Dataverse Network Project: http://thedata.org

mcrosas@iq.harvard.edu
IQSS, Harvard University