# Research Data Infrastructures
# - Challenges, Desires, Incentives -



**Maik Thomas & Jens Klump**

**Helmholtz Centre Potsdam**
**GFZ German Research Centre for Geosciences**
**Potsdam, Germany**

# Wikipedia's definition

➢ **Data infrastructure** is a digital infrastructure promoting data sharing and consumption.

Similarly to other infrastructures, it is a structure needed for the operation of a society as well as the services and facilities necessary for an economy to function, the data economy in this case.

# Paradigms of science

| | |
|---|---|
| experimental science | empirical understanding |
| ↓ | |
| theoretical science | theoretical description |
| ↓ | |
| computer simulations | explore domains unaccessible for experiments & theory |
| ↓ | |
| data intensive science | combine information sources (observations, simulations, …) |

after Bell et al. (Science, 2009)

GFZ
Helmholtz Centre
POTSDAM

HELMHOLTZ
| ASSOCIATION

# Example:
# German Research Centre for Geosciences (GFZ)

| Department 1 | Department 2 | Department 3 | Department 4 | Department 5 |
|---|---|---|---|---|
| **Geodesy and Remote Sensing** | **Physics of the Earth** | **Geodynamics and Geomaterials** | **Chemistry and Material Cycles** | **Earth Surface Processes** |
| Prof. H. Schuh | Prof. M. Weber | Prof. O. Oncken | Prof. J. Erzinger | Prof. B. Merz |
| 1.1 GPS/Galileo Earth Observation<br>Prof. H. Schuh | 2.1 Physics of Earthquakes and Volcanoes<br>Prof. T. Dahm | 3.1 Lithosphere Dynamics<br>Prof. O. Oncken | 4.1 Reservoir Technologies<br>Prof. E. Huenges | 5.1 Geomorphology<br>Prof. N. Hovius |
| 1.2 Global Geomonitoring and Gravity Field<br>Prof. F. Flechtner | 2.2 Geophysical Deep Sounding<br>Prof. M. Weber | 3.2 Geomechanics and Rheology<br>Prof. G. Dresen | 4.2 Inorganic and Isotope Geochemistry<br>Prof. J. Erzinger | 5.2 Climate Dynamics and Landscape Evolution<br>Prof. A. Brauer |
| 1.3 Earth System Modelling<br>Prof. M. Thomas | 2.3 Earth's Magnetic field<br>Prof. C. Stolle | 3.3 Chemistry and Physics of Earth Materials<br>Prof. W. Heinrich | 4.3 Organic Geochemistry<br>Prof. B. Horsfield | 5.3 Hydrogeology<br>Prof. M. Kühn |
| 1.4 Remote Sensing<br>Prof. H. Kaufmann | 2.4 Seismology<br>Prof. F. Tilmann | 3.4 Earth Surface Geochemistry<br>Prof. F. v. Blanckenburg | 4.4 Basin Analysis<br>Prof. M. Scheck-Wenderoth | 5.4 Hydrology<br>Prof. B. Merz |
| 1.5 Geoinformatics<br>Prof. D. Dransch | 2.5 Geodynamic Modelling<br>Dr. S. Sobolev | | 4.5 Geomicrobiology<br>Prof. D. Wagner | |
| | 2.6 Seismic Hazard and Stress Field<br>Prof. G. Grünthal | | | |

- 5 research areas (Departments) each with 4-6 sections and more than 1000 employees

- several research projects in every section
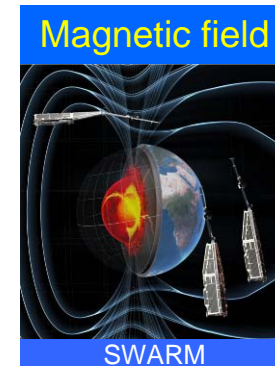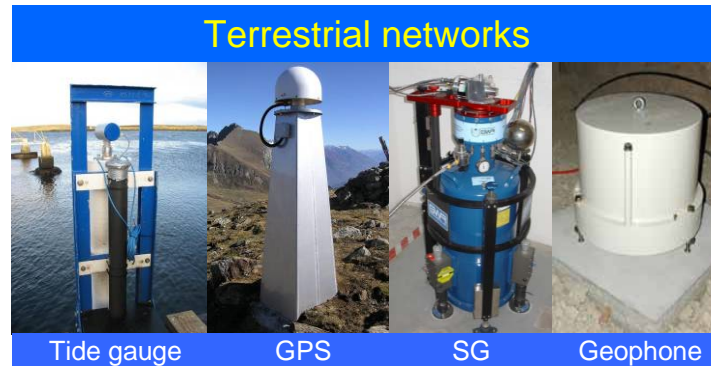
- more than 100 different data sources

GFZ
Helmholtz Centre
POTSDAM

HELMHOLTZ
|ASSOCIATION

# Example: Data sources at GFZ Potsdam

Many data sets come from …

- **Earth monitoring systems**
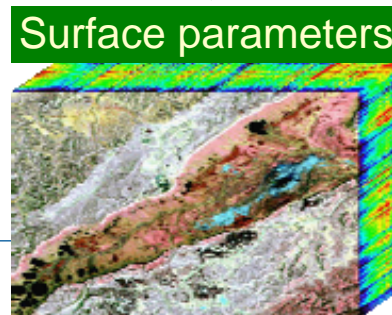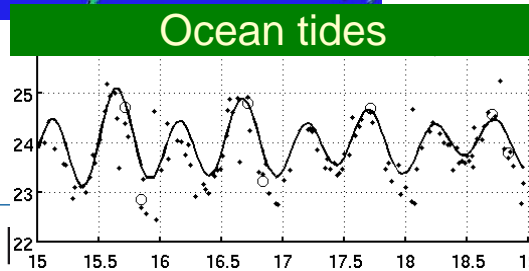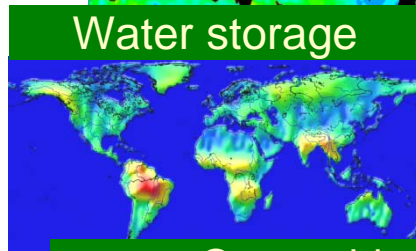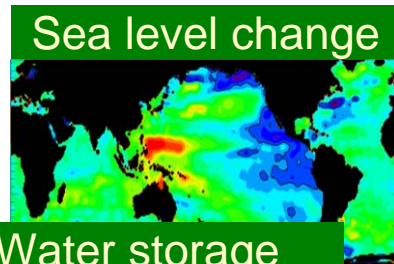  - generate large amounts of data
  - mostly homogeneous structures
  - partly automated workflows

- **Modelling systems**
  - produce medium-sized data amounts
  - partly homogeneous structures
  - but mostly no supporting workflows

- **Labs and field observations**
  - in most cases only small data amounts
  - but heterogeneous structures
  - data sets often generated
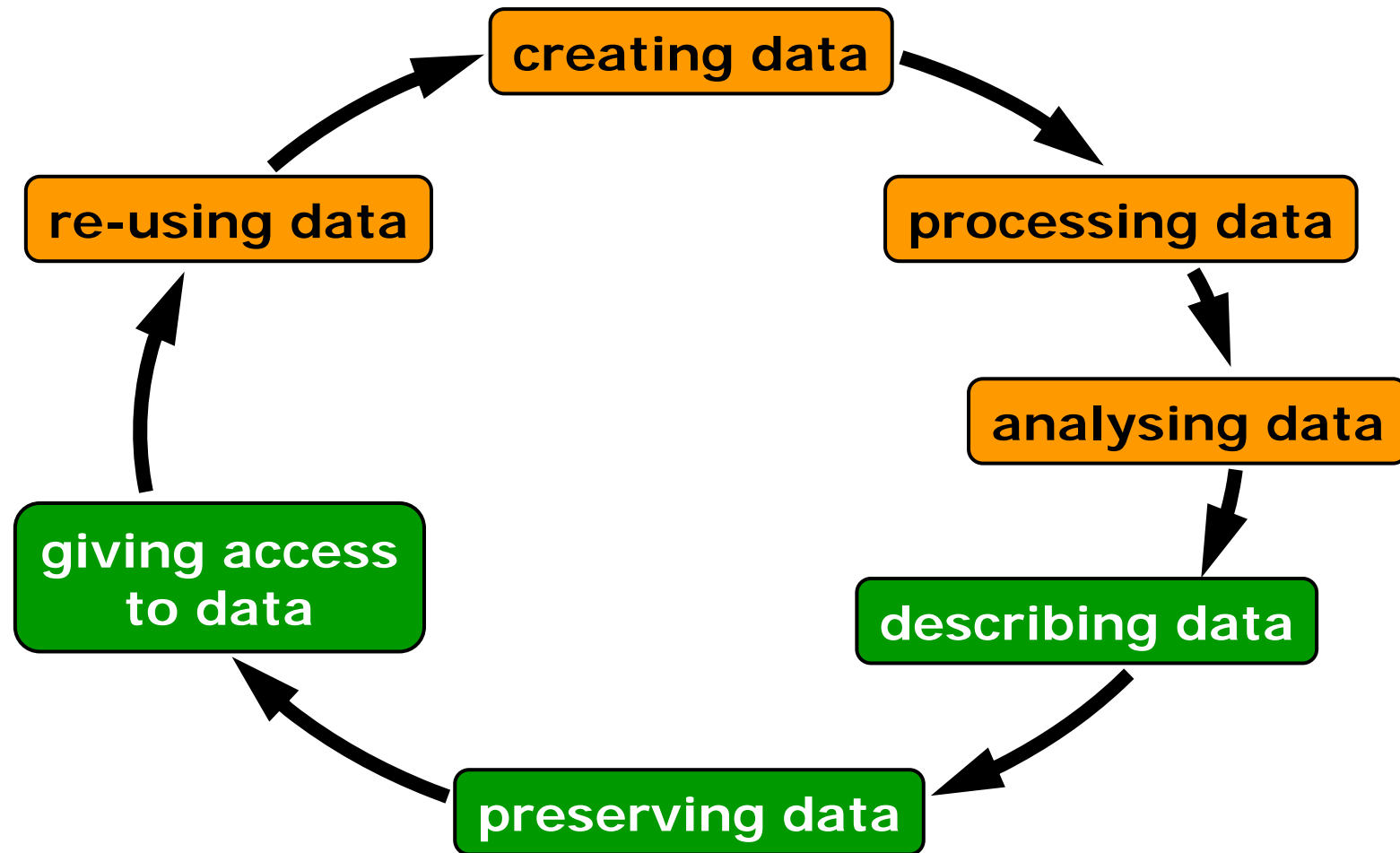    in ad-hoc created structures

# The Earth system: Monitoring, Analysis, Simulation

# Data Life Cycle Model



creating data

processing data

analysing data

describing data

preserving data

giving access to data

re-using data

... defines demands on data management !

GFZ
Helmholtz Centre
POTSDAM

HELMHOLTZ
|ASSOCIATION

Data Life Cycle Model: From private to public

Private Domain — Transfer → Group Domain — Transfer → Persistent Domain — Publication → Access Domain
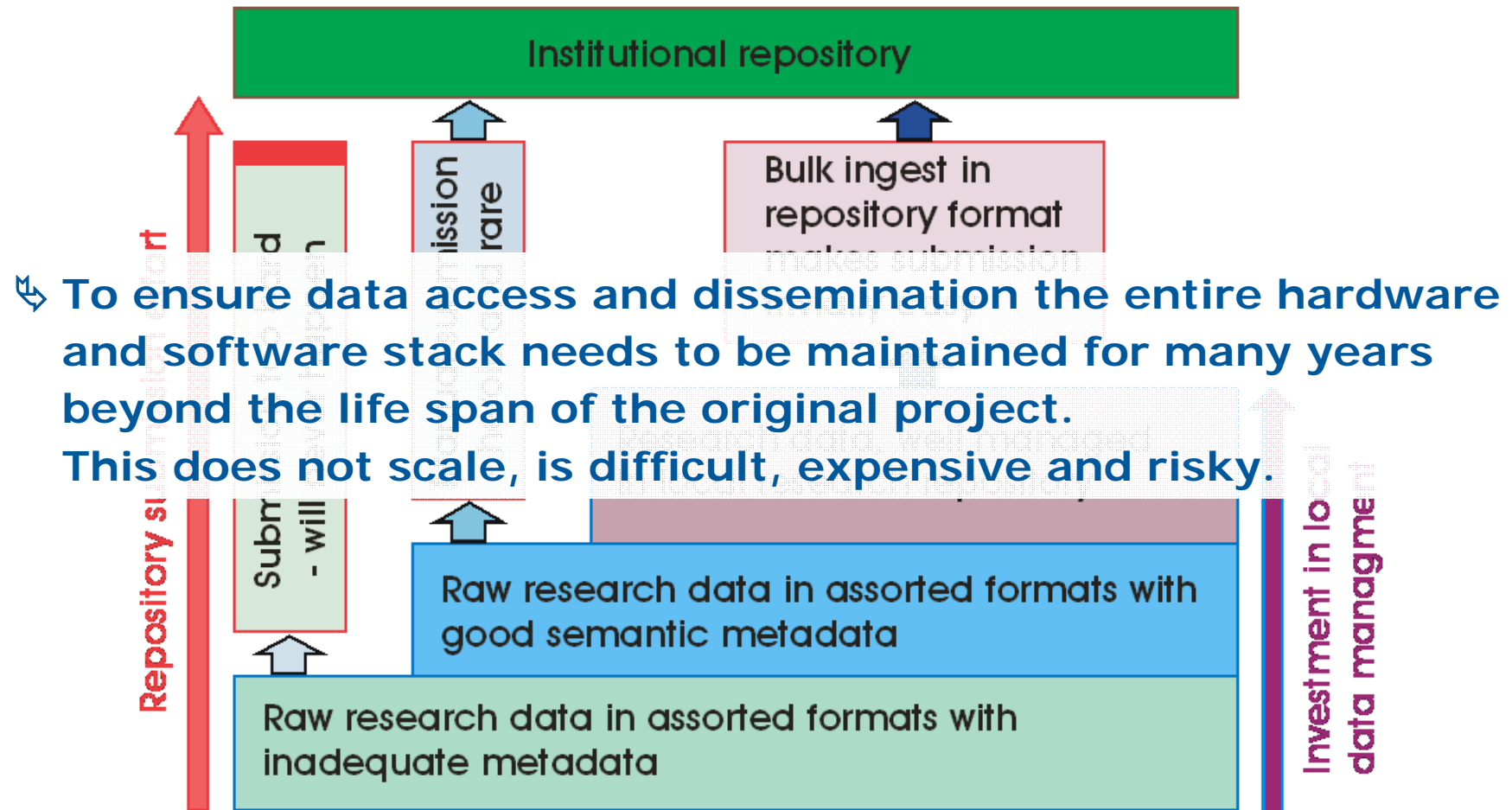
Simple Metadata → Enriched Metadata

↳ How to provide research data sets to external users ?

# Submitting data to repositories



D. Shotton, opencitations.wordpress.com

# Submitting data to repositories



To ensure data access and dissemination the entire hardware and software stack needs to be maintained for many years beyond the life span of the original project.
This does not scale, is difficult, expensive and risky.

D. Shotton, opencitations.wordpress.com

# An isolated application:
# Information System and Data Center (ISDC)

**Gravity fields & models**
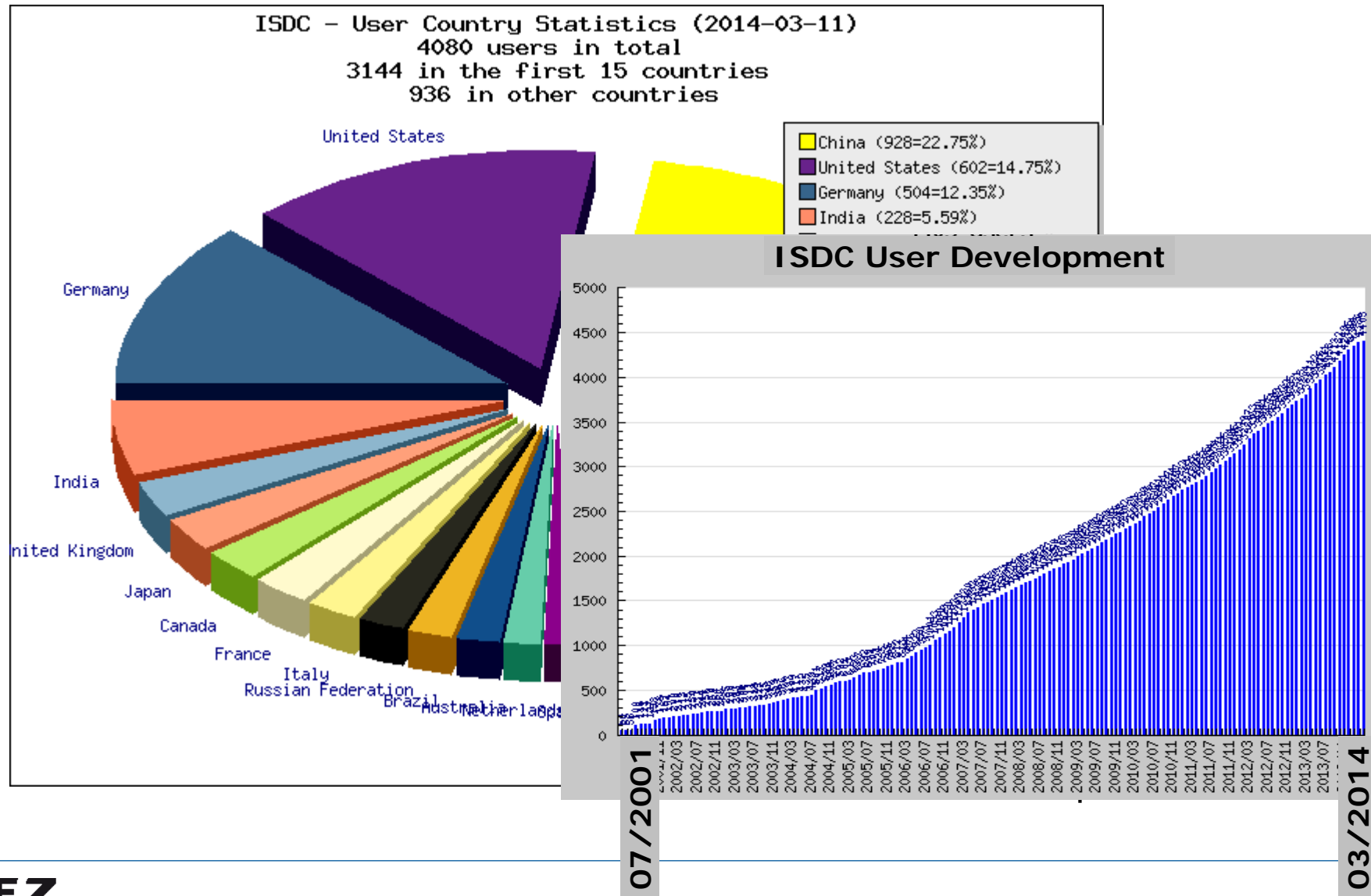
**terrestrial gravimetry**

**GPS data, Galileo-Service**



GGP

GNSS

GPS-PDR

gsp

INFORMATION SYSTEMS AND DATA CENTER

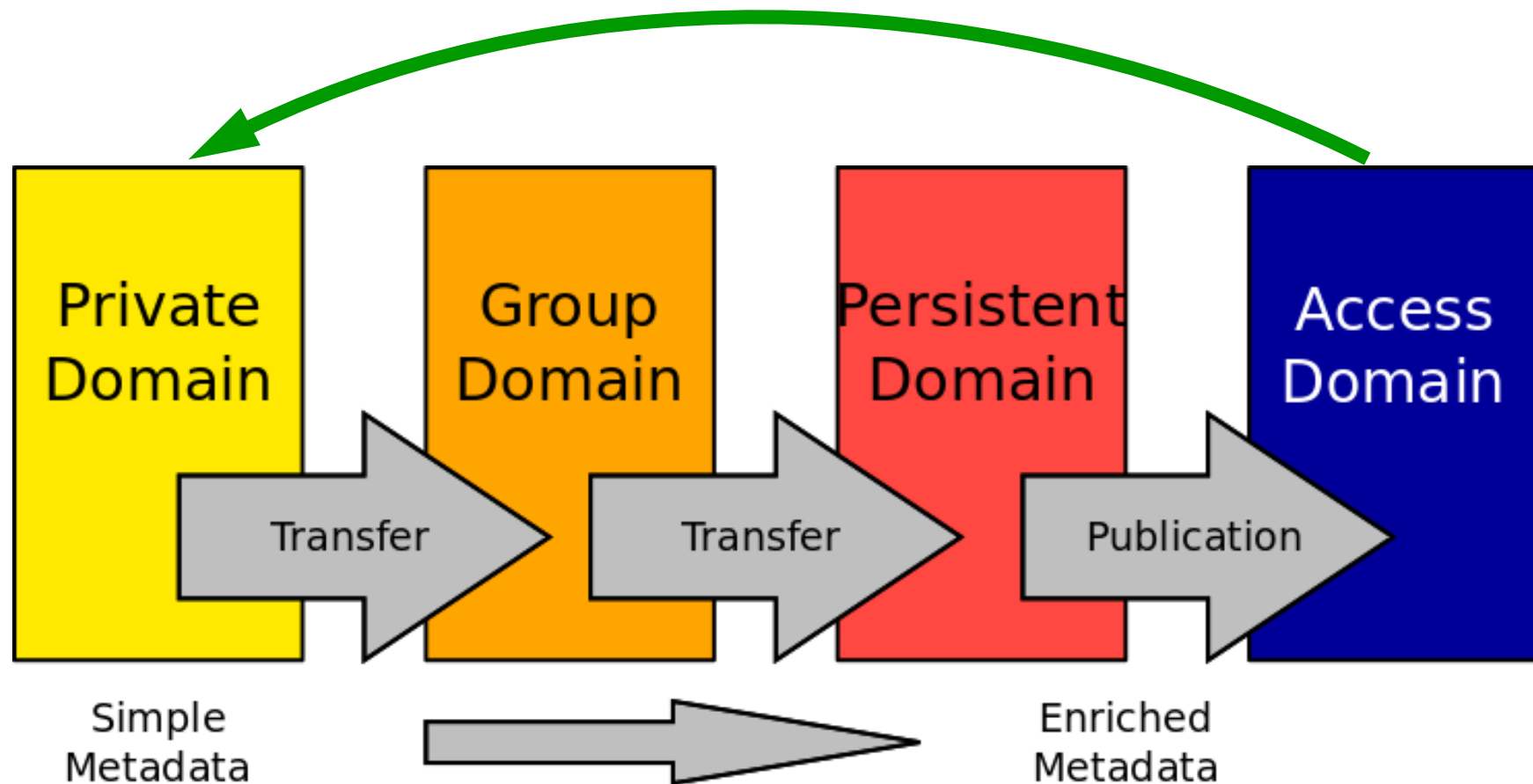isdc.gfz-potsdam.de            Global Earth Science Data

✍ **various data products with corresponding meta data, error estimates FAQs and user support, assigned to major research projects**

# ISDC: Statistics



ISDC – User Country Statistics (2014-03-11)
4080 users in total
3144 in the first 15 countries
936 in other countries

China (928=22.75%)
United States (602=14.75%)
Germany (504=12.35%)
India (228=5.59%)

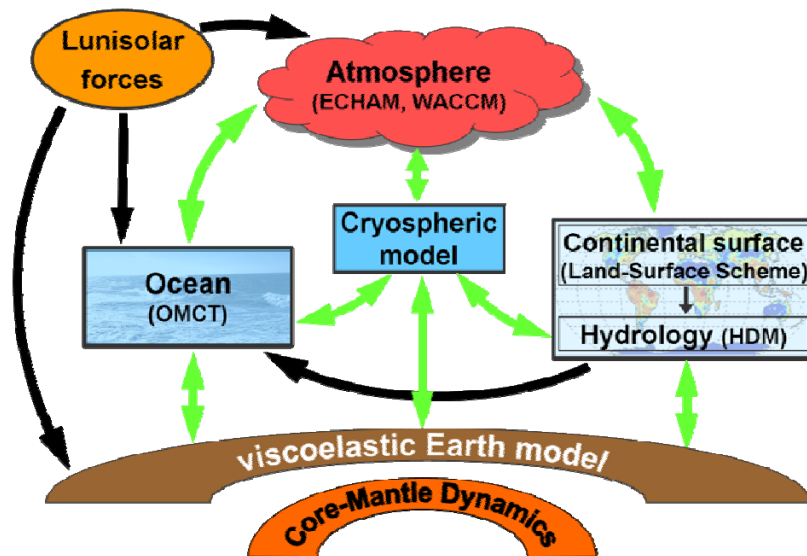ISDC User Development

07/2001          03/2014

# Data Life Cycle Model



What about the integration of data from neighboring research areas ?

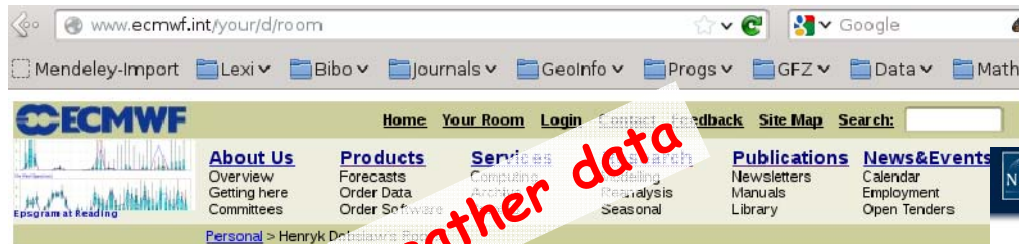# From multi- via inter- to transdisciplinarity

- **differentiation of disciplines** due to scientific specialization over the last decades

- significant **scientific progress** is presently expected from research across specialized disciplines

- recent **IT developments** principally provide excellent conditions **for rapprochement of scientific disciplines** (and research cultures ?)



**Example:**
climate change and its socio-economic impacts
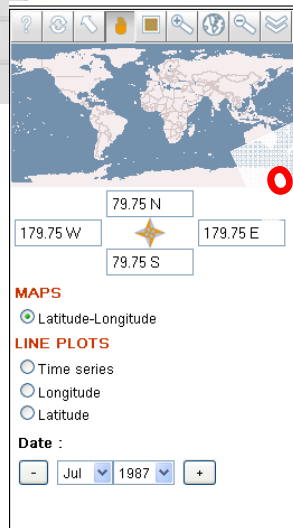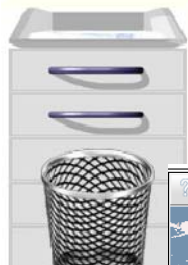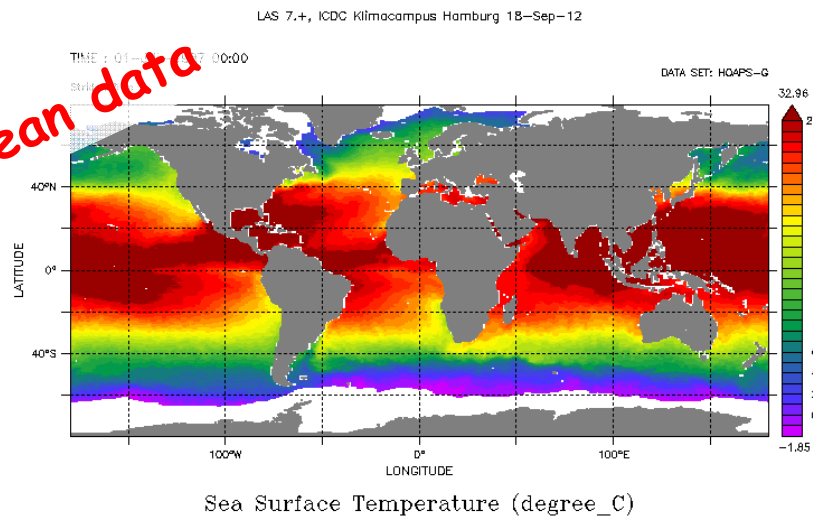
# The geoscientific archipelago

# Many communities and even more „solutions" …



- incompatible standards
- different interfaces
- diverse restrictions
- …

⇨ **poor conditions
for efficient interoperability**

# Expected Results?

EGS Abstract for
Nice, 2000

*accepted !*

## A NEW WAY OF PRODUCING AND ANALYZING INCONSISTENT, IRRELEVANT AND INCOMPREHENSIBLE DATA

A. Hagermann, V. Klemann, K. M. Seiferlin, E. K. Jessberger and T. Stephan ████████████ WWU Münster, D 48149 Münster, Germany. ███████i-muenster.de/Fax: +49 251 83 39083

The ways by which data are produced are sometimes not fully understood. However, data may sometimes bear no significance at all, especially with respect to coffee making or egg frying (████ et al., 1974). In this paper, we present a new method to produce and analyze data that are inconsistent, irrelevant and incomprehensible.

To produce a set of irrelevant data we sampled the distribution of frozen rabbit droppings in front of our institute as cometary analogues. However, due to the climatic peculiarities of the sampling site, none of the samples collected during the period of July 2nd through August 3rd contained the necessary quantities of water ice. The droppings distribution was then analyzed as a function of space and time by means of the Hazelnut transform, with no understanding of the production process whatsoever. These data may have been inconsistent with the findings of other researchers we have probably never heard of. As a method for analyzing these irrelevant data points, ████ (1999) recently proposed a method applying the standard deviation of curvatures in clustered data points as a measure for quantifying an arbitrary line. We cunningly extended this method by inventing an adaptation technique for a semi-empirical line-drawing law for almost any distribution of data points on a double logarithmic scale.

GFZ
Helmholtz Centre
POTSDAM

HELMHOLTZ
|ASSOCIATION

# Data Science: From Big to Small

- **Usually:** Rules for data handling focus on large data sets, i.e., *„Big Data Science"*.
  - ⇨ Existing data management due to regulations of funding bodies.

- *„***Big Data Science***"*: large data amounts, homogeneous structures



(P. Vettiyattil, 2012)

- **But:** Vast investments of time (**and brain work** !) to generate an immense number of small data sets, i.e., *„Small Science Data"*.
  - ⇨ Systematic data management missing !

- *„***Small Data Science***"*: small data amounts, heterogeneous structures

# Policies

**institutional policies**

**interdisciplinary policies**

- Alliance of German Research Organisations:
  - Berlin Declaration on Open Access to Knowledge
    in the Sciences and Humanities (2003)
  - Principles for the Handling of Research Data (2010)
- OECD: Guidelines for Access to Research Data ... (2007)
- DFG: Safeguarding Good Scientific Practice (2013)

**policies of funding organisations**

**journal policies**

- DFG: Proposal Preparation Instructions (2013)
- EU: Guidelines on Data Management in Horizon 2020
- ...

# Example: Policy at GFZ

- DFG-Rules of Good Scientific Practice (GSP)
  are standing instructions.

- Helmholtz-Association supports Open Access initiative
  (1st signee of Berlin Declaration!).

- GFZ's publication policy includes GSP and Berlin Declaration
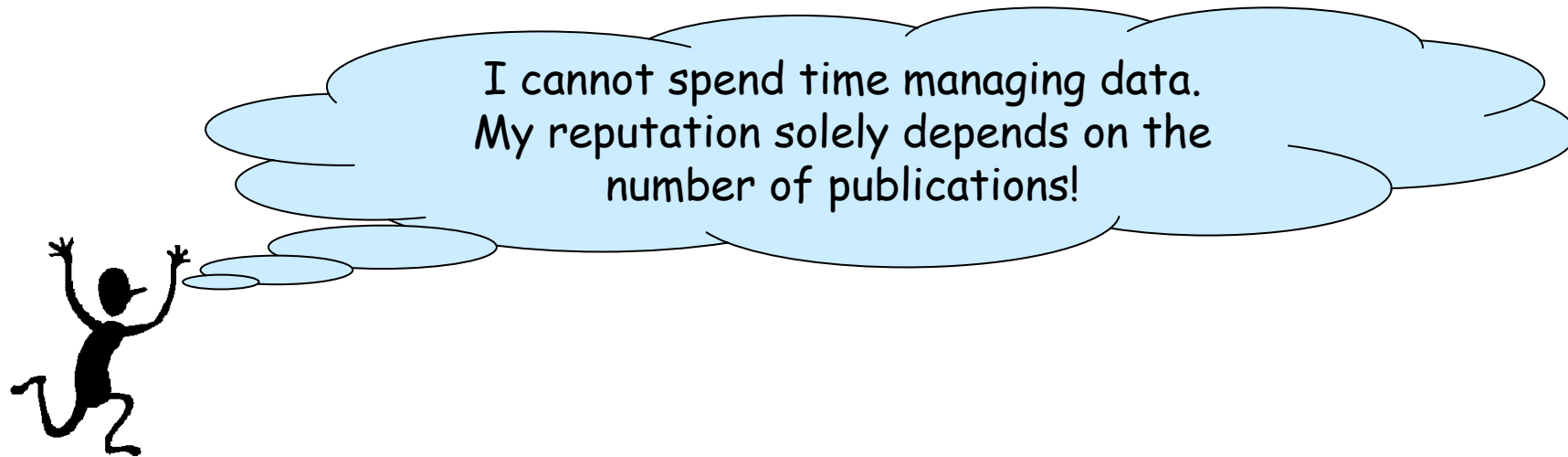  on Open Access and is part of the employment contract.

  ⇪ Since GSP and Open Access are components
  of employment contracts, **data should be accessible.
  But** ...

# The spirit is willing …

How can responsible data handling be established ?

- Imposition of sanctions fail in this context.
  - ↳ Incentives are indispensable for realization !

- Implementation of policies requires:
  - Organization of scientific workflows
  - Tools allowing integration into scientific workflows

- Currently, data handling is not a primary task of scientists. For the individual scientist responsible **data handling is more trouble than it's worth !**

**GFZ**
Helmholtz Centre
**POTSDAM**

**HELMHOLTZ**
**|ASSOCIATION**

# Trusting the cat to keep the cream ...

*I cannot spend time managing data.
My reputation solely depends on the
number of publications!*

**On the other hand:**

Every day several thousand eBay users add metadata
to objects - and they do it correctly.

↳ **Are the tasks of players in data life cycles
correctly assigned?**

GFZ
Helmholtz Centre
POTSDAM

HELMHOLTZ
|ASSOCIATION

# Research data handling ...



... is like driving a flock of cats over plains.

Example GFZ:

- ca. 150 simultaneous research projects

- annual fluctuation of ...
  ... research projects: $\approx 30$
  ... employees: $\approx 200$

↬ **We need generic tools !**

# Describing data



- Content useless without description
- How can I get metadata without sprawling „data bureaucracy" ?

- Collect data descriptions (semi-)automatically as part of the workflow

# Why does it work in other areas?

- Ebay:
  - In order to sell goods people spend much time on describing metadata.
  - motivation: money



- Facebook:
  - People are willing to make vast amounts of personal data public.
  - motivation: self-expression, networking



- Wikipedia:
  - People write and edit articles on specific topics without being acknowledged as authors.
  - motivation: prospect of becoming a Wikipedia administrator

## ... sometimes it also works in the scientific world:

# Institutional Workflows

**Examples:**

ZALF (Leibniz Centre for Agricultural Landscape Research)
- Requests for lab performances are fed online.
  - Metadata of experiments are already recorded.

Rutherford Appleton Laboratory
- Applications for using instruments are made online.
  - Metadata of experiments are automatically recorded.

University of Bremen
- Graduates in geology have to store their research data
  in the database PANGAEA before the certificate is delivered.

# Example: PANGAEA (www.pangaea.de)



- Data sets are not sorted, but searched for.

- At the time of incorporation a DOI is added to the data sets.

- No incorporation of operationally increasing data, so far.

# PANGAEA:
# Cross-linking papers with research data

# Basic institutional measures

- Implement infrastructure for systematic publication of research data by means of Digital Object Identifiers (DOI).

- Install virtual research environments to integrate data management into scientific work flows.

- Install adequate data repositories (centralized vs. decentralized; internal vs. external).



GFZ
Helmholtz Centre
POTSDAM

HELMHOLTZ
| ASSOCIATION

# Once again: Institutional workflows

| Notification |
| :---: |

↓

| Publication DB |
| :---: |

↓

| IR |
| :---: |

↓

| Data Repository |
| :---: |

- In general, required policy components already exist at institutions.

- Bibliometric indicators are crucial factors for scientific reputation !

- Institutional support of workflows (library, administration, data center, ...)

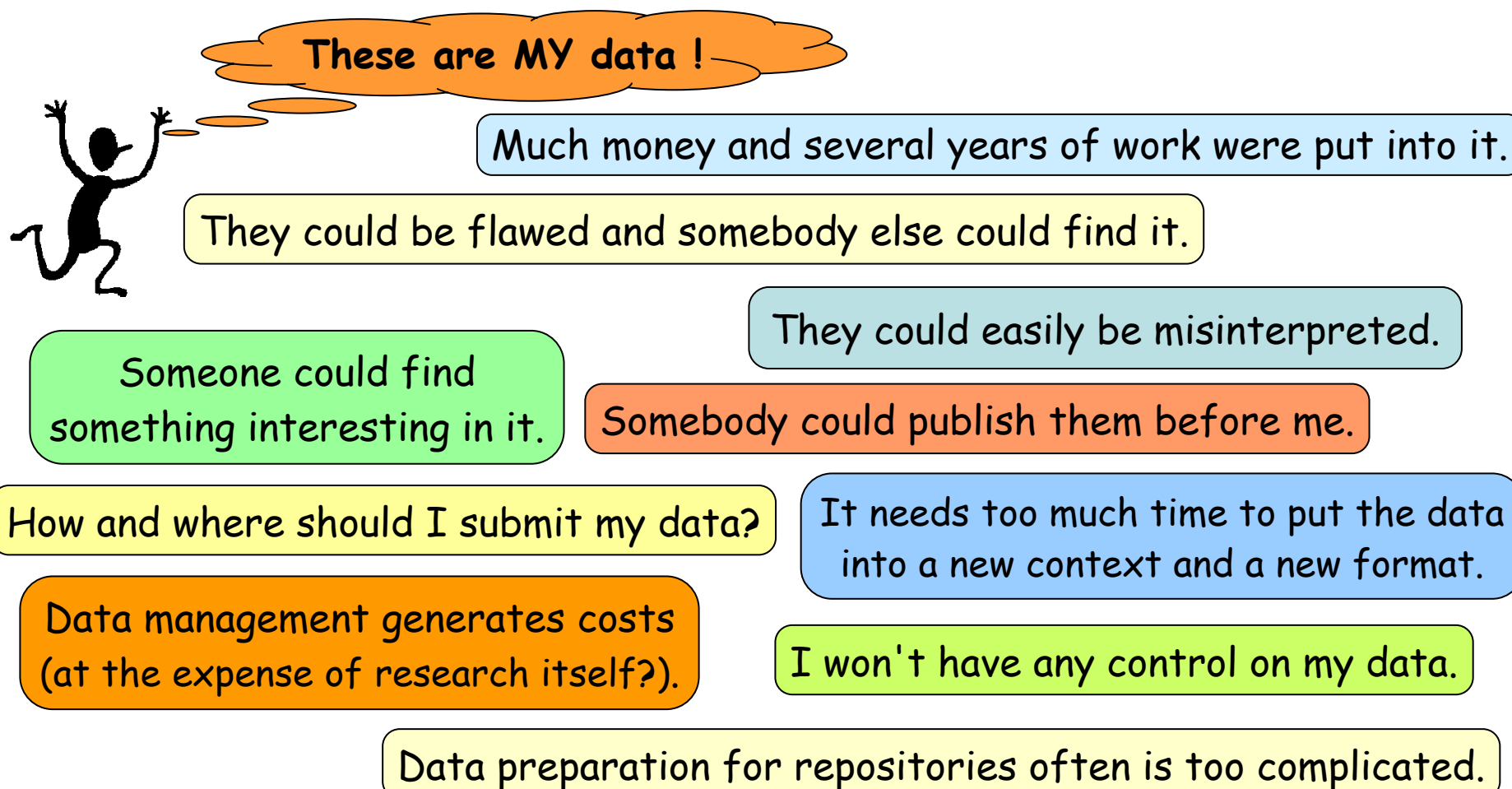- Platform for data publication (e.g., provided by library)

**GFZ**
Helmholtz Centre
**POTSDAM**

**HELMHOLTZ**
**ASSOCIATION**

# Vision


danblackonleadership.info

- **Vision** of research organisations:

  „*Research data should be openly accessible for the long-term and distributed without restrictions.*"

- **Intended purpose:** Open and long-term access ...
  - promotes new findings and research approaches,
  - allows new interpretations,
  - is of strategic importance for the competetiveness of science locations
  - enables syntheses of research results of different disciplines,
  ↬ accelerates the process of making scientific findings !

- **But ...**

# Challenges


123RF Limited

- "Temporary scientists generate the major share of research data." (and are gone ...)

- Integration of data management into work flows

- Conflicting areas of „Big Data Science" vs. „Small Data Science"

- Interoperability of discipline specific, international and across discipline structures

- Data management as obligatory component of training ?

- Establishment of incentive schemes

- Cultural change: availability of research data as a matter of course

# Desires

- Visualization of mutual „*benefits*"

- Adapting data management, tools and services into typical work flows in disciplines

- Generic tools that are easily transferable (e.g. within temporary projects)

- Flexible interfaces for adaptation to needs and previous knowledge of users

- Transparency: *What happens with my data ?*

- Balance between information science and service



I should have wished for more wishes

sam brown, explodingdog

# Incentives:
# On honour, fame and fortune

- **Reputation**:
  - Establishment of a culture of data publication
  - Integration in reputation systems (h-index, …)

- **New professions**:
  - Bridging the discontinuity between science and information management

- **Provision of science-related tools**:
  - generic (and comprehensible) interfaces, virtual research environments, …

- **Service-oriented infrastructures**

⇨ Recognition of the data provision for further use as a natural component in science

RESERVED FOR EMPLOYEE OF THE MONTH

© SmartSign

Great Job!

© 2014 - Positioning Systems

# Imagine …

… that backup and archiving, sharing and reuse
of research data would be possible and is reality !

✓ **We would have more data per discipline.**

✓ **We could compare many differing data.**

✓ **Doubling of studies could be avoided.**

✓ **Quality control would be guaranteed.**

✓ **Other data sets could support own findings.**

✓ **Own results could be supported by others.**

✓ **Unique, non-reproducible results would be documented.**

✓ **New interpretations would be possible.**

➯ **Make your data usable, not just accessible !**

Thank you for your attention !