

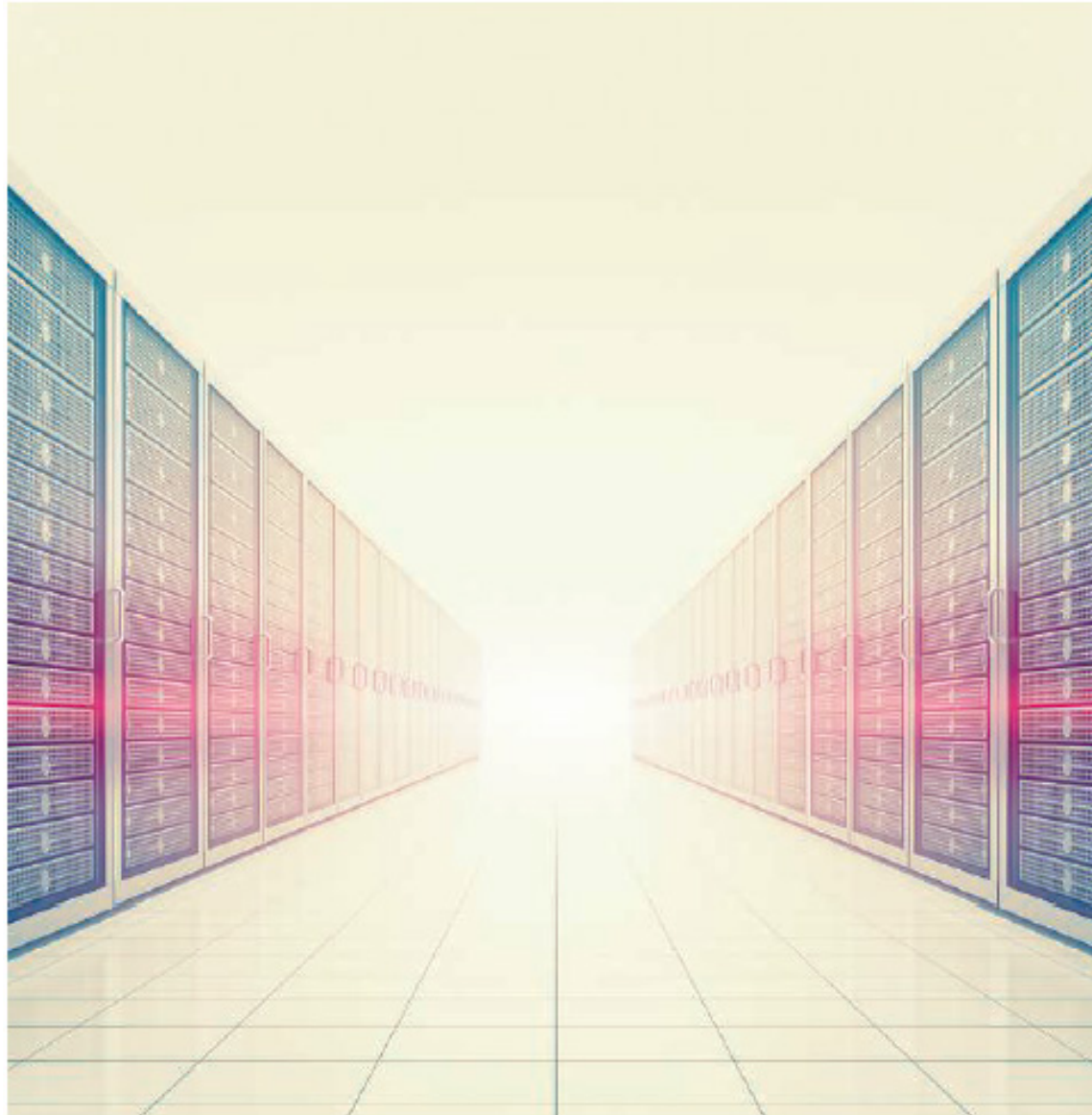
Visualizing data integrity and management through the lens of transparency

© 2016 fpm/iStockphoto & AAAS



© 2016 Sébastien Thibault

C. S. Raman
UMB School of Pharmacy
15 September 2016

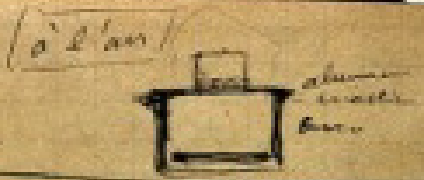


***“...interpretations come and go,
but data are forever.”***

*Marcia McNutt
Editor-in-Chief
Science Journals*

Marie Curie's radioactive notebook (1899 - 1902)

19 janvier
1 jour
rien



1000 1750 chauffe
cable étalonné
Cable



id 1000 chauffe
Cable



après 2 jours
après 3^e
tube chauffe 1000 - 10^e
faucilles - chauffe 2000 - 20

faucilles verre tube chauffe 100 - 10^e
chauffe 1000 - rien

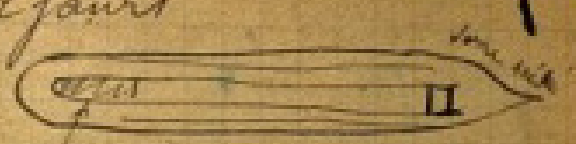
21 janvier - (boîte métallique)
rien
2 jours



papier AB - rien
papier dilaté 2000 - 4^e app e
(boîte étalée app e)
ce AB est rien, cabl' A D rien
(desm' cuivre CD - 2000 - 10^e à 5^e app e)

boîte ouverte AB devant
intérieur ou extérieur rien app e
mastic actif collé sur plaque
ou la rend très active
(c'est la cause de l'erreur de dernière page)

définis 2 jours
21 jour



A 1000 1750 chauffe
lanchan parafuse 3^e
trappe fait app e
amy fort app e

C ₁₆	500 - 40	(1) 2000 - 10
M ₁	509 - 17	(5) 1000 - 1
C ₁₂ - fac	10 - 8	

après 2 jours
traverse tube aluminium
B



100	13
100	14
100	18

Rien

Alexander Fleming's data recordings: 1928

Oct 30. 28

12

Staph. inhibiting mould.

Cult

Sarcina

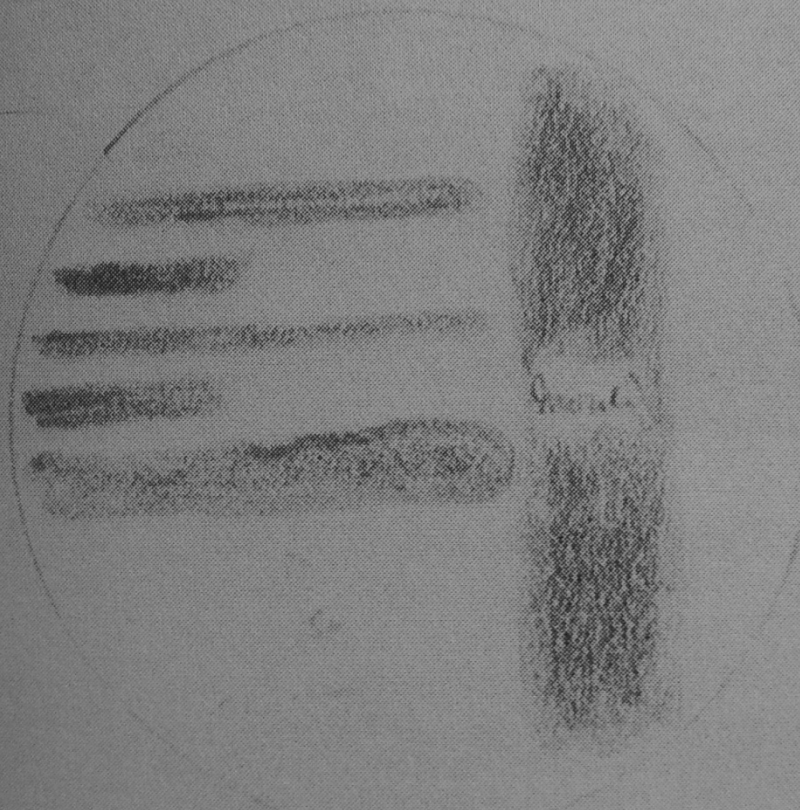
Cult

Staph

Staph

Staph

inhibiting
growth of Staph. aureus
by
mould



Alexander Fleming's data recordings: 1928

Oct. 30. 28.

Extract of staph inhibiting mould made in Yeast Centrifuged
(24 hours at 37°C)

To 0.5 cc of this 0.25 cc of staph emulsion added and
incubated at 45° and 56°C

In 3 hrs @ 45° considerable lysis of staph in tube with mould extract
 56° no visible difference from control.

\therefore Mould culture contains a bacteriolytic substance for staphylococci

Life is mostly composed of the elements carbon, hydrogen, nitrogen, oxygen, sulfur, and phosphorus. Although these six elements form nucleic acids, proteins, and lipids and thus the bulk of living matter, it is theoretically possible that some other elements in the periodic table could serve the same function. We describe a bacterium, strain GFAJ-1 of the Halomonadaceae, isolated from a hydrothermal vent in the Salton Sea, California, that is able to substitute arsenic for phosphorus to sustain its growth. Our data show evidence for arsenate in macromolecules that normally contain phosphorus, notably nucleic acids and proteins. Exchange of one of the major bio-elements may have profound evolutionary and geochemical importance.

Science

SCIENCE

AAS

A Bacterium Arsenic Instead

Felisa Wolfe-Simon,^{1,2*} Jodi Swanson,¹
Shelley E. Hoefft,² Jennifer Pett
Paul C. W. Davies,^{1,7} Ariel D. A

NATURE

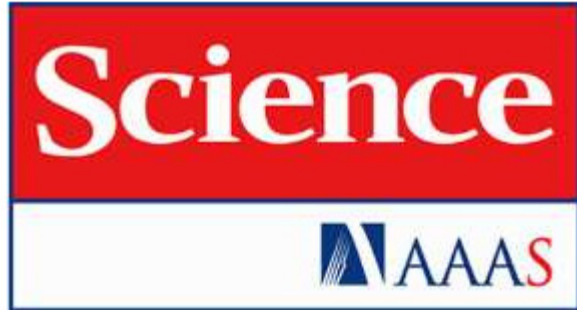
Methods: Face up to false positives

Daniel MacArthur

Nature **487**, 427–428 (2012) doi:10.1038/487427a

‘Scientists and journals must work together to ensure that eye-catching artefacts are not trumpeted as genomic insights’

‘hunting for biological surprises without due caution can easily yield a rich crop of biases and experimental artefacts, and lead to high-impact papers built on nothing more than systematic experimental ‘noise’.



12 December 2014



When contact changes minds: An experiment on transmission of support for gay equality

Michael ... and Donald P. Green²

Can a single (20 minute) conversation change minds on divisive social issues, such as same-sex marriage?

“LaCour has not produced the original survey data from which someone else could independently confirm the validity of the reported findings.” --Marcia McNutt, Editor-in-Chief, May 2015

RETRACTED

Global data integrity crisis

nature 2012

Must try harder

Too many sloppy mistakes are creeping into scientific papers. Lab heads must look more rigorously at the data – and at themselves.

Believe it or not: how much can we rely on published data on potential drug targets?

See also: [News and Analysis by Arrowsmith](#)

Florian Prinz¹, Thomas Schlange² & Khusru Asadullah³

Bayer Healthcare

Global data integrity crisis

“...we will soon issue an international good practice for regulatory authorities and inspectors that can help reduce incidents of incomplete presentation of data by manufacturers or deliberate data falsification.”

March 2015

Dr Margaret Chan, Director-General, World Health Organization

<https://goo.gl/LHQKdW>

FDA's draft guidance on data integrity

<https://goo.gl/4sAhPR>

April 2016

Theater makeup artist turned image forensics expert



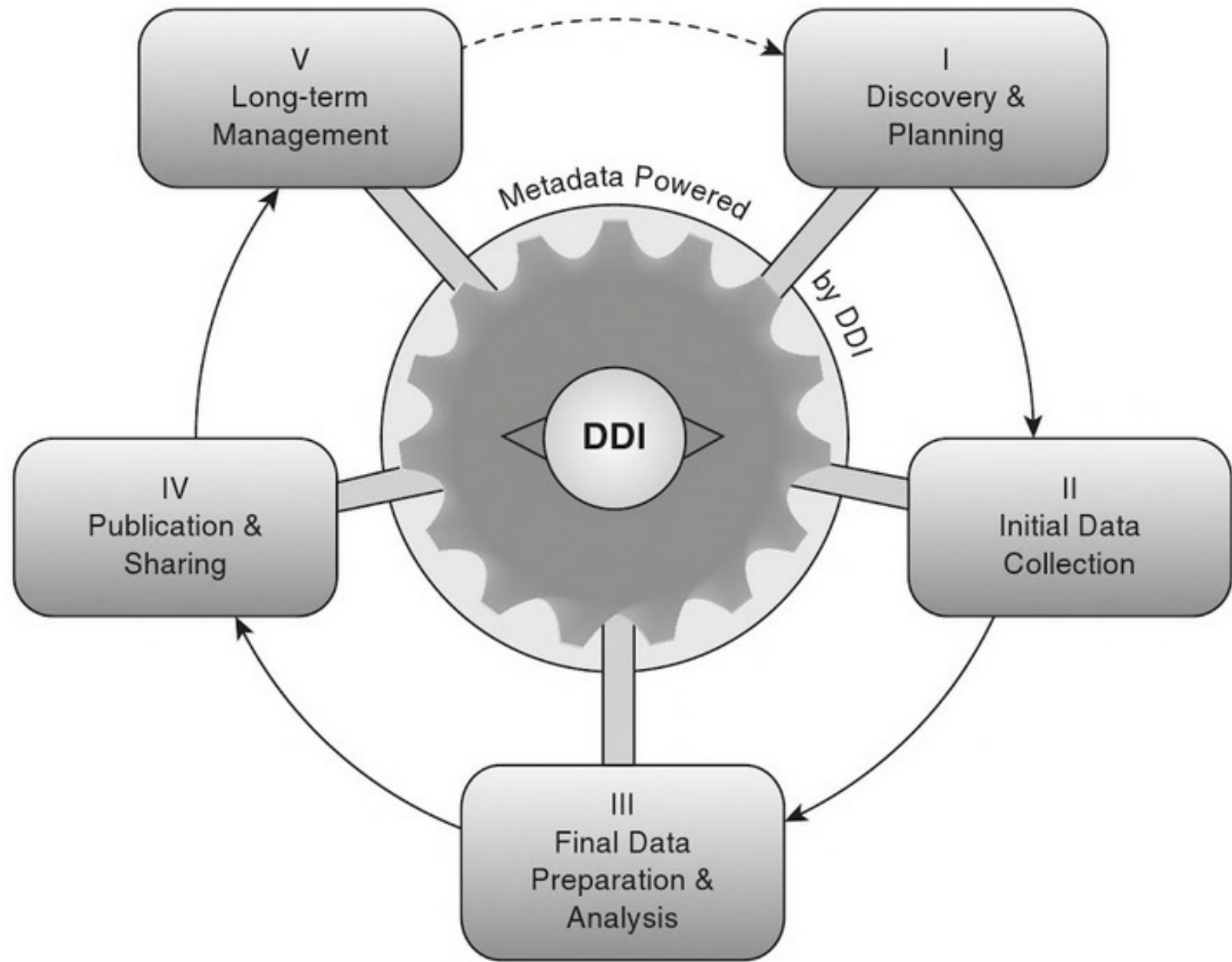
Jana Christopher

Ensuring data integrity in science

Framework proposed by Committee on Science, Engineering, and Public Policy (COSEPUP) of the National Academy of Sciences (NAS), National Academy of Engineering, and Institute of Medicine (2010).

- ❖ Urgent attention required on two fronts:
 - (a) “general practice of science” and
 - (b) “personal behaviors of scientists”
- ❖ “Reinforce clarity and transparency to build and maintain trust in science.”
- ❖ “Teach scientists to describe experiments, data, and calculations fully so that other scientists can replicate the research”

Research data lifecycle: Data documentation initiative

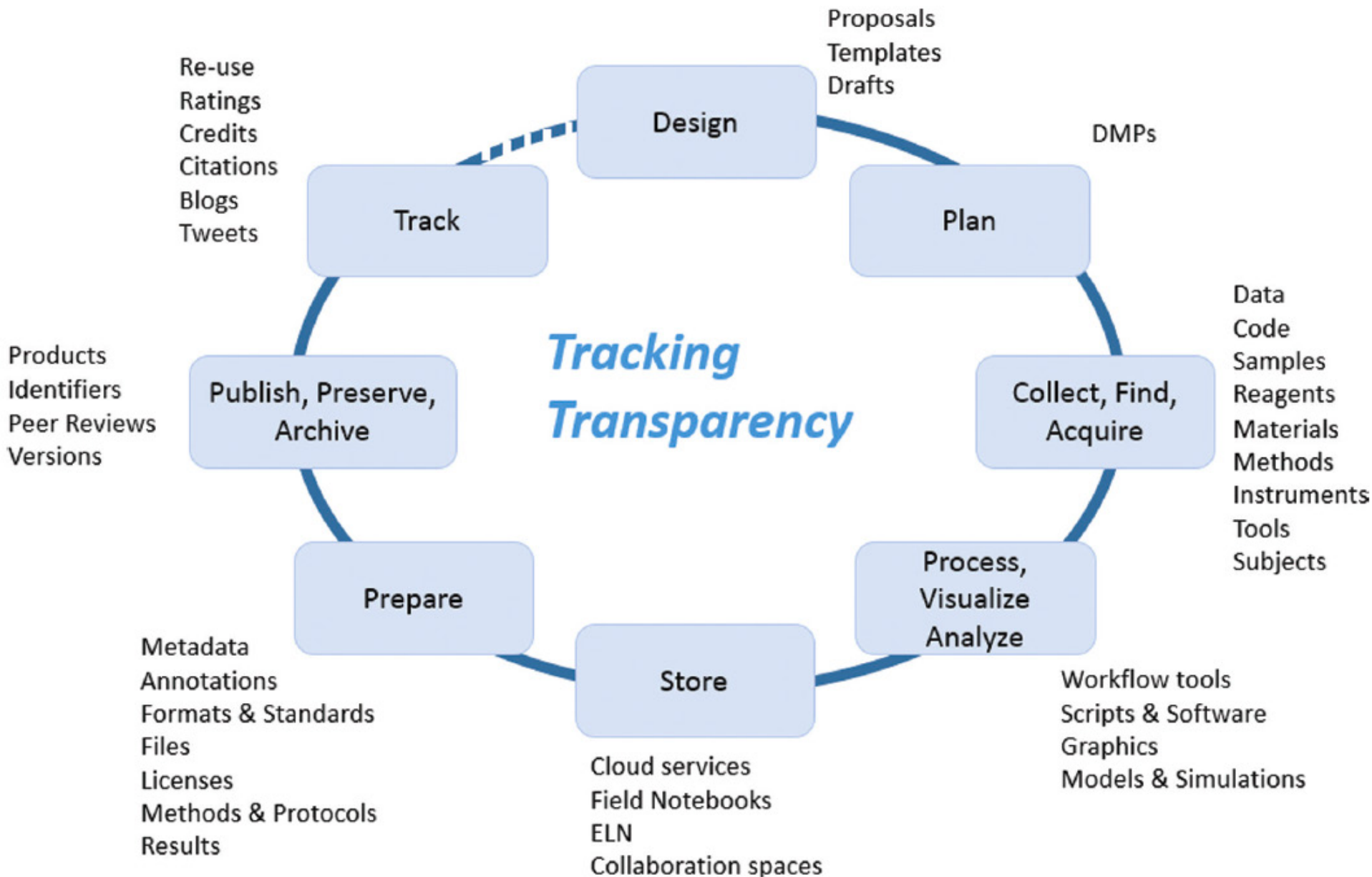


Source: DDI (2013)

© 2014 Corti

Urgent need for establishing explicit analytical workflows

Research data lifecycle: Transparency



Key resource: Office of Research Integrity RDM site



RCR Administrators

topics resources topic contents

Data Management

Tutorial

- Opening Case

Ownership of Data

- Overview
- Federal Policies
- Institutional Policies

Access to Data

- Stewardship Responsibilities
- Who has access?
- Data Control

Retention

- Data Retention

Problems

- Examples of Problems

Case Studies

- Not My Job...
- But We Have a Policy!
- Data Wars

Assessment

- Quiz

Resources

- Bibliography
- Regulations & Policies
- Glossary

Whose Data Is It?

case study

Prof. Faith Promise has taken a faculty position at another institution. Upon learning of her final decision to leave, her Department Chair, Duncan Pheef, asks her whether she will be taking the original data generated under her Department of Defense contract or whether she will be taking a copy. Promise is confused and says, "I'll be taking the original, of course. It's my data! DoD expects me to keep it so I can continue my work." Pheef explains that the university has the responsibility to retain the data generated under the contract. Promise can take a copy and leave the original, or she can take the original and leave a copy. If she takes the original data, then she has to promise that she will give the university access to the original data if it needs such access. Promise becomes increasingly irritated and says angrily, "Look Duncan, this is my contract and my data! I am not going to copy all of my data just so you can file it away or maybe even give it to someone else. I'm taking my data and if anyone needs it, they can call me and if I have time, I'll make them a copy!"

This case illustrates common assumptions and problems involved in the management of data resulting from the conduct of sponsored projects. Discussions over who owns and who has access to research data can be very contentious. As will be seen in this tutorial, research data is only one of three types of information that institutions must manage appropriately, responsibly, and in compliance with both their own and sponsors' policies.

How is data defined?

The word Data is defined in Webster's 11th Collegiate Dictionary as "factual information (as measurements or statistics) used as a basis for reasoning, discussion, or calculation." That definition confirms what is commonly thought about data when referring to research projects. That is, "data" means the all of the information collected and generated in the course of a research project.

Data Management: Opening Case



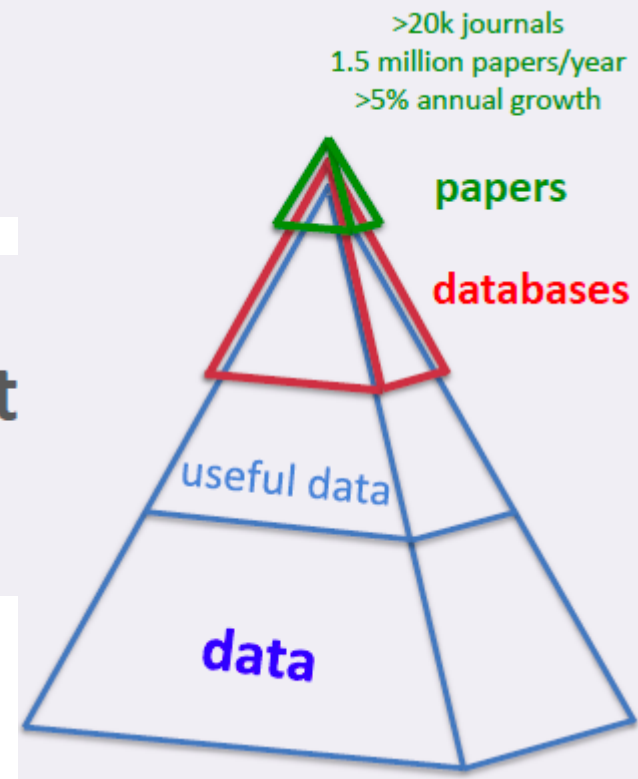
Administrators and the Responsible Conduct of Research

Got data?

Data published in papers represents a small fraction of all *useful* data generated in labs
peer review; stable; citable; usually unstructured

Data deposited in databases represents a bigger fraction of the data
usually curation; (stable); citable; structured

Data deposited in repositories may capture a large fraction of the data
some curation; often unstructured
>> validation, citability, stability



Useless data

- Raw
- Unstructured
- Unreproducible
- Flawed

Useful data

(post validation and curation)

If you plan on using Adobe Photoshop®...

- ❖ Changing brightness and contrast is allowed **ONLY** when you apply it equally across the entire image and equally to all controls.
- ❖ Do **NOT** alter contrast to make your data disappear (e.g. masking)
- ❖ Do **NOT** enhance or emphasize one region of your image while leaving others unchanged
- ❖ Do **NOT** use touch up tools (e.g. cloning, healing, etc) on your data
- ❖ Do **NOT** use features that deliberately obscure manipulations
- ❖ **Be prepared** to supply the Journal editors with original data on request

Enhancing transparency and including source data



SOURCE
DATA



TRANSPARENT
PROCESS



OPEN
ACCESS

Published online: September 13, 2016

Article



SOURCE
DATA



TRANSPARENT
PROCESS



OPEN
ACCESS

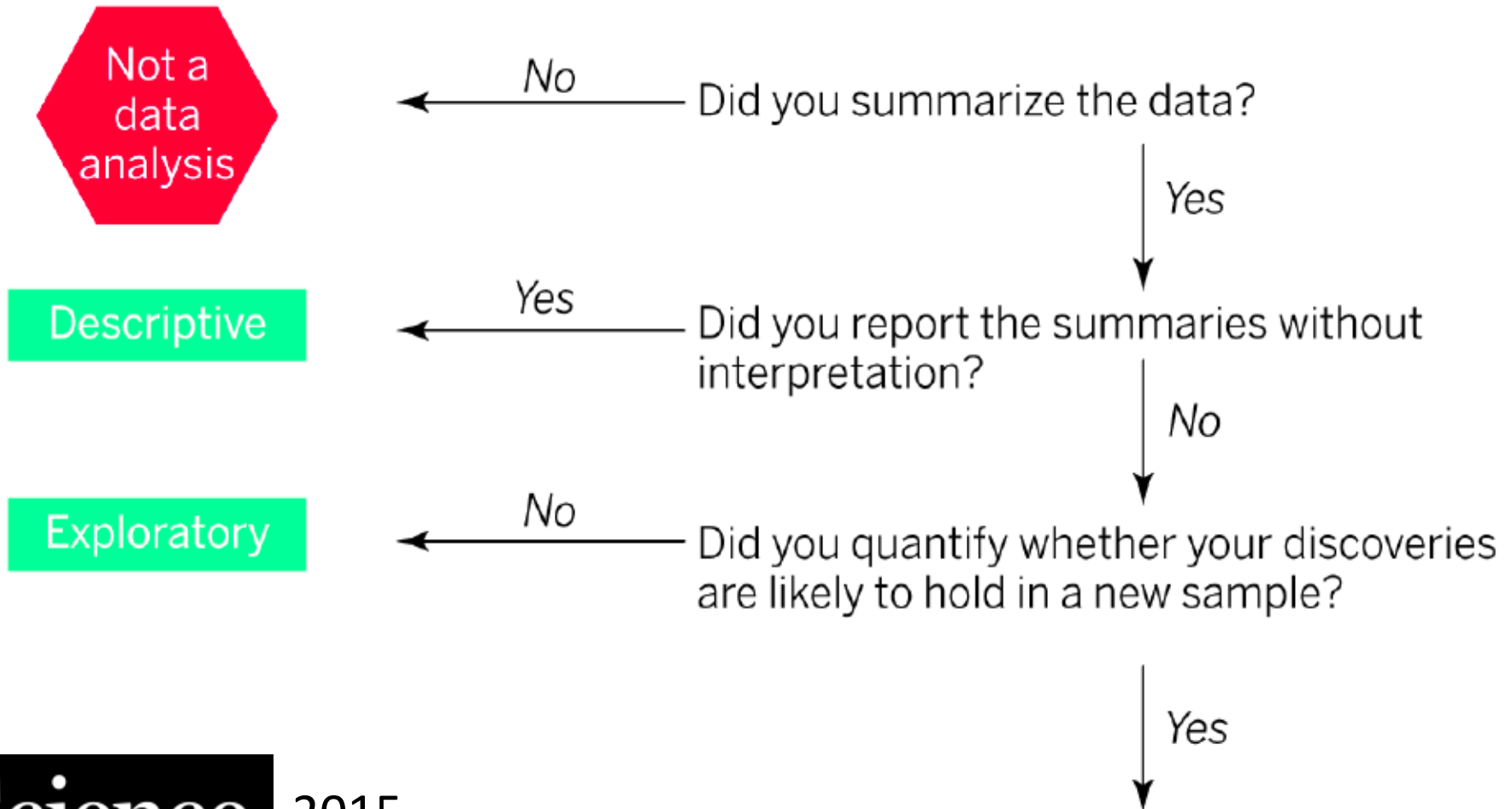
THE
EMBO
JOURNAL

FBW7 suppression leads to SOX9 stabilization and increased malignancy in medulloblastoma

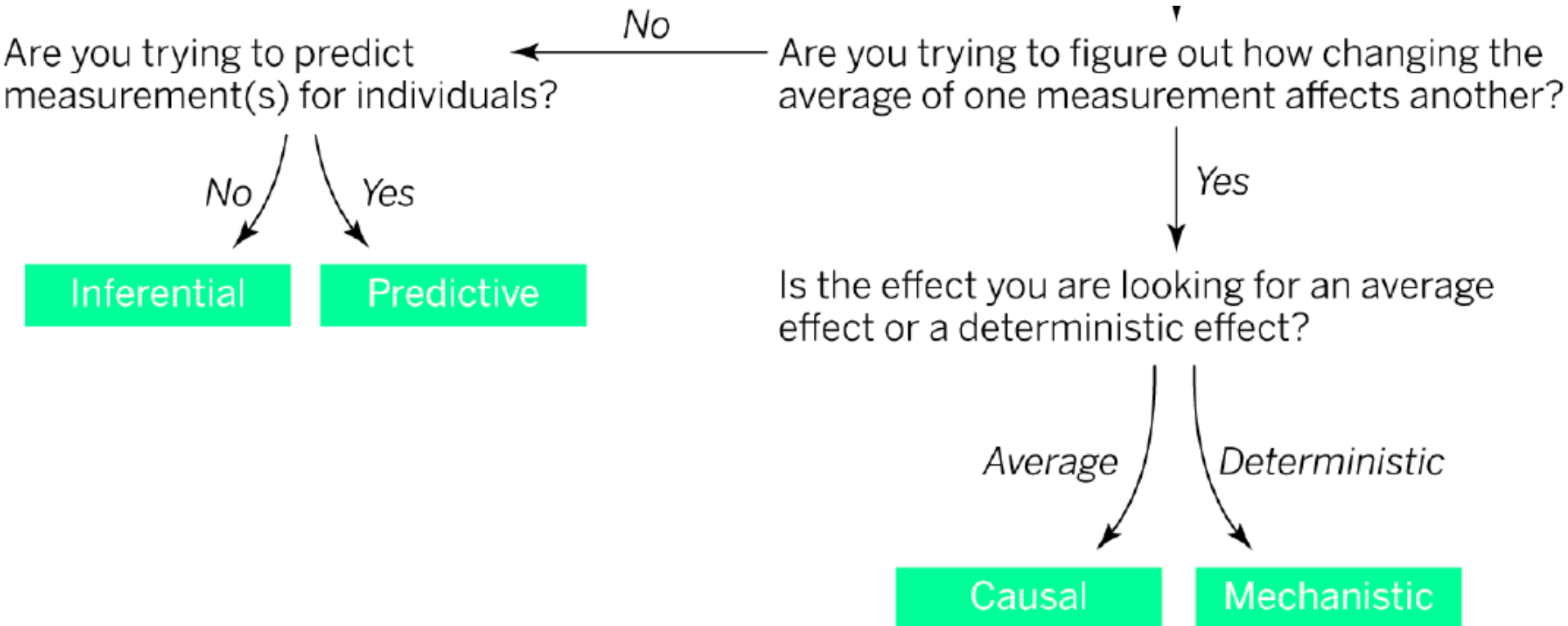
Aldwin Suryo Rahmanto^{1,†}, Vasil Savov^{2,†}, Andrä Brunner^{1,‡}, Sara Bolin^{2,‡}, Holger Weishaupt^{2,‡}, Alena Malyukova¹, Gabriela Rosén², Matko Čančer², Sonja Hutter^{2,3}, Anders Sundström², Daisuke Kawauchi³, David TW Jones³, Charles Spruck⁴, Michael D Taylor⁵, Yoon-Jae Cho⁶, Stefan M Pfister^{3,7}, Marcel Kool³, Andrey Korshunov^{3,8}, Fredrik J Swartling^{2,*,§} & Olle Sangfelt^{1,**,§}

Data analysis: know thy question

“Mistaking the type of question being considered is the most common error in data analysis” Jeff Leek, JHU SPH



Data analysis: know thy question



Data analysis: know thy question

REAL QUESTION TYPE	PERCEIVED QUESTION TYPE	PHRASE DESCRIBING ERROR
Inferential	Causal	"correlation does not imply causation"
Exploratory	Inferential	"data dredging"
Exploratory	Predictive	"overfitting"
Descriptive	Inferential	"n of 1 analysis"

FAIR principles for data stewardship

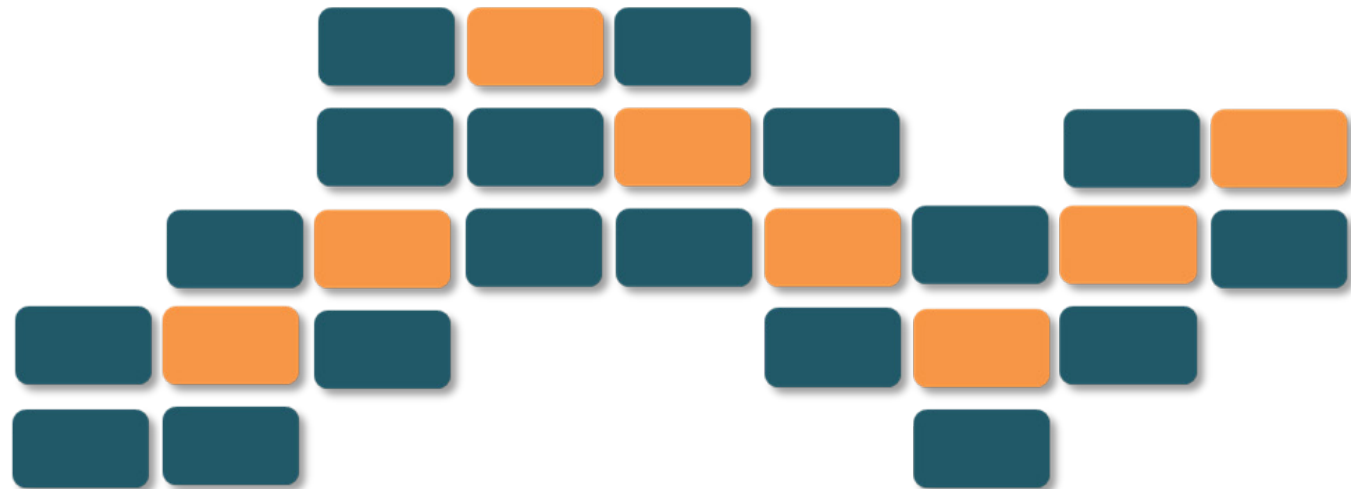
Find

Access

Interoperate

Re-use

Data



FAIR principles for data stewardship

Equally important to good scholarship is the publication of **non-data research objects**. Explicit analytical workflows, for example, are essential to most forms of knowledge generation. Publication of these according to FAIR principles is essential to ensure transparency of the work as well as maximal use to the community. **The key to working with data is to realize that the human touch**, the urge to annotate tables with footnotes and cram multiple elements and data types into every cell of a table, **gets in the way of computation, automation and scaling up**. And this impedes the usefulness of your work for other people. **All research objects should be findable, accessible, interoperable and reusable (FAIR) both for machines and for people.**

A curated, informative and educational resource on inter-related data standards, databases, and policies in the life, environmental and biomedical sciences

Find

 **Recommendations**

Standards and/or databases recommended by journal or funder data policies.

Discover

 **Collections**

Standards and/or databases grouped by domain, species or organization.

Learn

 **Educational**

About standards, their use in databases and policies, and how we can help you.



FAIRDOM



Findable



Accessible



Interoperable



Reusable

FAIRDOM helps you to be in control of collecting, managing, storing, and publishing your data, models, and operating procedures.

The
Dataverse
Project 

Open source research data repository software

Top 10 tips for research data management (RDM)

- ❖ Imagine the worst case scenario
- ❖ Check UMB and funder RDM policy
- ❖ Identify information sources and people who can help you
- ❖ Where are the useful checklists and workflows?
- ❖ Data Management Plan (DMP)
- ❖ Think about what to archive and how to describe and cite it
- ❖ Check out suitable data archives
- ❖ Does your dataset require special software?
- ❖ What are the RDM-associated costs at each stage?
- ❖ Who owns your data and can you make it publicly available?

Data management plan



[Home](#)

[DMP Requirements](#)

[Public DMPs](#)

[News](#)

[Help](#)

Data Management Planning Tool

Create, review, and share data management plans that meet institutional and funder requirements.

Repositories for appropriately consented data

nature
genetics

No impact without data access

A considerable proportion of the usefulness and interest of research publications in our field comes from the data and associated metadata. We therefore insist that data be available for peer reviewers to see and readers to use. Authors should use public permanent repositories designed for appropriately consented data.

Data use under the NIH GWAS Data Sharing Policy and future directions

Dina N Paltoo^{1,10}, Laura Lyman Rodriguez^{2,10}, Michael Feolo^{3,10}, Elizabeth Gillanders⁴, Erin M Ramos², Joni L Rutter⁵, Stephen Sherry³, Vivian Ota Wang², Alice Bailey², Rebecca Baker¹, Mark Caulder⁵, Emily L Harris⁶, Kristofor Langlais¹, Hilary Leeds⁷, Erin Luetkemeier¹, Taunton Paine¹, Tamar Roomian^{2,9}, Kimberly Tryka³, Amy Patterson¹ & Eric D Green² for the National Institutes of Health Genomic Data Sharing Governance Committees⁸

In 2007, the US National Institutes of Health (NIH) introduced the Genome-Wide Association Studies (GWAS) Policy and the database of Genotypes and Phenotypes (dbGaP) to facilitate 'controlled' access to GWAS data based on participants' informed consent. dbGaP has provided 2,221 investigators access to 304 studies, resulting in 924 publications and significant scientific advances. Following on this success, the 2014 Genomic Data Sharing Policy will extend the GWAS Policy to additional data types.

Repositories for appropriately consented data

VOLUME 47 | NUMBER 7 | JULY 2015 | **NATURE GENETICS**

The European Genome-phenome Archive of human data consented for biomedical research

Ilkka Lappalainen¹, Jeff Almeida-King¹, Vasudev Kumanduri¹, Alexander Senf¹, John Dylan Spalding¹, Saif ur-Rehman¹, Gary Saunders¹, Jag Kandasamy¹, Mario Caccamo^{1,5}, Rasko Leinonen¹, Brendan Vaughan¹, Thomas Laurent¹, Francis Rowland¹, Pablo Marin-Garcia^{1,5}, Jonathan Barker¹, Petteri Jokinen¹, Angel Carreño Torres², Jordi Rambla de Argila², Oscar Martínez Llobet², Ignacio Medina¹, Marc Sitges Puy², Mario Alberich², Sabela de la Torre², Arcadi Navarro²⁻⁴, Justin Paschall¹ & Paul Flicek¹

The European Genome-phenome Archive (EGA) is a permanent archive that promotes the distribution and sharing of genetic and phenotypic data consented for specific approved uses but not fully open, public distribution. The EGA follows strict protocols for information management, data storage, security and dissemination. Authorized access to the data is managed in partnership with the data-providing organizations. The EGA includes major reference data collections for human genetics research.

Whither statistics?

Nature Reviews Neuroscience

Power failure: why small sample size

Undermines the reliability of neuroscience

Button, Ioannidis, Mokrysz, Nosek, Flint, Robinson & Munafò
Nature Reviews Neuroscience **14**, 365-376 (2013) doi:10.1038/nrn3475

‘the average statistical power of studies in the neurosciences is very low. The consequences of this include overestimates of effect size and low reproducibility of results. There are also ethical dimensions to this problem, as unreliable research is inefficient and wasteful.’

NATURE

Error prone

Nature **487**, 406 (2012) doi:10.1038/487406a

‘biologists fail to design experiments properly, and so submit underpowered studies that have an insufficient sample size and trumpet chance observations as biological effects....

Researchers ...must agree on standards that will protect against avoidable errors. ’

“There are three kinds of lies: lies, damned lies,
and statistics” Mark Twain

“If your experiment needs statistics, you ought to do a
better experiment.” Ernest Rutherford (Works in fields
with high signal-to-noise)

“If you are going to analyze your data using
statistical methods, plan the methods carefully, do
the analyses seriously, and report the data,
methods, and results completely.” Harvey J.
Motulsky (2015)

10 simple rules: effective statistical practice

Rule 1: Statistical Methods
Should Enable Data to
Answer Scientific
Questions

Rule 2: Signals Always
Come with Noise

Rule 3: Plan Ahead, Really
Ahead

Rule 4: Worry about Data
Quality

Rule 5: Statistical Analysis
Is More Than a Set of
Computations

Rule 6: Keep it Simple

Rule 7: Provide
Assessments of Variability

Rule 8: Check Your
Assumptions

Rule 9: When Possible,
Replicate!

Rule 10: Make Your
Analysis Reproducible

Blind analysis requires serious consideration



Illustration by Dale Edwin Murray

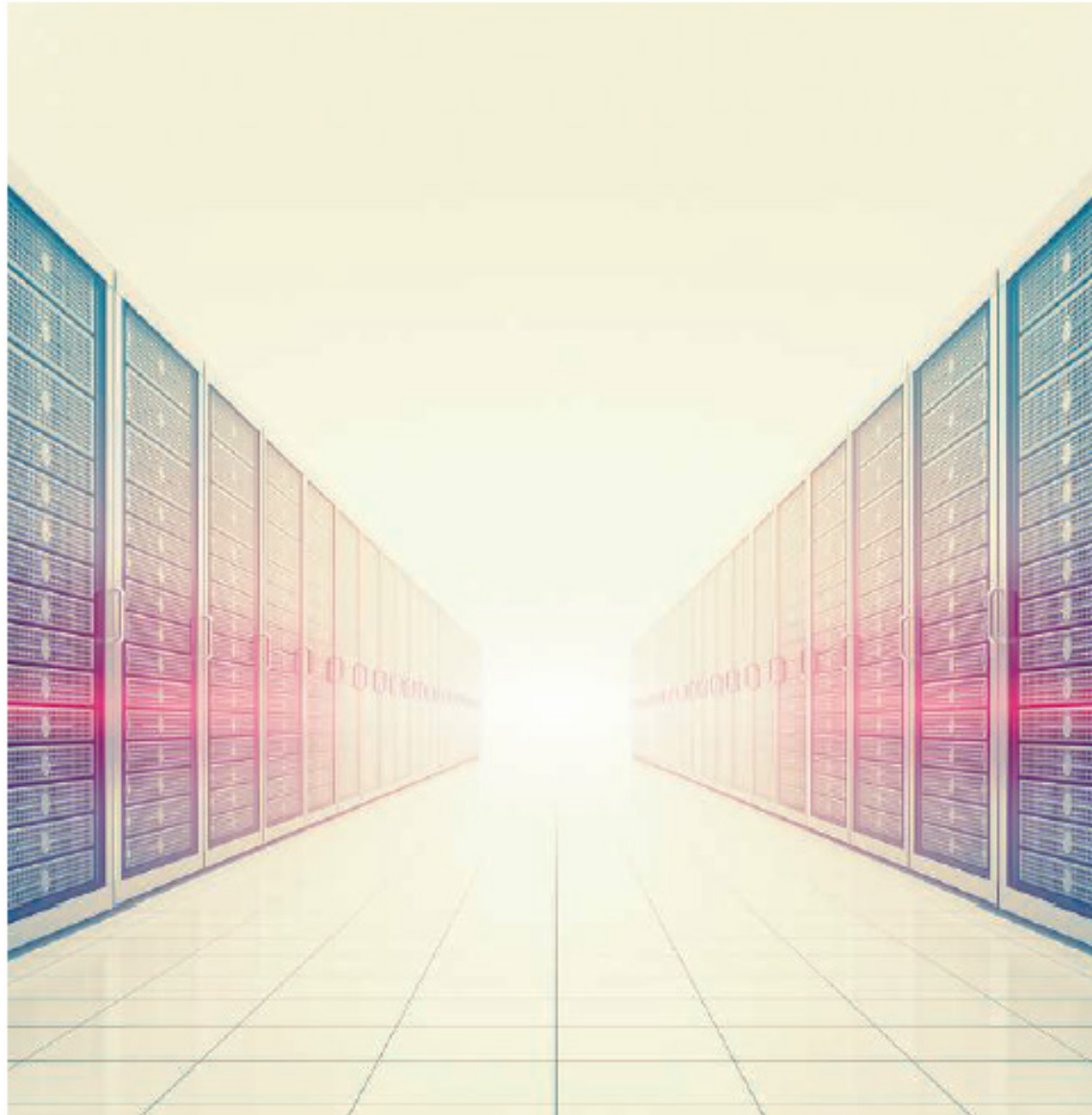
Confirmation bias

“The tendency to look for and perceive evidence consistent with our hypotheses and to deny, dismiss or distort evidence that is not.”

© 2010 Scientific American

“[P]erhaps the best known and most widely accepted notion of inferential error to come out of the literature on human reasoning” Evans, J., (1990) *Bias In Human Reasoning: Cause and Consequences*.

Even “[g]ood scientists are not immune from confirmation bias. They are aware of it and avail themselves of procedural safeguards against its pernicious effects.” (Lilienfeld, S., *Sci. Am.* 303, 18 (2010). *Fudge Factor: A Look at a Harvard Science Fraud Case*)



***“...interpretations come and go,
but data are forever.”***

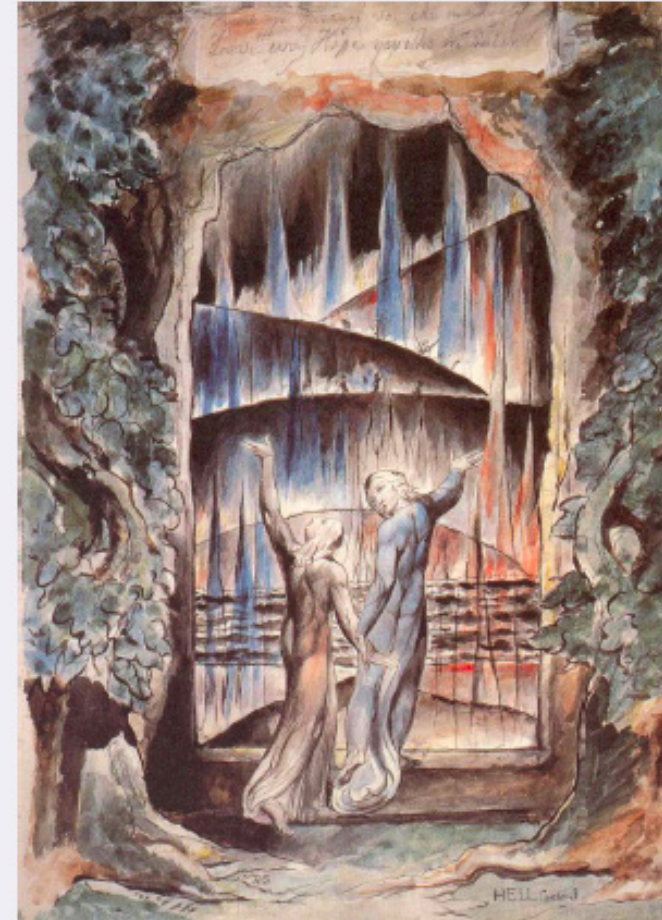
*Marcia McNutt
Editor-in-Chief
Science Journals*

Open data

Opening
the
Gates...



Jacob's Dream - William Blake



Gate of Hell - William Blake