# Embracing research data
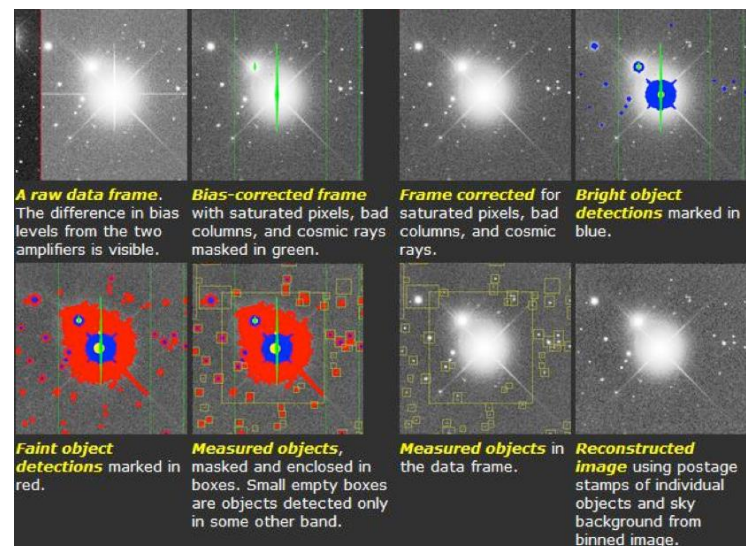## Lucie Boudova, Elsevier

- **6 December 2016**

## Embracing research data

- **Impact of research data sharing**

- **Fundamentals of research data**

- **Components of effective research data**

- **Tools and programmes supporting research data**
  - Linking-data programme
  - Industry standards
  - Data search
  - Research protocols (HiveBench)
  - Data repository (Mendeley Data)
  - Data journals

- **Research Data Policy**

# What are we really after: astronomy

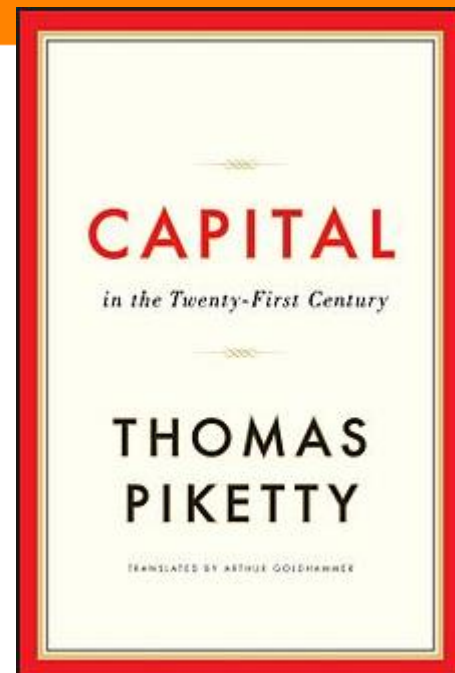**Extracts from "the top 10 benefits of data sharing in astronomy", from Sloan Digital Sky Survey:**

- **Early data releases greatly improve the final product, e.g. more people "looking" at the data increases the chance of finding subtle problems, especially important for space missions with finite lifetime, e.g. the ESA's Gaia mission**

- **More science is extracted from the same dataset, e.g. diversity of ideas: many of the most visible SDSS results were unanticipated in the original project proposal**

- **Sometimes the only way to secure scarce resources,** "easy things" (e.g. those that can be put together by a small number of groups/institutions) have been done in the last century; the "road ahead" requires more substantial merging of research resources, like HST Deep Field, UKIDSS, LSST

- **Results in more citations and prestige to the team who produced data;** practically all postdocs from the first phase of SDSS hold faculty-level positions today



*A raw data frame.* The difference in bias levels from the two amplifiers is visible.

*Bias-corrected frame* with saturated pixels, bad columns, and cosmic rays masked in green.

*Frame corrected* for saturated pixels, bad columns, and cosmic rays.

*Bright object detections* marked in blue.

*Faint object detections* marked in red.

*Measured objects,* masked and enclosed in boxes. Small empty boxes are objects detected only in some other band.

*Measured objects* in the data frame.

*Reconstructed image* using postage stamps of individual objects and sky background from binned image.

# What are we really after: social sciences

**Capital in the Twenty-First Century is a 2013 book by French economist Thomas Piketty.**

- **It focuses on wealth and income inequality in Europe and the United States since the 18th century**

- **Central thesis is that when the rate of return on capital (r) is greater than the rate of economic growth (g) over the long term, the result is concentration of wealth, and this unequal distribution of wealth causes social and economic instability**

- **All raw data, normalized data, the analysis, and methods have all been made publicly available on a dedicated website**

*"Here are enormous quantities of information distilled from tax rolls, inheritance records, and various other public data sources, laid out in charts that should be readily accessible to the layest of lay readers. Not all of the information in these sections is novel or startling. Having it together in one place, however, is valuable, and even most of the book's fiercest critics respect this achievement."* [1]

It also shows data sharing can lead to issues [2]:
- Chris Giles, economics editor of the Financial Times (FT), identified what he claims are "unexplained errors" in Piketty's data, in particular regarding wealth inequality increases since the 1970s. "contain a series of errors that skew his findings"
- Subsequently, Piketty wrote a response defending his findings; the accusation and responses received wide press coverage
- E.g. Scott Winship, a sociologist at the MIPR, claims the allegations are not "significant for the fundamental question of whether Piketty's thesis is right or not"

4

# When we talk about data, we really talk about the following:



Machine & environment settings



Scripts & analyses



Raw data



Processed data



Protocols, methods, algorithms

5

# When You Leave Your Institution, What Happens To Your Data?



| Category | Value | Label |
|---|---|---|
| Datenverbleib | 58 | Stays at institution |
| Datenmitnahme | 49 | Take it with me |
| Weiß nicht | 20 | Don't know |
| Datenlöschung | 7 | Data is lost |
| Sonstiges | 2 | Other |

e-infrastructures austria

„Forschende und ihre Daten. Ergebnisse einer österreichweiten Befragung (eBook)"
E-infrastructures Austria
Bauer, B. (Bruno) et all
Oct 2015
https://phaidra.univie.ac.at/detail_object/o:407736

# The 10 components for effective research data



10. Integrate upstream and downstream – make metadata to serve use.

9. Re-usable (allow tools to run on it)

8. Reproducible

7. Trusted (e.g. reviewed)

6. Comprehensible (description / method is available)

Use

5. Citable

4. Discoverable (data is indexed or data is linked from article)

3. Accessible

Share

2. Preserved (long-term & format-independent)

1. Stored (existing in some form)

Save

# Tools and programmes supporting research data

# Data-linking programme

- Elsevier has an extensive programme with 60+ leading domain-specific data repositories to interlink articles and data
- Makes it easier to find relevant data and place data into the right context
- Linking through in-article accession numbers, data DOI's, or data banners



Linking through **in-article data accession numbers**



**Database banners** shown next to the article on ScienceDirect

See http://www.elsevier.com/databaselinking

# Data-linking programme – example Pangaea



- Supplementary data at PANGAEA
- Bidirectional links between PANGAEA & ScienceDirect
- Data visualized next to the article

## Research Data Working Groups and Development of Industry Standards example www.Scholix.org

**SCHOLIX**

- ICSU/WDS/RDA Publishing Data Service Working group

- Creating linked-data model for exposing DOI to DOI links outside publisher's firewall

- Collaboration between CrossRef, DataCite, Europe PubMed Central, ANDS, Thompson Reuters, Elsevier, OpenAire

*Objective: move from*

a plethora of (mostly) bilateral arrangements between the different players…

*.. to ..*

.. **a one-for-all cross-referencing service** for articles and data

ICSU WORLD DATA SYSTEM

RDA RESEARCH DATA ALLIANCE

# <u>Data</u>search engine!

- Many (broad) datasearch examples already available

| BASE | BioCaddie/ DataMED | Datacite | Datahub.io | DataONE | EbiSearch | OneRepo |
|---|---|---|---|---|---|---|

| Quandl | RE3Data.org | Semantic Scholar | OSF\|SHARE | TR Data Citation Index | Zanran |
|---|---|---|---|---|---|

- Some common themes:
  - search of metadata only (i.e. ranking based on metadata only)
  - And/or federated search (i.e. no ranking)
  - And/or focused on giving credit (citation) rather than on discoverability

- Uncommon (because difficult):
  - Deep indexing of datasets (so real ranking and filtering)
  - Search engine really focused on data discovery

**Elsevier Data Search**
**E.g. search for "Temperature viscosity ionic liquids"**

# Research Protocols – capturing and sharing



www.hivebench.com

# Manage, store: Mendeley Data

**An open repository for posting & reusing research data**

# Manage, Store: Mendeley Data



Linked to published papers – or not

Linked to Github – or not

Versioning and provenance

# Data journals: SoftwareX

Home > Books & Journals > SoftwareX

## SoftwareX

Editors-in-Chief: Dr. Kate Keahey, Dr. Frank Seinstra, Dr. David Wallom
View full editorial board

Open Access

ISSN: 2352-7110

### Code metadata

| | |
|---|---|
| Current code version | v0.6 |
| Permanent link to code/repository used of this code version | https://github.com/ElsevierSoftwareX/SOFTX-D-15-00005 |
| Legal Code License | NCSA open source license |
| Code versioning system used | git |
| Software code languages, tools, and services used | C, C++, Python, Bash; MPI, OpenMP, CUDA |
| Compilation requirements, operating environments & dependencies | Compilers: GNU/Intel/Cray; OS: Linux (RedHat, Debian, Ubuntu, CentOS, SUSE); Dependencies: GDAL, GEOS, PROJ4, SPRNG, PySAL, OpenGeoDa, etc. |
| If available Link to developer documentation/manual | https://github.com/cybergis/cybergis-toolkit |
| | http://cybergis.cigi.uiuc.edu/cyberGISwiki/doku.php/ct |
| Support email for questions | CyberGIS Helpdesk (help@cybergis.org) |

Table options ▾

**AWARD FOR INNOVATION IN JOURNAL PUBLISHING**

🏆

Elsevier
*SoftwareX*
Editors Dr. Kate Keahey, Dr. Frank Seinstra and Professor David Wallom

AMERICAN PUBLISHERS AWARDS — FOR PROFESSIONAL AND SCHOLARLY EXCELLENCE — 40 YEARS

Commits ⓘ          14

Powered by GitHub and Scopus

# The 10 components for effective research data

# Elsevier initiatives

10. Integrate upstream and downstream – make metadata to serve use.

9. Re-usable

8. Reproducible

7. Trusted

6. Comprehensible

5. Citable

4. Discoverable

3. Accessible

2. Preserved

1. Stored

Research Protocols (Hivebench)

Mendeley data repository

Data journals

Data Linking

Data Search

# Efficiency – integration is the building stone

## Research Data Policy

Elsevier will:

- Encourage and support researchers and research institutions to **share data where appropriate and at the earliest opportunity**.
- Provide **guidance to authors regarding the deposit and sharing** of data.
- **Encourage and enable two-way linking** of relevant datasets and publications **using permanent standard identifiers.**
- Encourage and **support proper data citation practices** so that researchers can be cited and credited for their work.
- Work closely with the scientific community to **establish data review practices** to ensure that published research data is valid, properly documented and can be re-used.
- **Develop tools and services** to support researchers to **discover, use and reuse** data to further their research.

> *"Raw research data should be made freely available to all researchers wherever possible"* – STM Brussels Declaration 2007

# Thank you ! Questions?

## Contact me at L.Boudova@Elsevier.com

ELSEVIER

Empowering Knowledge