

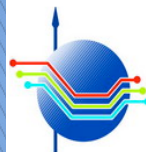
The name is Science! d-Science!

Yannis Ioannidis

MaDgIK Lab

Dept. of Informatics & Telecom, Univ. of Athens

ATHENA Research Center



Ερευνητικό Κέντρο Αθηνά

Ερευνητικό Κέντρο Καινοτομίας στις Τεχνολογίες
της Πληροφορίας, των Επικοινωνιών, της Γνώσης



National and Kapodistrian
UNIVERSITY OF ATHENS

DEPT. OF INFORMATICS
& TELECOMMUNICATIONS

Dedication / Credits

Jim Gray (+)

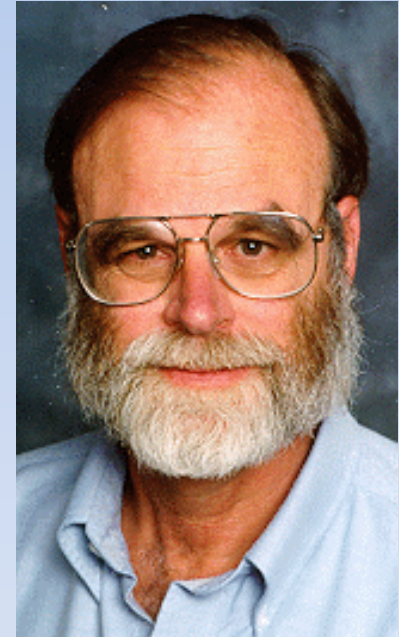
Missing at sea since 28 January 2007

Turing Award 1999

Sigmod Innovations Award

US National Academy of Engineering

...




Science Paradigms

- ▶ 1st - Thousand years ago:
science was empirical
describing natural phenomena
w/ some models, generalizations
- ▶ 2nd - Last few hundred years:
theoretical branch
using models, generalizations



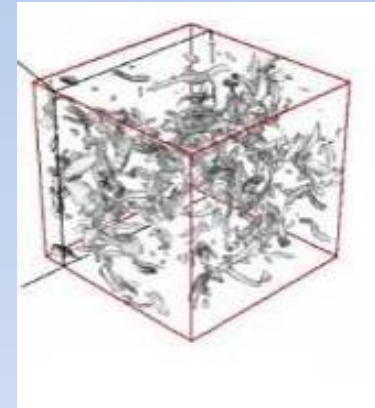
$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K \frac{c^2}{a^2}$$

Really Early Times


- ▶ One scientist
 - ▶ One location
 - ▶ One discipline
 - ▶ One phenomenon
- 
- ▶ One pencil (... carver ...)
 - ▶ One paper (... stone ...)
 - ▶ Street announcements, e.g., Εύρηκα!

Science Paradigms

- ▶ 3rd - Last few decades:
a computational branch
simulating complex phenomena



Recent Times


- ▶ One small group of scientists
 - ▶ One location
 - ▶ One discipline
 - ▶ One phenomenon
- 
- ▶ One file system
 - ▶ One local disk with custom files
 - ▶ Publications at refereed forums

Science Paradigms

- ▶ 4th - Today:
 - data exploration (eScience)
 - unify theory, experiment, and simulation



Current Times

- ▶ Many/large teams of scientists
 - ▶ Many locations
 - ▶ Many disciplines
 - ▶ Many phenomena
- 
- ▶ Many data management systems
 - ▶ Many data forms
 - ▶ Web uploads for publications, data, processes,
...

d~~e~~-Science

*“... computationally-intensive science that is carried out in highly distributed **network** environments,*

or

*science that uses immense **data sets** that require **computing.**”*

d-Science

As the **amount**, **variety**, and **complexity** of data increases, it becomes more and more apparent that the construction of well articulated information technology is needed for scientific discovery

Data's shameful neglect

Research cannot flourish if data are not preserved and made accessible. All concerned must act accordingly.

More and more often these days, a research project's success is measured not just by the publications it produces, but also by the data it makes available to the wider community. Pioneering archives such as GenBank have demonstrated just how powerful such legacy data sets can be for generating new discoveries — especially when data are combined from many laboratories and analysed in ways that the original researchers could not have anticipated.

All but a handful of disciplines still lack the technical, institutional and cultural frameworks required to support such open data access (see pages 168 and 171) — leading to a scandalous shortfall in the sharing of data by researchers (see page 160). This deficiency urgently needs to be addressed by funders, universities and the researchers themselves.

Research funding agencies need to recognize that preservation of and access to digital data are central to their mission, and need to be supported accordingly. Organizations in the United Kingdom, for instance, have made a good start. The Joint Information Systems

also the software that will help investigators to do this. One important facet is metadata management software: tools that streamline the tedious process of annotating data with a description of what the bits mean, which instrument collected them, which algorithms have been used to process them and so on — information that is essential if other scientists are to reuse the data effectively.

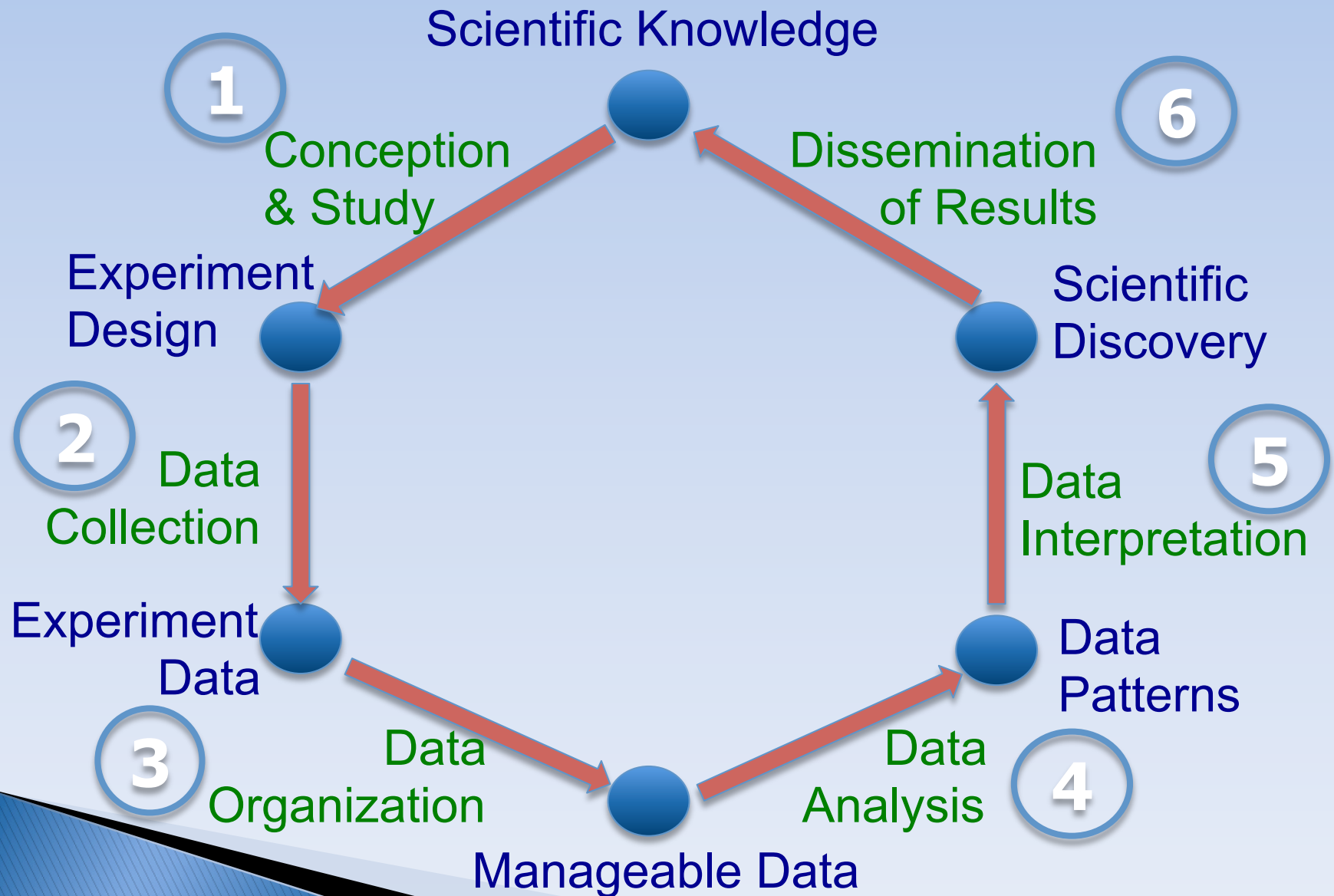
Also necessary, especially in an era when data can be mixed and combined in unanticipated ways, is software that can keep track of which pieces of data came from whom. Such systems are essential if tenure and promotion committees are ever to give credit — as they should — to candidates' track-record of data contribution.

Who should host these data? Agencies and the research community together need to create the digital equivalent of libraries: institutions that can take responsibility for preserving digital data and making them accessible

“Data management should be woven into every course in science.”

More and more often these days, a research project's success is measured not just by the publications it produces, but also by the data it makes available to the wider community. Pioneering archives such as GenBank have demonstrated just how powerful such legacy data sets can be for generating new discoveries — espe-

Research Life Cycle



Scientists' Role

- ▶ Early times: scientist at all stages

- ▶ Current times:

 - Specializations at certain stages

 - Data collection (2) and Result dissemination (6)

 - Scientist at scientific stages

 - Conception and study (1) and Data interpretation (5)

 - Scientist at data management stages (often)

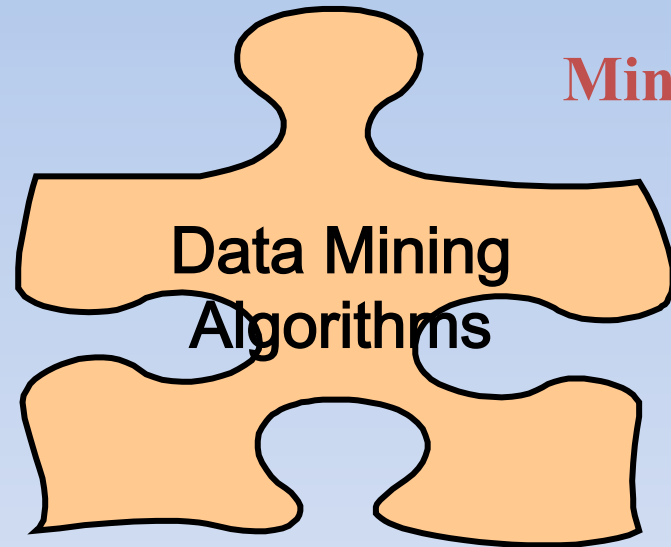
 - Data organization (3) and Data analysis (4)

Roles wrt Informatics

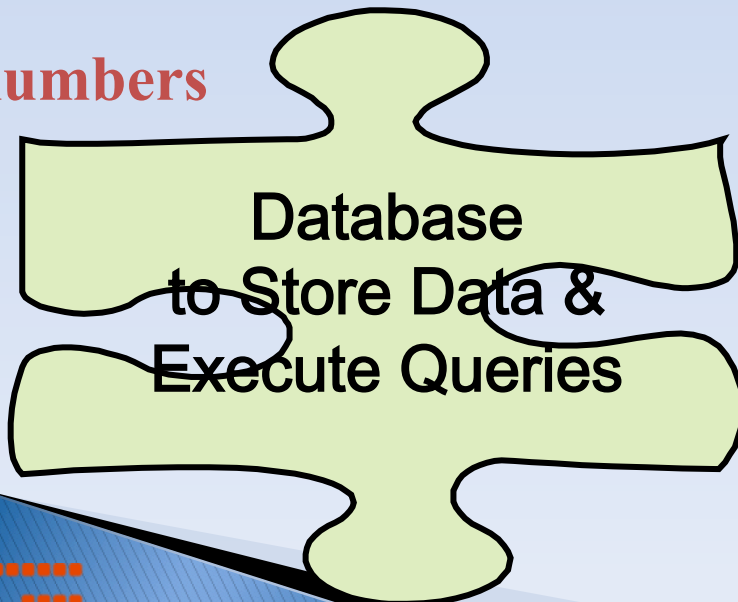
Scientists



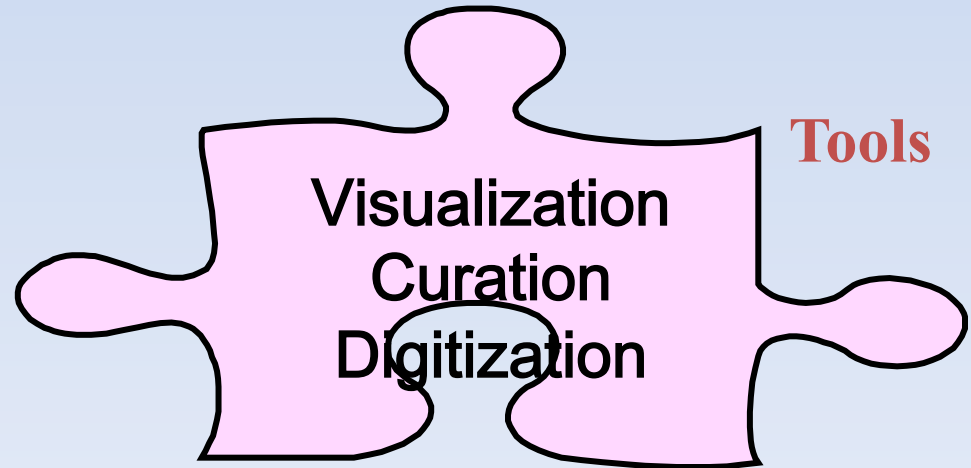
Miners



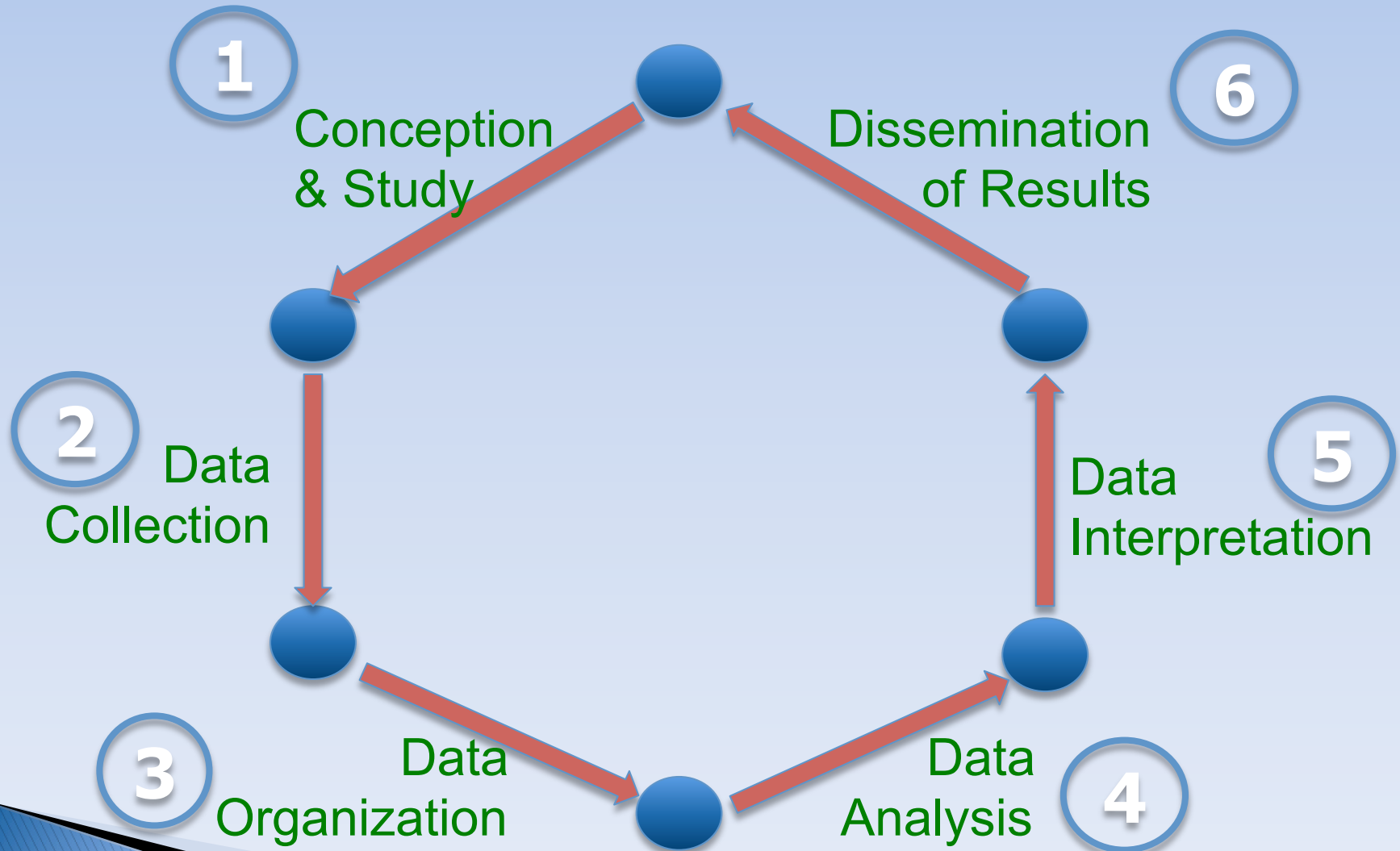
Plumbers



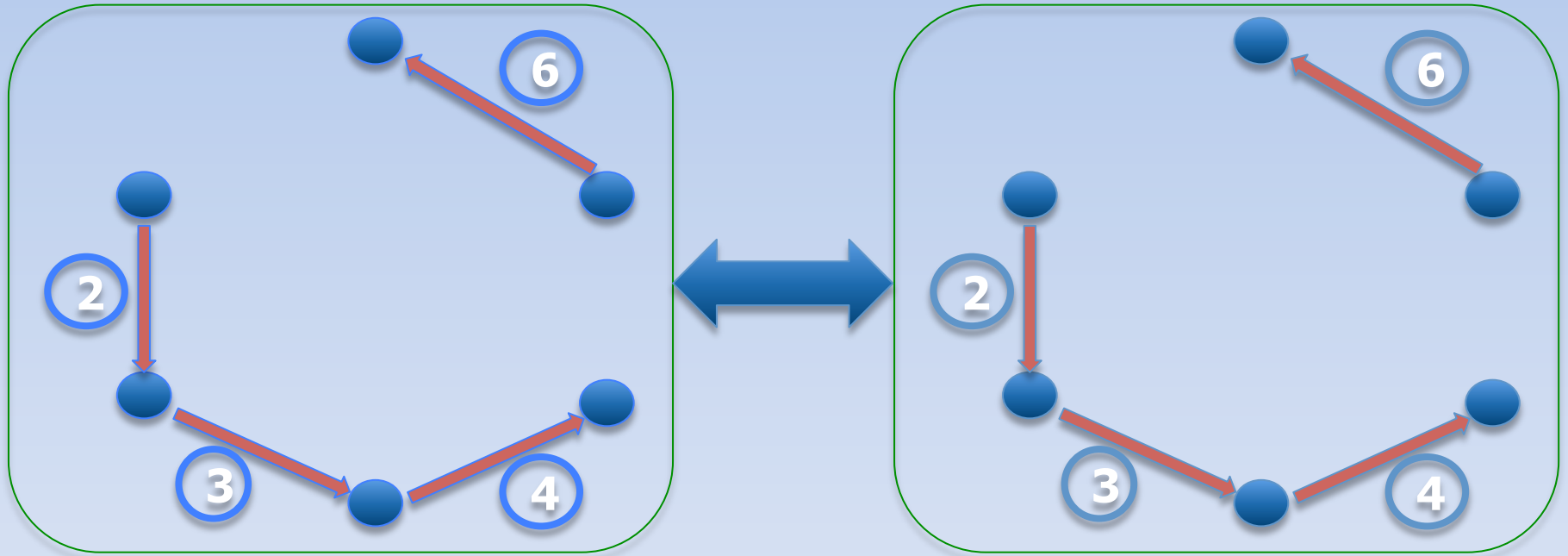
Tools



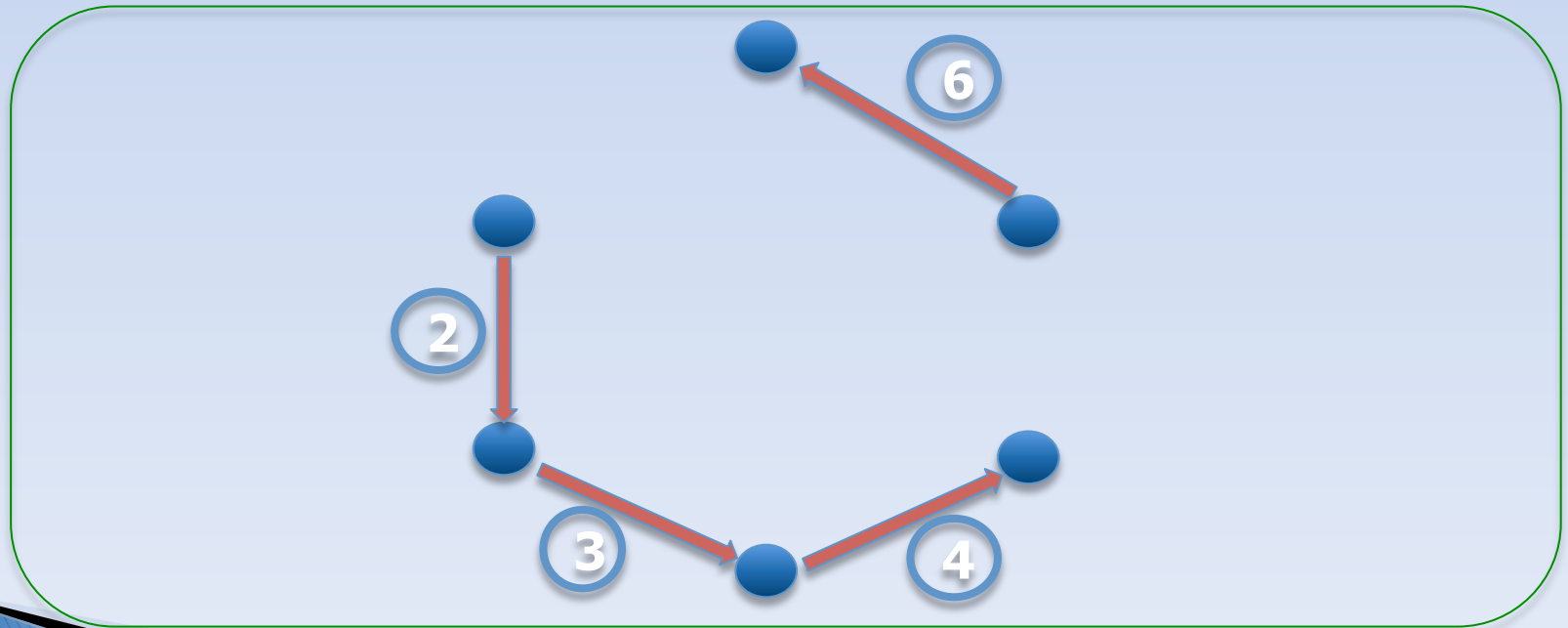
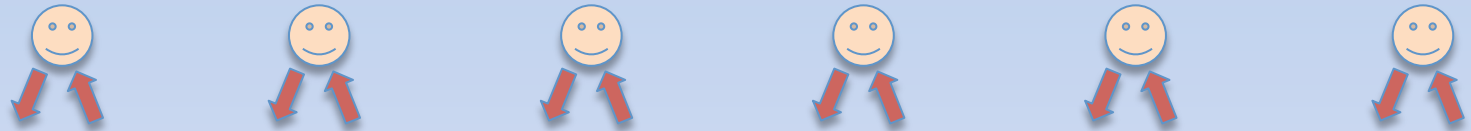
Research Life Cycle



Agenda: Data Interoperability



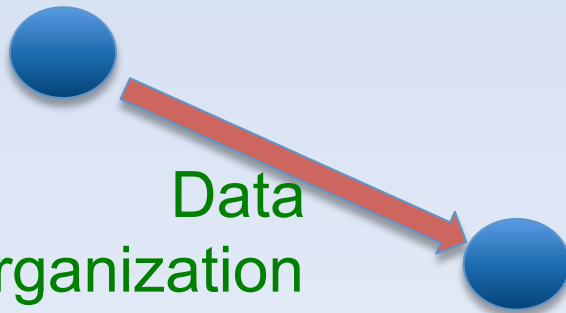
Agenda: Data Infrastructures



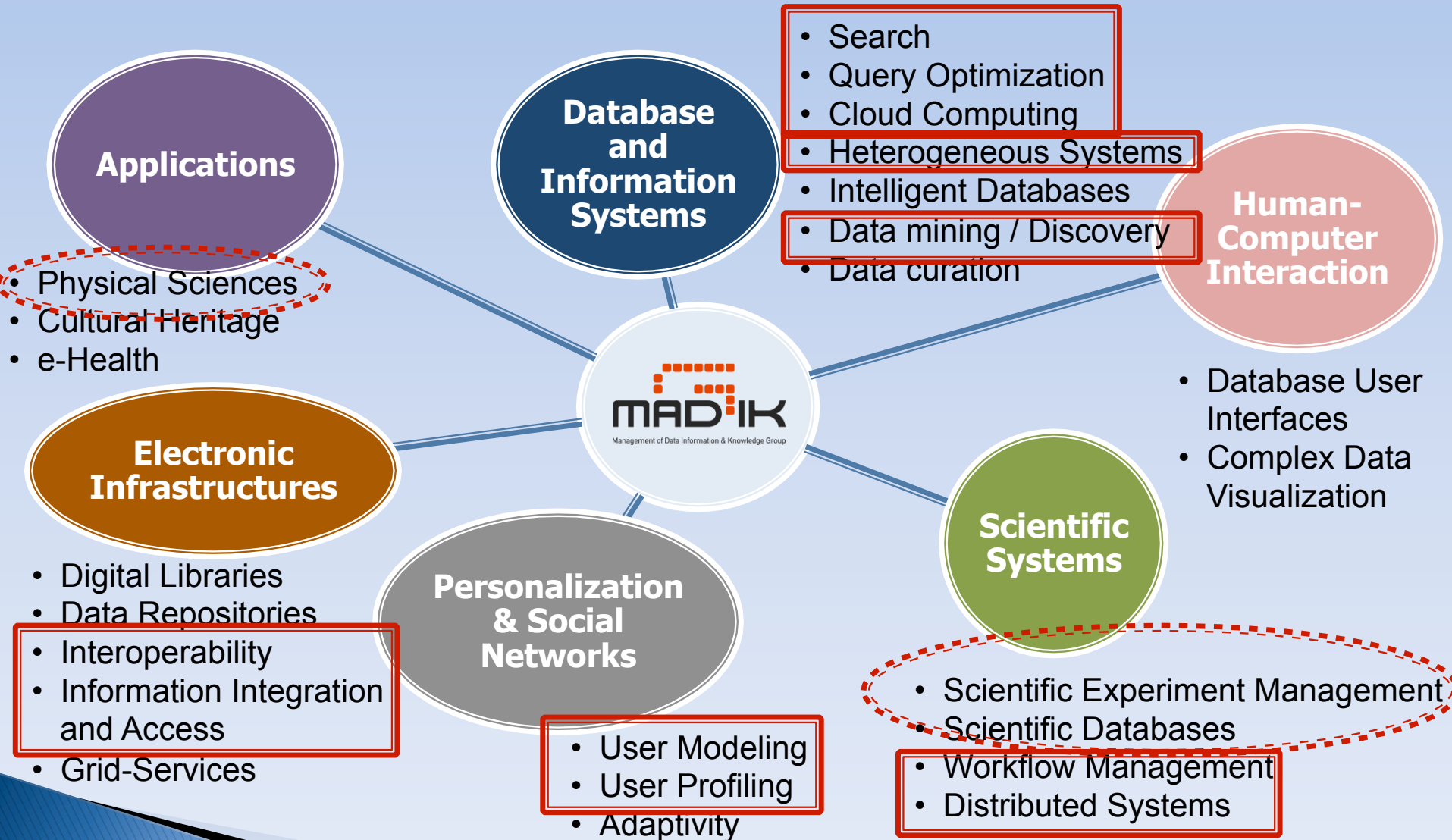
Agenda: Data Management

3

Data
Organization



Research Areas



Spectrum of d-Science Projects

- ▶ ESPAS: Near-earth Space
- ▶ D4Science II: Environmental Monitoring
- ▶ iMarine: Fisheries and Marine Life
- ▶ EarthServer: Earth Sciences
- ▶ Interact: Arctic Environment
- ▶ EFG2: European Film Gateway
- ▶ CHES: Museum Storytelling
- ▶ BMRB: Bio-magnetic Resonance (NMR)
- ▶ OpenAIRE/OpenAIREplus: Publication Monitoring

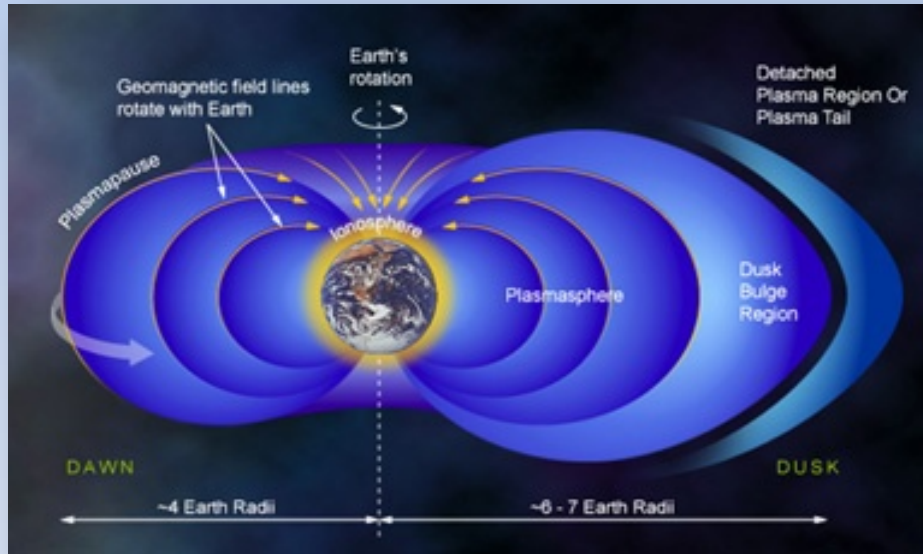


ESPAS

near earth space data infrastructure for e-science

Near-Earth Space

Disparate interrelated scientific communities



- Space weather
- Space climate
- Ionosphere
- Magnetosphere
- Lithosphere-Ionosphere Coupling
- Thermosphere

Industrial users and communities

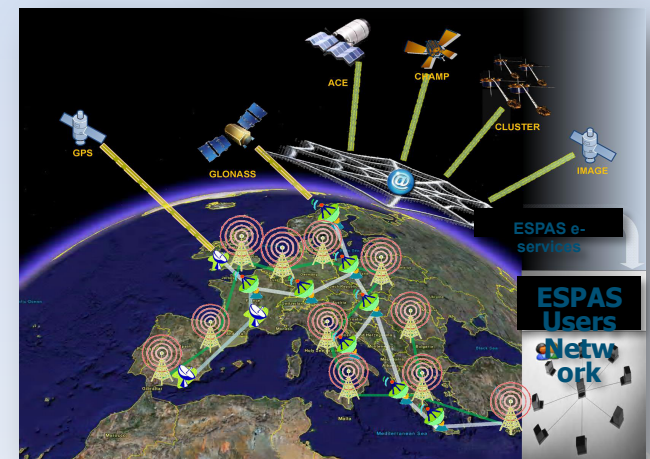


- ▶ Space communications
- ▶ Satellite Operation
- ▶ Navigation and Surveillance

Key issues and challenges

- ▶ Near-Earth space domain is yet to be established
 - Diverse organizations with different (or no) data policies
 - No formal data conceptualization
 - Different scopes, different use by researchers
- ▶ Complexity of data/metadata
 - Rich set of distributed and moving sensors
 - In-situ & active/passive remote sensing

22 partners



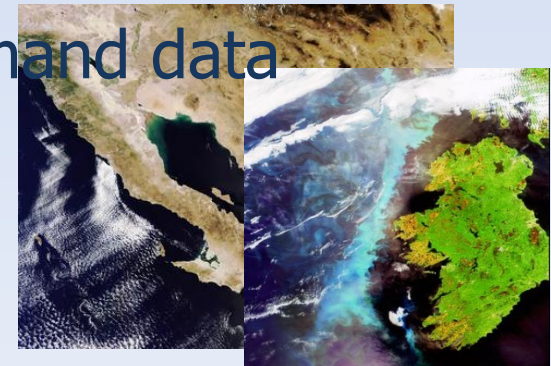
Key issues and challenges

- ▶ Data modeling
 - Diverse communities
- ▶ Data policies
 - Nothing exists
- ▶ Infrastructure to provide access to instrumentation
 - Different policies, different views
- ▶ Infrastructure experiment monitoring
 - Dynamic scientific workflows



D4Science: Environmental Monitoring and Fisheries

- ▶ European Space Agency, FAO, World Fish Center, ...
- ▶ Global environmental issues: marine environment, forest ecosystem, air quality
- ▶ Sensor data analysis, **integration and correlation** of data sources; reasoning, information/knowledge mgnt
- ▶ Large amount of information (→ 1TB), added-value applications and services
- ▶ Seamless workflow definition & on-demand data processing





Fisheries & Marine Life

- ▶ Fishery@FAO, Unesco, ...
- ▶ Worldwide researchers from **many disciplines** (biologists, climatologists, socio-economists, fishery managers, ...)
- ▶ Sustainable environmental management
- ▶ Climate change mitigation, marine biodiversity loss containment, poverty alleviation, disaster risk reduction
- ▶ Use of global fishery ecosystem and aquaculture, e.g., species, aquatic resources, hydrological changes
- ▶ Extreme data diversity





Earth Server **Earth Sciences**

- ▶ European Space Agency
- ▶ Atmospheric, oceanography, geology, and general earth observation communities
- ▶ Open access and ad-hoc analytics on Earth Science
- ▶ RASDAMAN raster data management system
- ▶ Queries on cross-domain data of 20+ Tbytes

Current limitations & gaps

- ▶ **Application-specific software** of limited long term value
- ▶ Absence of **consistent informatics perspective** for data management
- ▶ Discipline-specific solutions: science **re-builds** rather than **re-uses** software
- ▶ No set of **common requirements**

Common Challenges

- ▶ Managing and processing exponentially-growing volumes of non co-located data
- ▶ Dealing with time-sensitive stream data from arrays of sensors and instruments or simulations
- ▶ Significantly reducing data analysis cycles so that researchers can make timely decisions
- ▶ “Data cleaning” is much harder than data loading

Common Research Issues

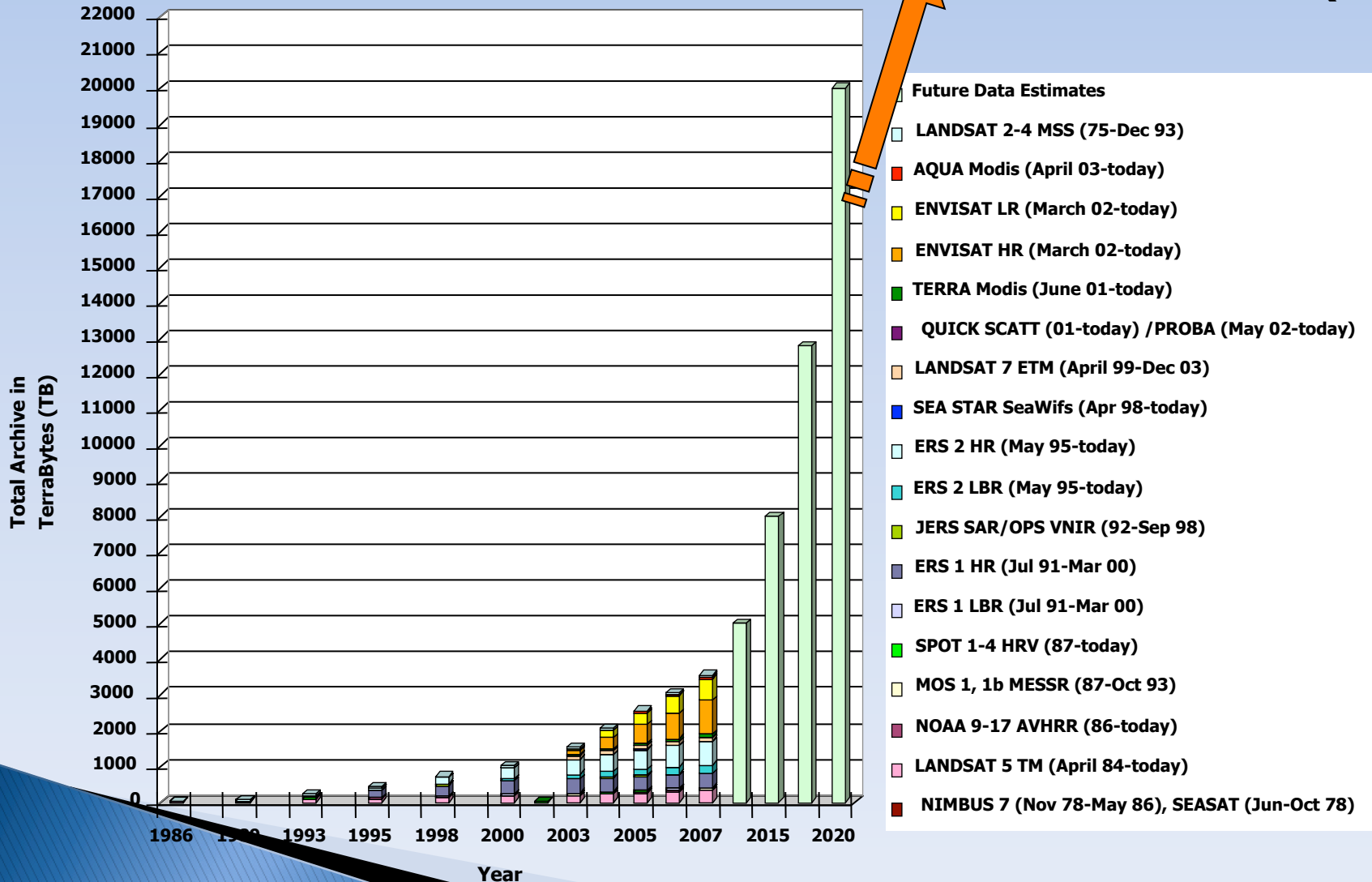
- ▶ Specialized storage and data structures
- ▶ Filtering and fusion
- ▶ Efficient querying and distribution
- ▶ Parallelism
- ▶ Very few predefined query patterns
 - Everything goes....
 - Rapidly extract small subsets of large data sets
 - Geospatial everywhere

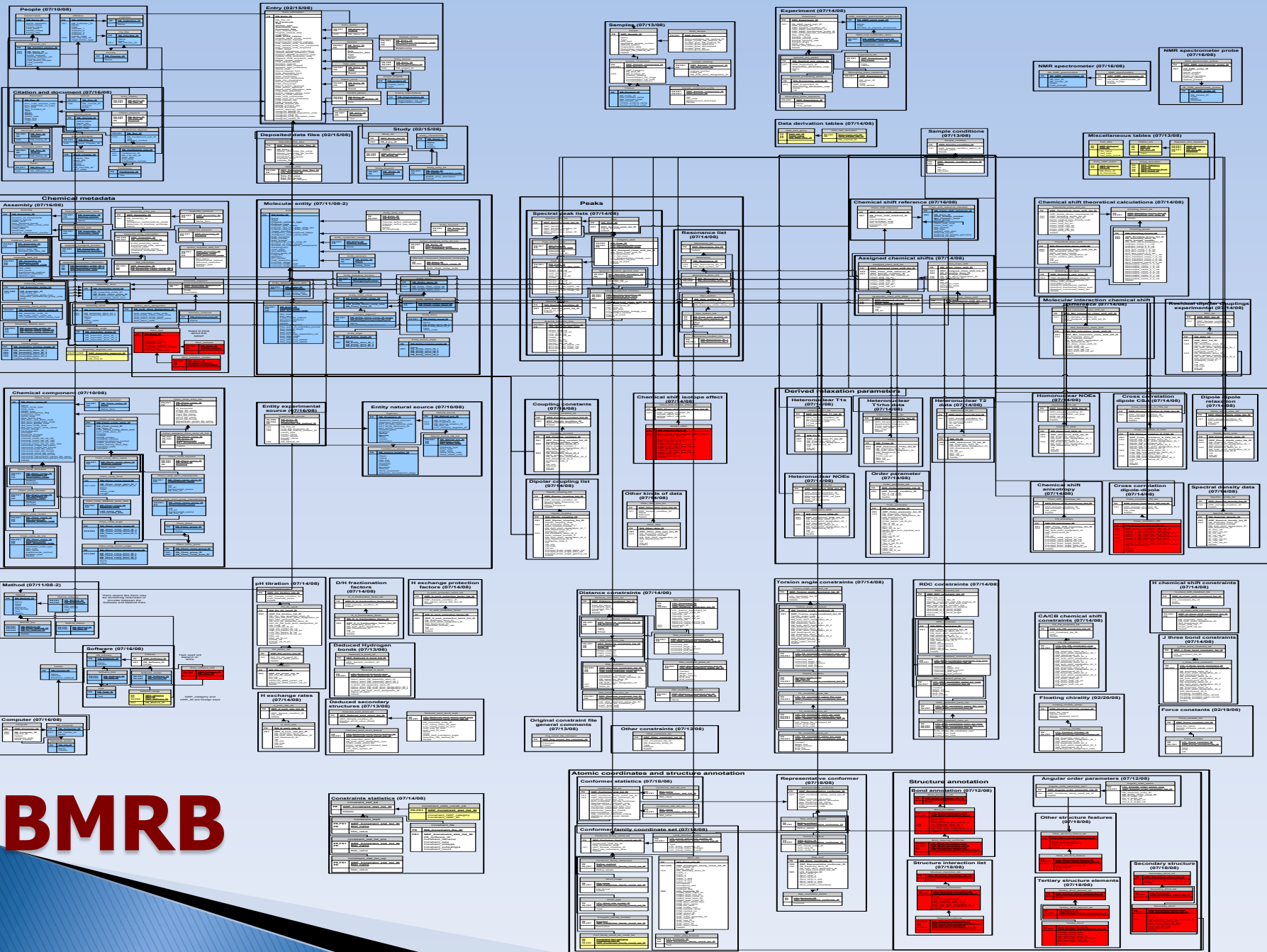
Recent initiatives in informatics

- ▶ Open-source data management and analytics software specifically data-intensive scientific problems
- ▶ Main characteristics
 - scalability up to petabytes and beyond and down to megabytes (consistent interface)
 - built-in support for provenance (lineage), workflows, uncertainty
 - support for "external" data objects
 - optimizing query language
 - use of Open Standards
- ▶ **SciDB** (Stonebraker@MIT et al.), **Rasdaman** (Baumann@Jakobs U.), **MonetDB/SciLense** (Kersten@CWI)

Evolution of Space Info

Evolution of ESA's EO Data Archives 1986-2007+future estimates (<2020)



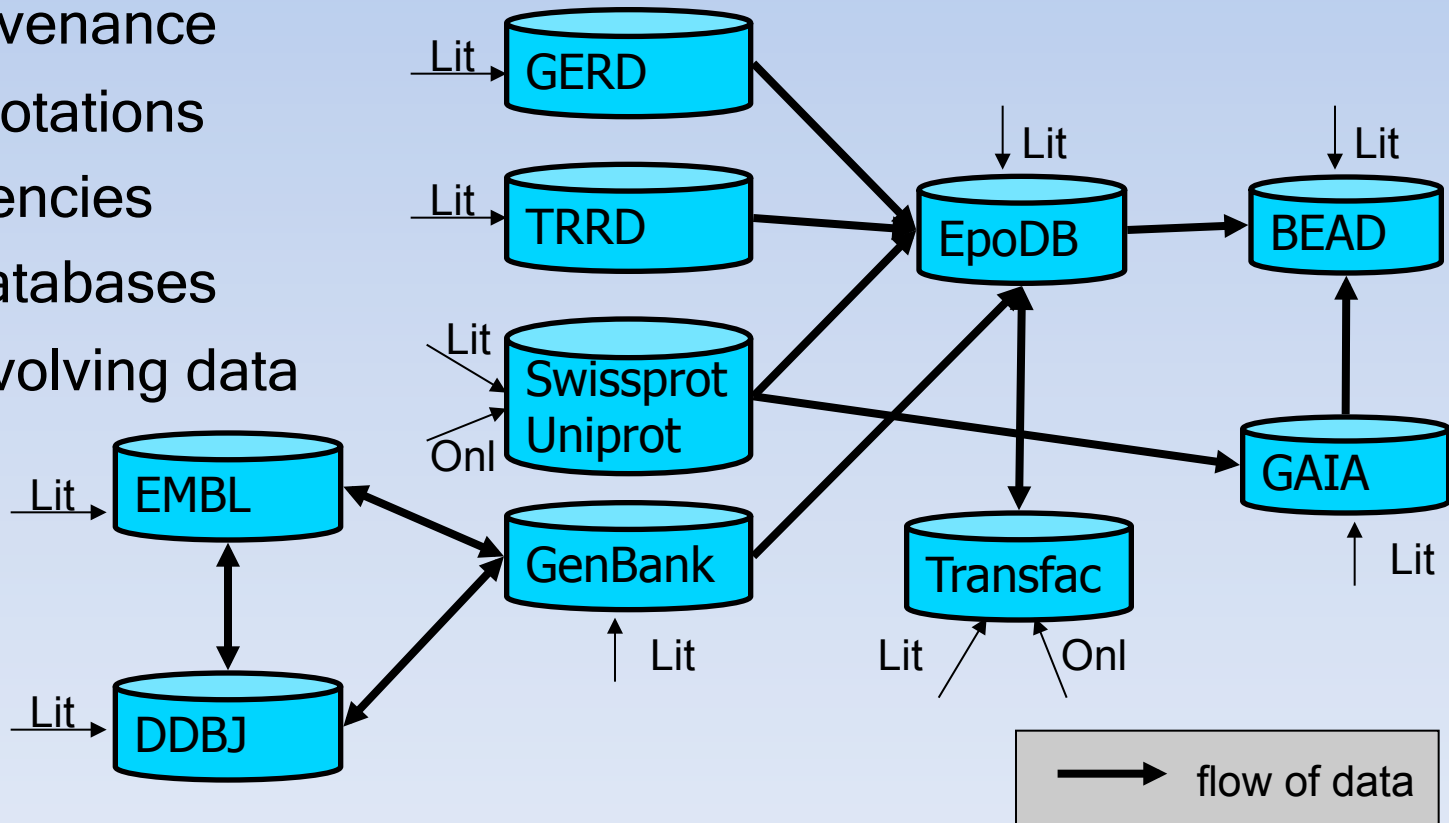


BMRB

Complex Data Interdependencies

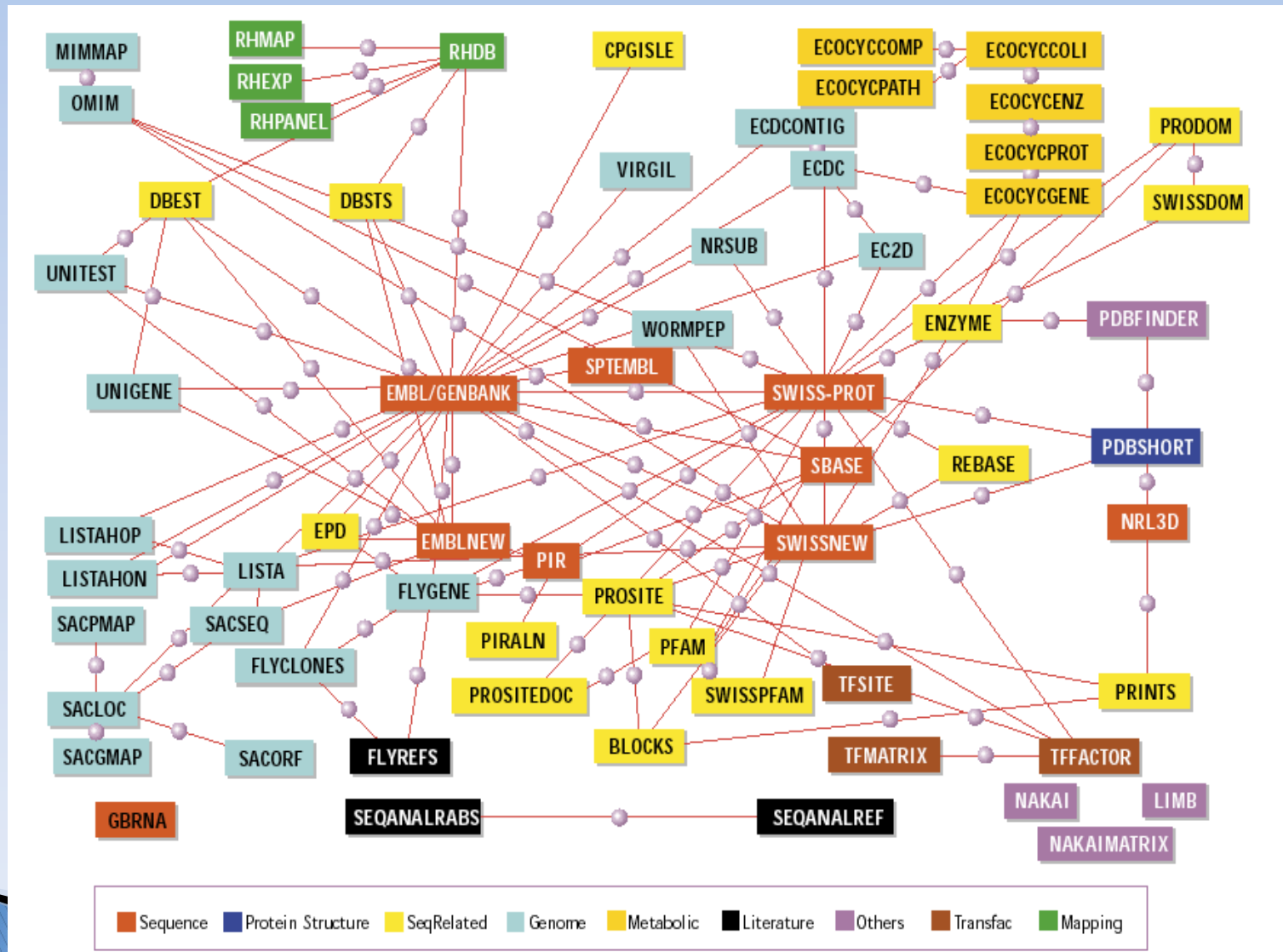
Problems

- Trace data provenance
- Propagate annotations
- Cyclic dependencies
- Cite curated databases
- Synchronize evolving data



Lit: Data from Literature (human intervention)
Onl: Online Data Updates

Complex Data Interdependencies



Scientific Workflow Management

Query Formulation and Results Presentation

Data & Knowledge Management

Personalized Intelligent Applications

Analytics

Knowledge Discovery & Data Mining
Reasoning & Decision Support
Simulation
Visualization

Declarative Language

Learning Inference

Personalization

Distributed Planning & execution Engine

Information Integration

Distributed Search

SRL based probabilistic Model

Domain Models

Configuration Profiles & Models

Service Related Metadata

Data modelling & Knowledge Representation

Dynamic Process...

Distributed

In-Situ Querying

Distributed Processing Middleware for Query Answering and Analytics



Streaming Data



Relational Data



Un/Semi-Structured



Domain Expert



Digital Libraries



Ontologies

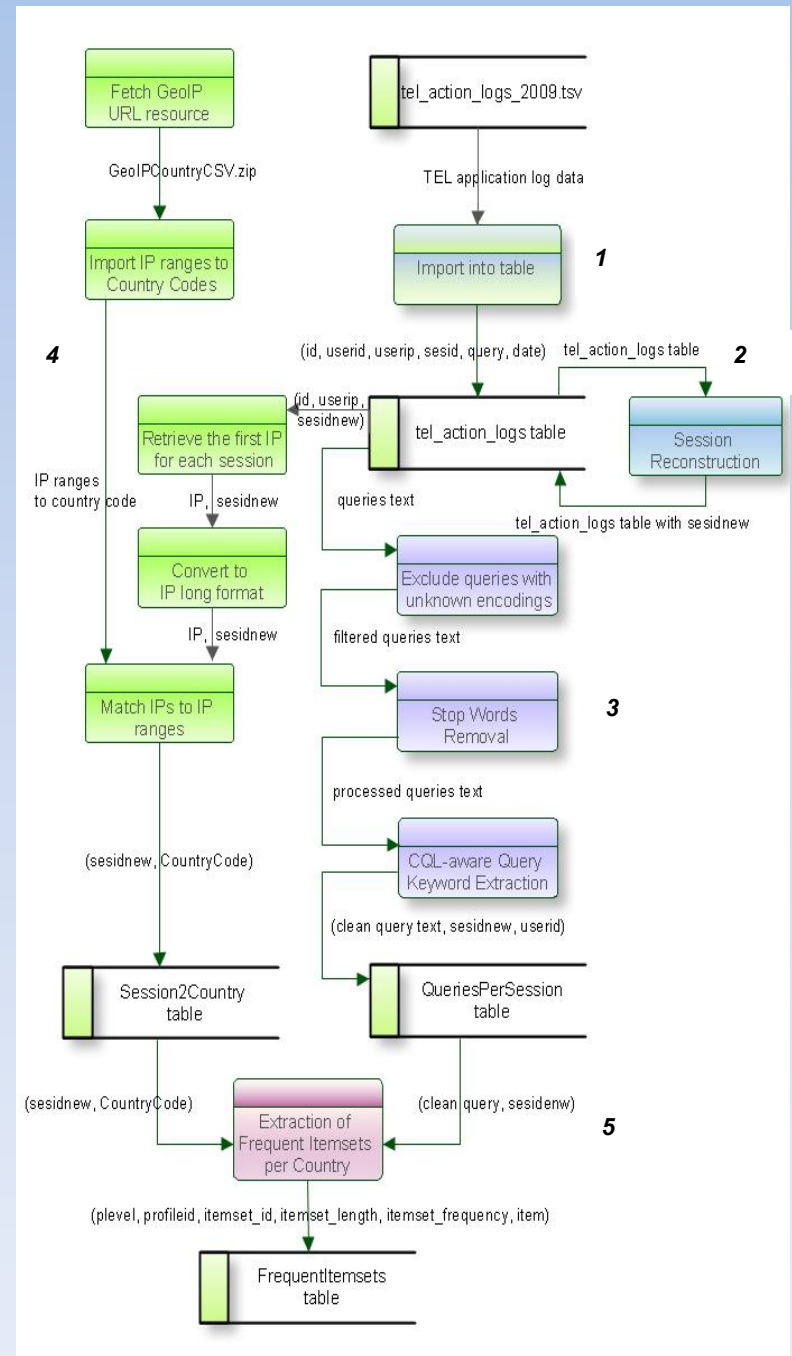
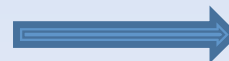
Heterogeneous, distributed data and information sources

Emerging Query Optimization

- ▶ Query: graph of **arbitrary** operators
- ▶ Optimality: response time or completion time and money
- ▶ Environment: **cloud** of hosts (elasticity)

Motivation

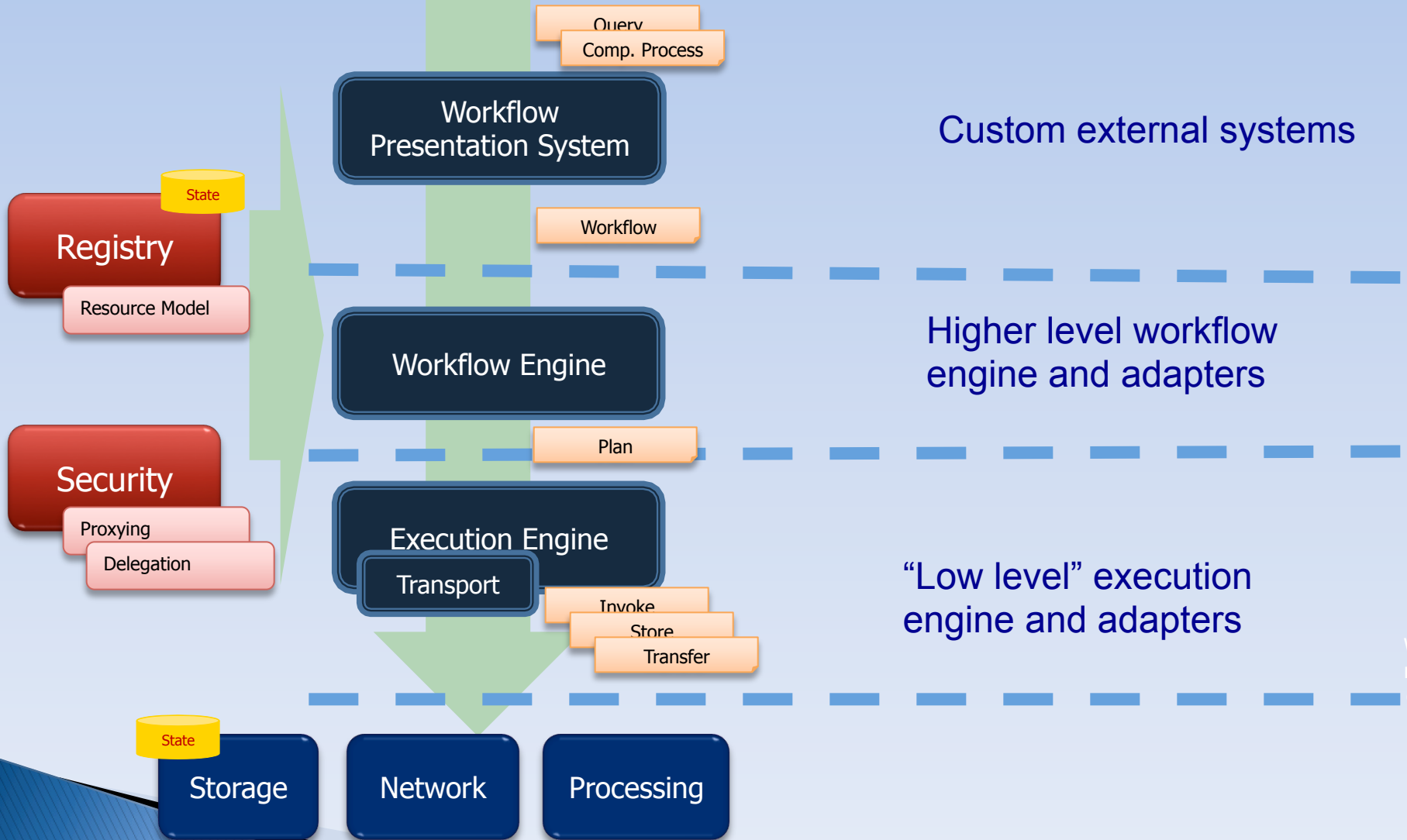
- ▶ Graph of **arbitrary** operators
- ▶ Non-relational data analytics
 - Query log analysis
 - Data mining
 - Simulation model composition
 - ...
- ▶ User behavior analysis for European national libraries
 - One of sixteen flows



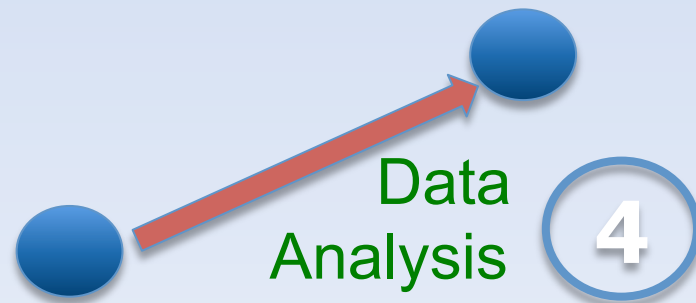
Objectives

- ▶ Execution of complex processing flow
- ▶ Staging of data among different storage providers
- ▶ Expressive and powerful execution plan language
- ▶ Technologies and disciplines abstractions - architecture of loosely coupled layers
- ▶ Unbound extensibility for integration with different environments (Storage, Security, Execution, ...) – plugability
- ▶ Execution of ... everything (SOAP/REST & WSRF web services, local executables and scripts, POJOs, native components) and in any manner (in-process, intra-process, intra-node)
- ▶ Align with cloud computing principles & data stream pipelining

Architecture Overview



Agenda: Data Management



Analysis and Databases

- ▶ Much statistical analysis deals with
 - Creating uniform samples
 - Data filtering
 - Assembling relevant subsets
 - Estimating completeness
 - Censoring bad data
 - Counting and building histograms
 - Generating Monte-Carlo subsets
 - Likelihood calculations
 - Hypothesis testing
- ▶ Traditionally these are performed on files
- ▶ Most tasks are much better done inside a database
 - Indices, parallel data search & analysis

AITION Data Mining Process

1. Qualitative dependency analysis: Learning the structure
2. Quantitative analysis: Learning the parameters (CPD)

Make inferences based on the model (Graph + CPD) using different known / unknown schemes

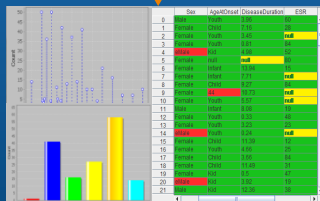
Raw Data

ID	Gender	Age	DisDur	CHAQ	ESG	ESG
1	Male	65	12	10	10	10
2	Female	72	15	12	12	12
3	Male	58	10	8	8	8
4	Female	68	18	15	15	15
5	Male	75	20	18	18	18
6	Female	60	14	11	11	11
7	Male	70	16	13	13	13
8	Female	62	13	10	10	10
9	Male	73	19	16	16	16
10	Female	64	15	12	12	12
11	Male	71	17	14	14	14
12	Female	66	14	11	11	11
13	Male	74	21	19	19	19
14	Female	61	13	10	10	10
15	Male	76	22	20	20	20
16	Female	63	14	11	11	11
17	Male	72	18	15	15	15
18	Female	65	15	12	12	12
19	Male	77	23	21	21	21
20	Female	67	16	13	13	13
21	Male	78	24	22	22	22
22	Female	69	17	14	14	14
23	Male	79	25	23	23	23
24	Female	70	18	15	15	15
25	Male	80	26	24	24	24
26	Female	71	19	16	16	16
27	Male	81	27	25	25	25
28	Female	72	20	17	17	17
29	Male	82	28	26	26	26
30	Female	73	21	18	18	18
31	Male	83	29	27	27	27
32	Female	74	22	19	19	19
33	Male	84	30	28	28	28
34	Female	75	23	20	20	20
35	Male	85	31	29	29	29
36	Female	76	24	21	21	21
37	Male	86	32	30	30	30
38	Female	77	25	22	22	22
39	Male	87	33	31	31	31
40	Female	78	26	23	23	23
41	Male	88	34	32	32	32
42	Female	79	27	24	24	24
43	Male	89	35	33	33	33
44	Female	80	28	25	25	25
45	Male	90	36	34	34	34
46	Female	81	29	26	26	26
47	Male	91	37	35	35	35
48	Female	82	30	27	27	27
49	Male	92	38	36	36	36
50	Female	83	31	28	28	28

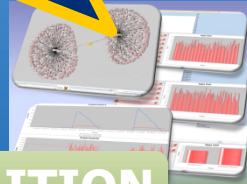


DCV

Data Curation & Validation

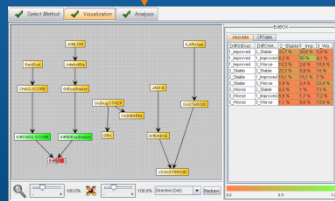


Curated Data

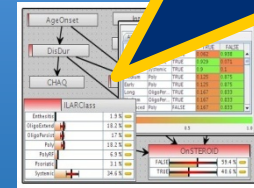


AITION DESK

Model Building

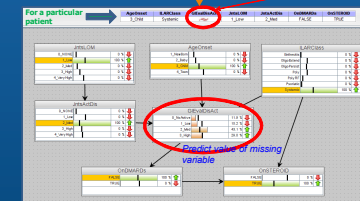


Model



AITION VIZ

Model Exploration & Reasoning

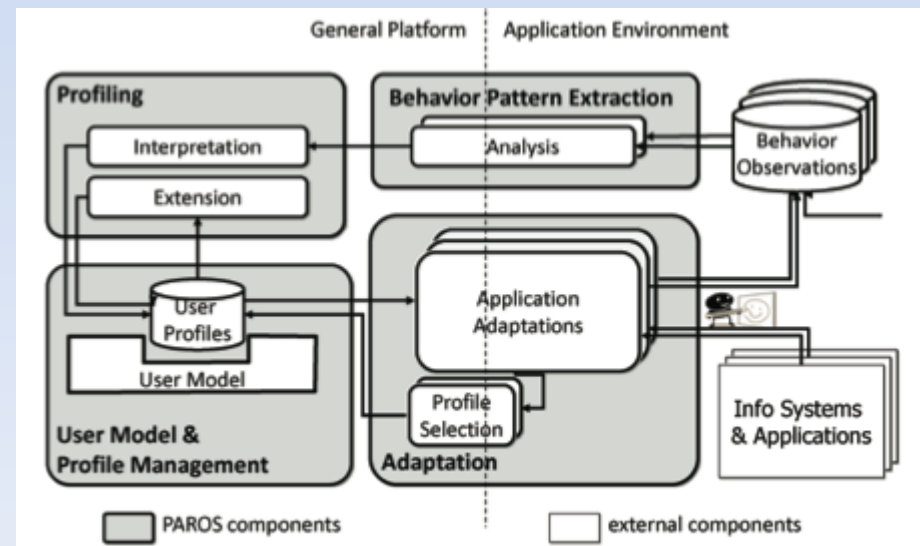


Decisions & Simulation

One way & static ...

PAROS Personalization Framework

- ▶ Provide:
 - Modelling user preferences & attitudes
 - Personalized querying, context & services
 - Graph mining for usage pattern analysis



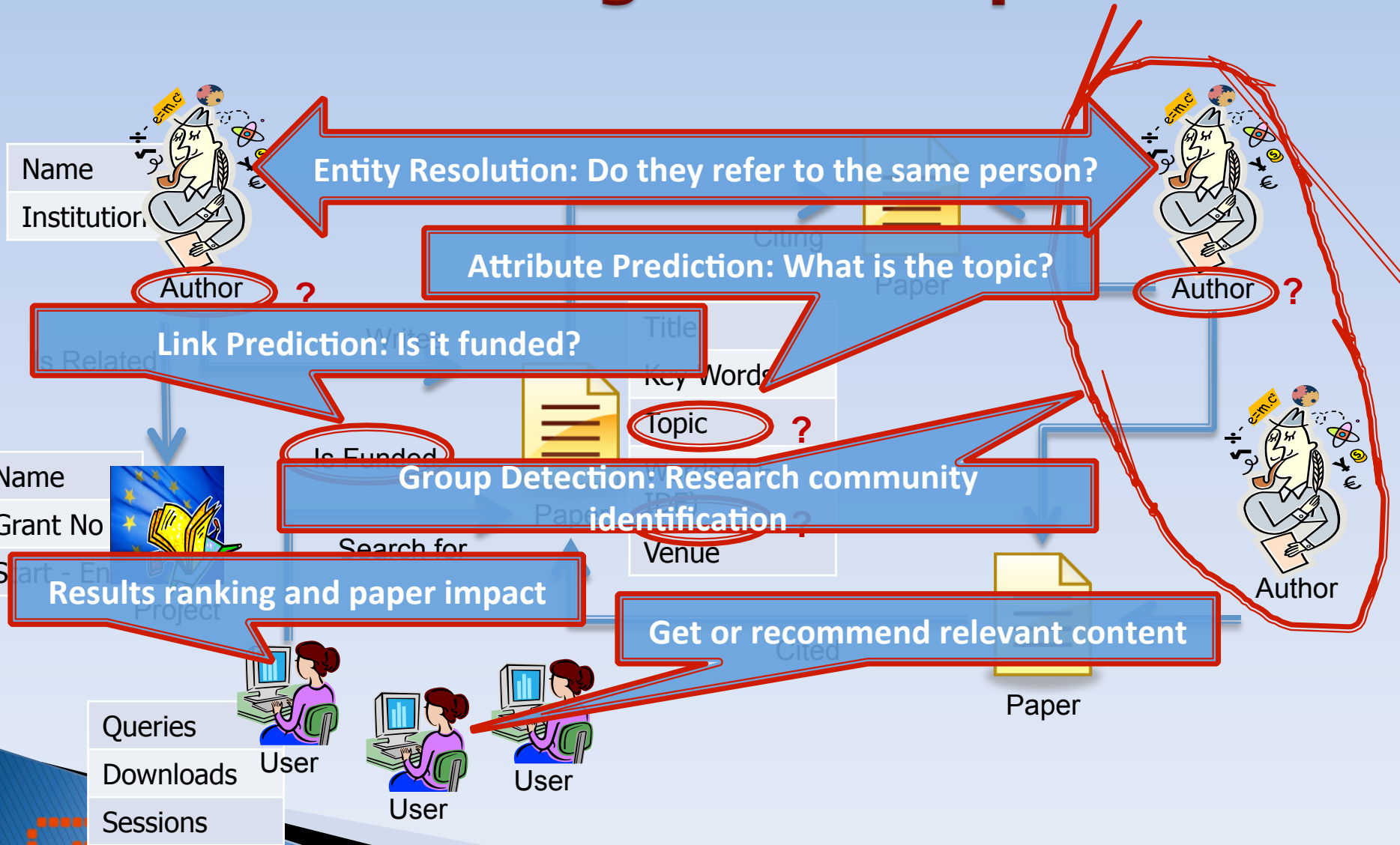
Phase II: Holistic approach...

- ▶ Evolutionary information processing and knowledge discovery framework able to provide highly accurate predictive and simulation models combining:
 - **Bottom-up data driven process:** analyzing huge volume of high distributed, federated, multi-type, vertically integrated data and streams,
 - **Top-down model driven process:** incorporating domain knowledge represented as relational & semantic theoretic models and First-Order Logic (FOL) rules & constraints.

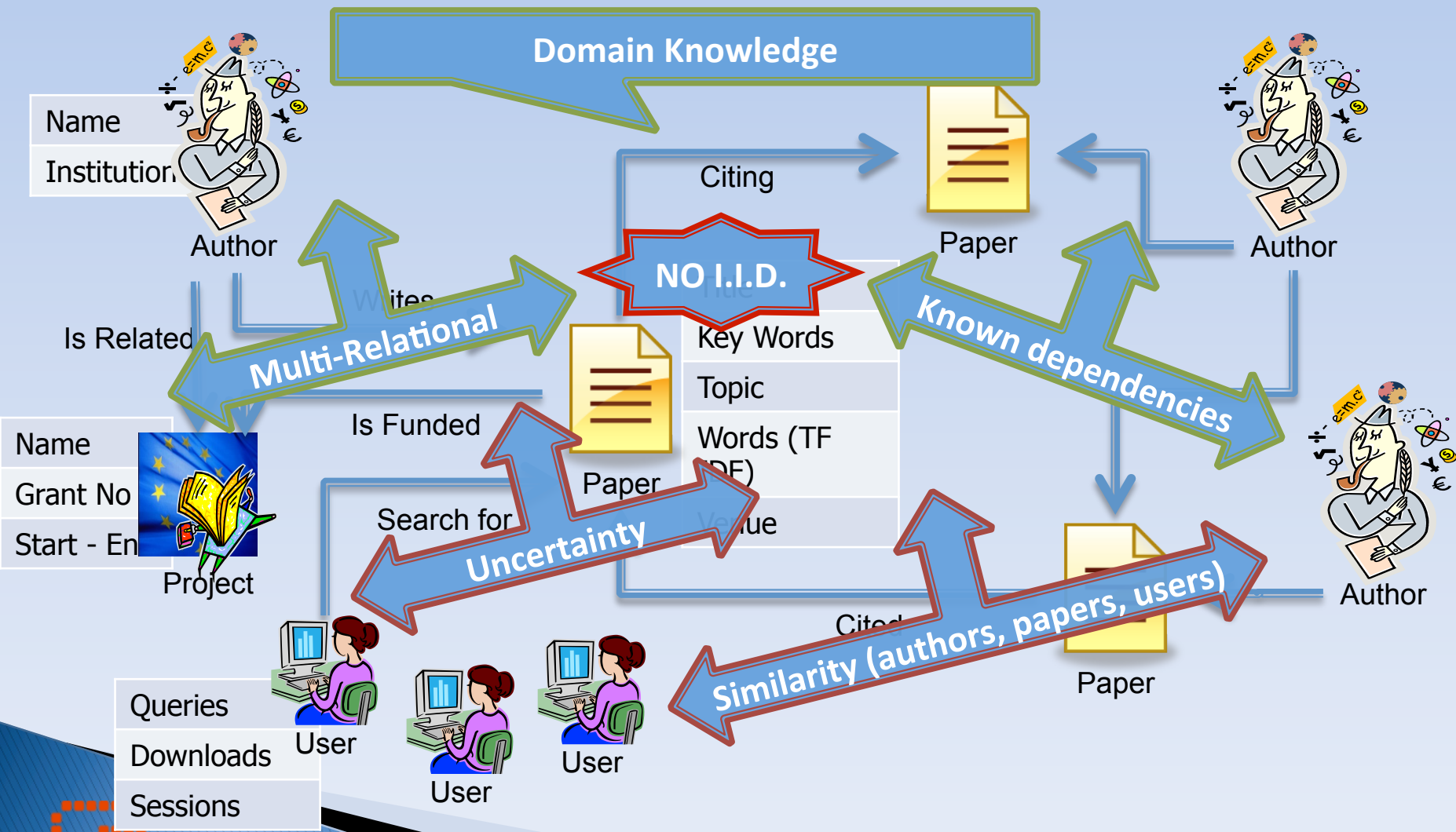
Basic Problem....

- ▶ Traditional Machine Learning (ML) approach based on single relation (propositional) and consider i.i.d. data
- ▶ Traditional Data Mining (DM) cannot handle uncertainty in attributes, entities, relations
- ▶ But:
 - Uncertainty is everywhere:
 - No complete knowledge
 - Noise/outliers in the data due to measurement procedures, human involvement in the loop, multi-cite, etc.
 - Real world has objects, properties, known relations
 - Instances are not always independent

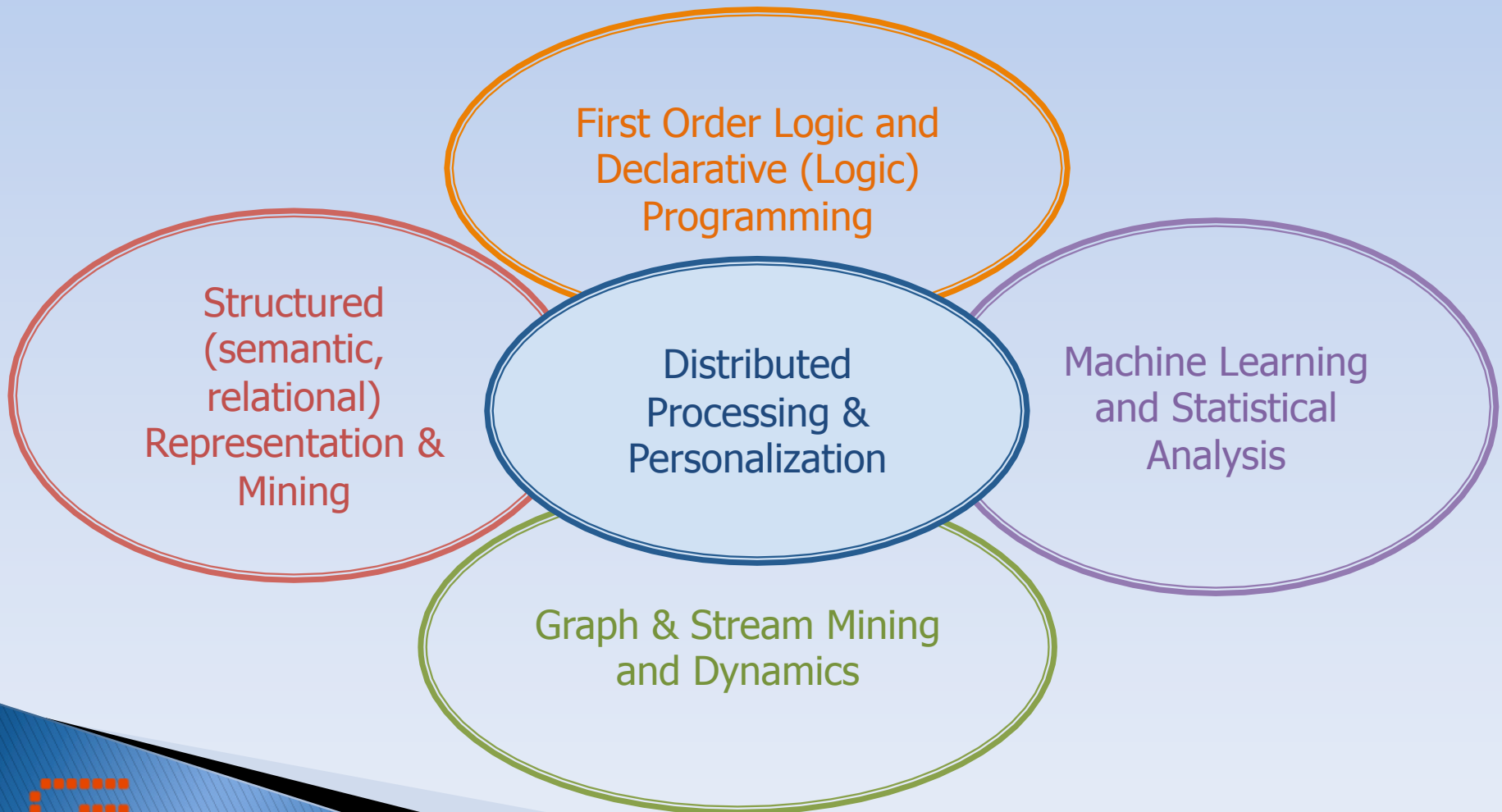
Hollistic Intelligence in OpenAIRE



Intelligence in Open AIRE



Holistic approach...



ICT Results

[Features & Publications](#)[Press Desk](#)[Investors Room](#)[Help & Links](#)[ICT Results for You](#)**High-impact
ICT research**[Learn more](#)**Features**[Print friendly](#)[Send this feature](#)[Send Feedback](#)**Software for solving life-threatening medical puzzles**

New software is under development that doctors hope will help them identify brain tumours in children that will grow aggressively.

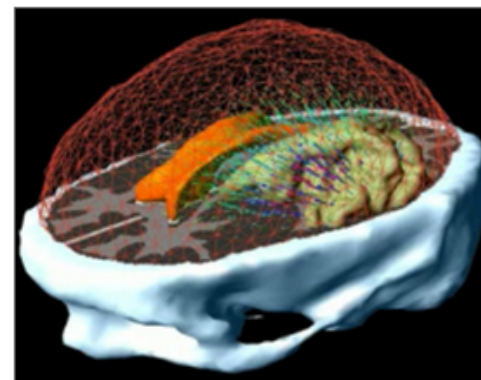
Some **brain tumours** in children remain benign and doctors choose not to operate. But a small percentage of those will suddenly start to grow aggressively.

Doctors have not identified what triggers that aggressive tumour growth, despite the vast array of data they hold on their child patients – demographic, environmental, genetic and clinical data, as well as images such as MRI and CAT scans of the developing tumours.

But a new software tool called **AITION** can integrate all the medical data from a tumour patient and then analyse it to calculate the probable factors that are stimulating tumour development, combining up to 30 correlated variables. AITION provides an overview of the causal relationship across all factors.

Graphical network of causal relationships

AITION's conclusions are displayed as a 'knowledge model', a graphical network of medical factors with links that represent the correlations between them. Strongly interdependent concepts are directly connected, loosely dependent concepts are not connected at all. The patient's doctors can play around with the knowledge model. They can improve the model by adding information they know to be true about the patient. They can use the model to test the likely effects of different types of medication, surgery



Agenda: Data Management



Dissemination
of Results



Open Access

Open Access brings the innovation cycle to a new open world, contributing to the **Digital Agenda for Europe**

"Scientific data has the power to transform our lives for the better – it is too valuable to be locked away."

Neelie Kroes,

Vice-President of the EC, responsible for the Digital Agenda

Open Science

- **Open Science** enables researchers to conduct **efficiently** and **effectively** their research activities
- **Open Science** facilitate multidisciplinary **collaboration** and data **sharing** entailing **Open Access** to
 - research data
 - data services
 - tools
 - analyses
 - methods

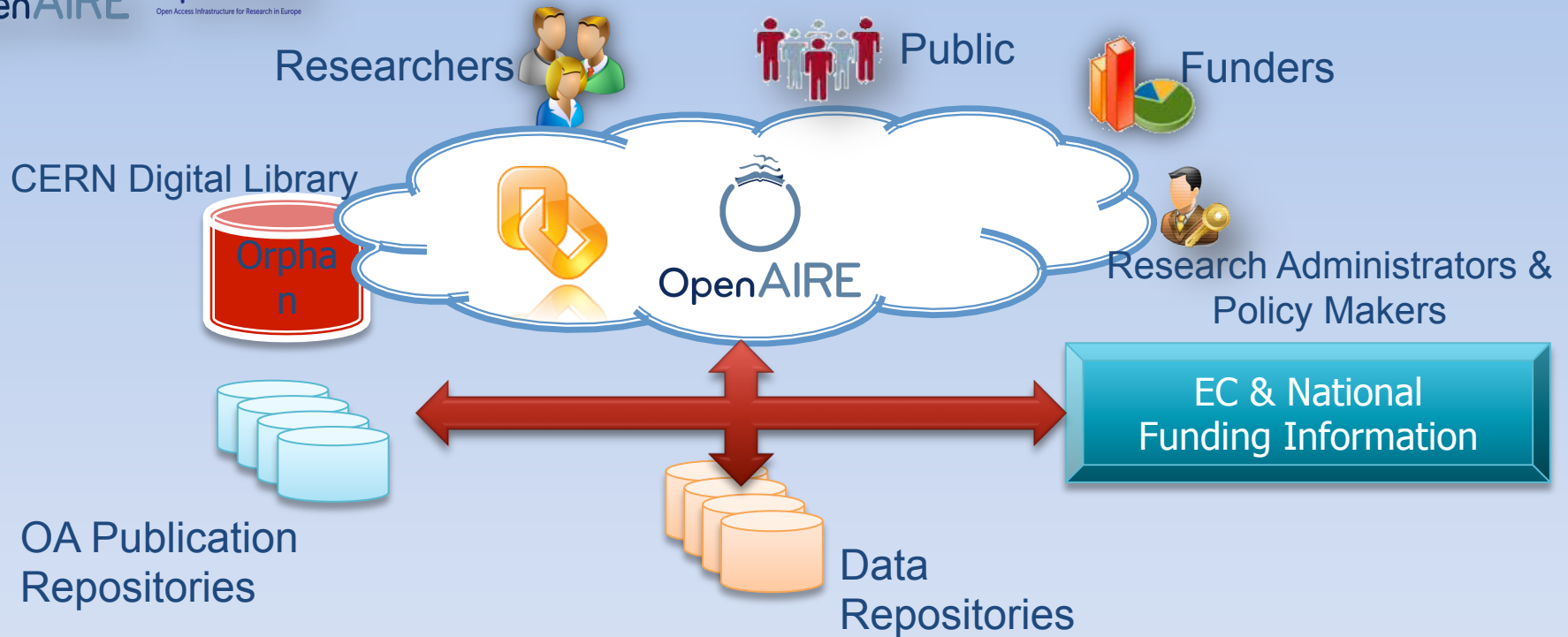
Open Access & Scientific Freedoms

- **Open Access** improves **information access, use and reuse** - easier & more reliable
- **Open Access** improves **free info movement** in **Research** and **Public Sector**, benefiting
 - researchers, students, citizens
 - funding officers and industry

Investing in Open Data Infrastructures

- **Relax** and **bridge** geographic, organizational, disciplinary boundaries
- A domain with a **global** horizon

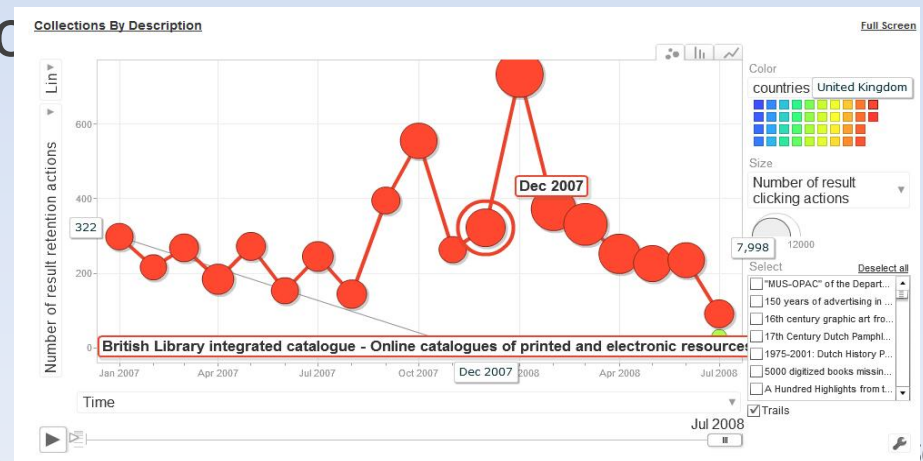
Unifying European Research



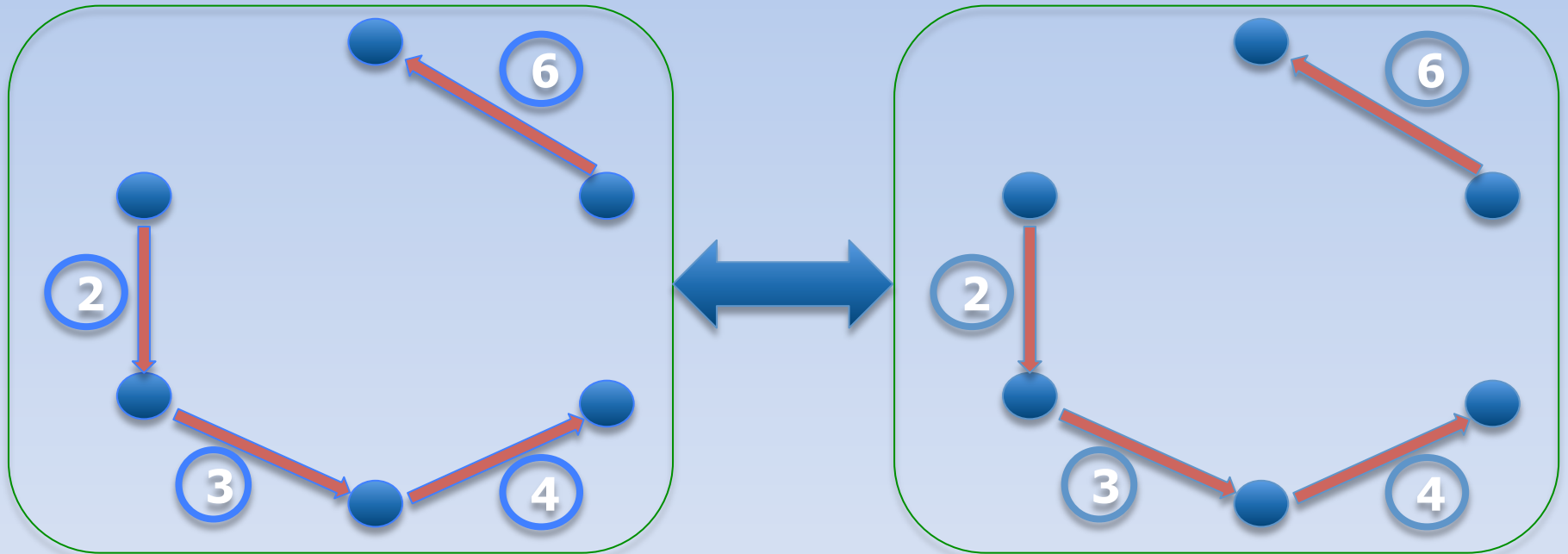
- Links OA repositories with FP7 / ERC /national funding
- Links publications with research data
- Provides repository for "orphan" publications & data
- Calculates statistics on OA/FP7/ERC /... deposition

Usage data

- ▶ Analysis to measure **research impact**
- ▶ Usage pattern analysis tools for content linking
- ▶ Challenges, research direction
 - Big volumes of data to process
 - Definition of metrics
 - Researcher behavior and so
 - Organizational issues



Agenda: Data Interoperability



Interoperability

DEFINITION

- ▶ *“The ability of two or more systems or components to **exchange** information and to **use** the information that has been exchanged” (IEEE definition)*
- ▶ *“A property referring to the ability of diverse systems and organizations to work together (inter-operate)” (Wikipedia)*
- ▶ *“The capability to communicate, execute programs, or transfer data among various functional units in a manner that requires minimal knowledge of the unique characteristics of those units” (ISO/IEC 2382-2001 Information Technology Vocabulary, Fundamental Terms)*

Data interoperability

- ▶ The ability to exchange data between different systems and use them
- ▶ Indicative areas: information representation - profiling, discovery, security
 - Multidisciplinary information – Scientific heterogeneity
 - Inconsistency of intended use of information

Technical Challenges

- ▶ **Heterogeneity** of the exchanged information objects
 - Syntactic heterogeneity
 - Structural heterogeneity
 - Semantic heterogeneity
- ▶ **Inconsistency of the intended use of the information object**
 - by the producer entity and the intended exploitation of this object by the consumer entity
 - the exchanged information must be complemented with some **“descriptive”** information, such as contextual, provenance, quality, security, privacy, etc. information
 - The descriptive information should be modeled by **purpose-oriented descriptive data models / metadata models.**

Organisational and policy challenges

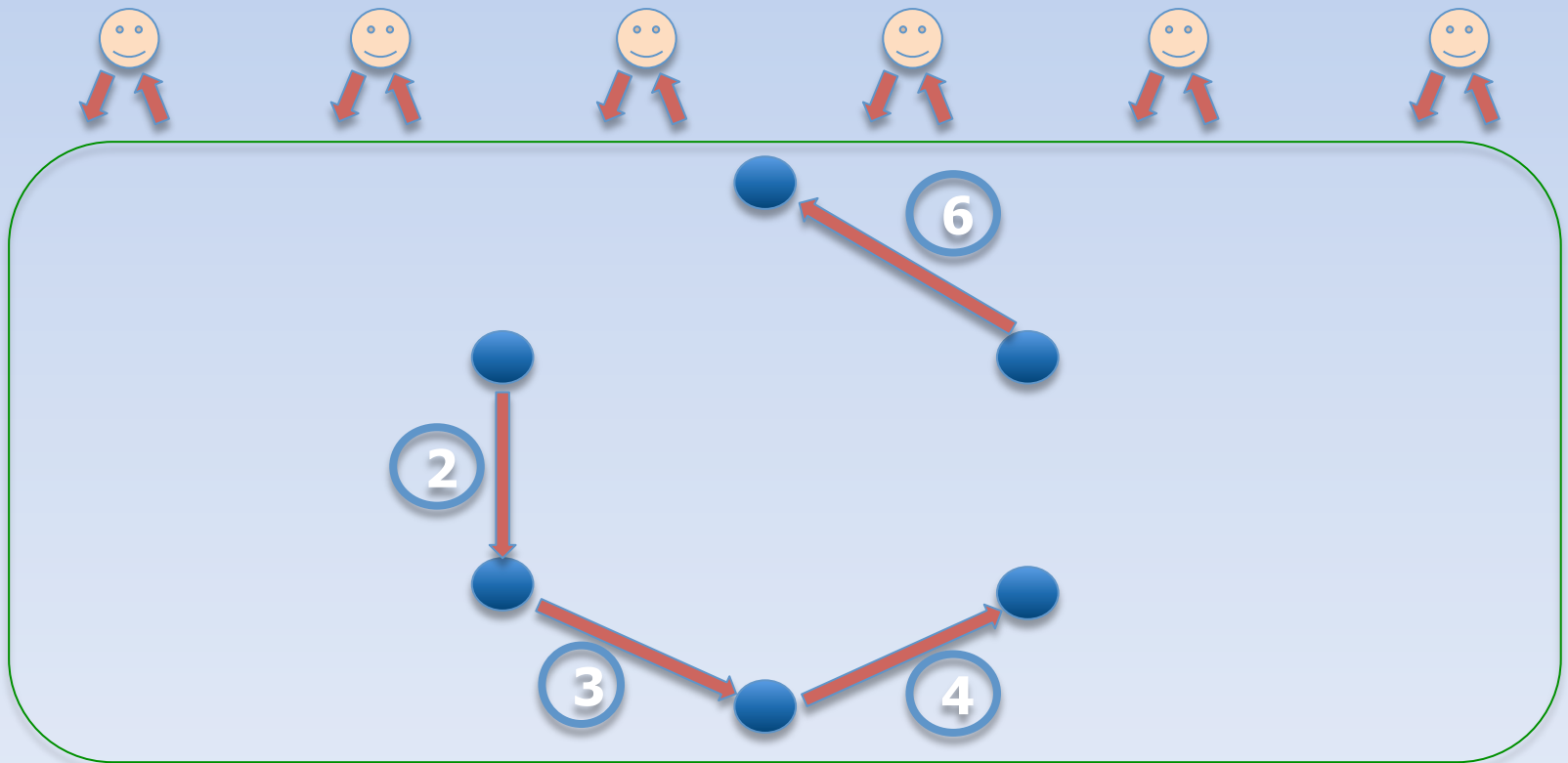
- ▶ Interoperability is not only about technical challenges
- ▶ Policy and organisational challenges exist:
 - Interoperability entails (global / con-federated) governance and trust
 - Interoperability entails certification and standardisation
 - Interoperability entails sustained funding and support!

Solution

- ▶ Quick and dirty: Mediation
- ▶ Slow and clean: Standardization

- ▶ Top-down, strong-armed will never be final
- ▶ Bottom-up, grass-roots will always exist
- ▶ Never-ending issue

Agenda: Data Infrastructures



What is “e-Infrastructure”?

“an environment where research resources (**hardware, software and content**) can be readily shared and accessed wherever this is necessary to promote better and more effective research”.

Networks, grids and middleware, comp. resources, exp workbenches, data repos, tools and instruments and the op support that enable global virtual research collabs.

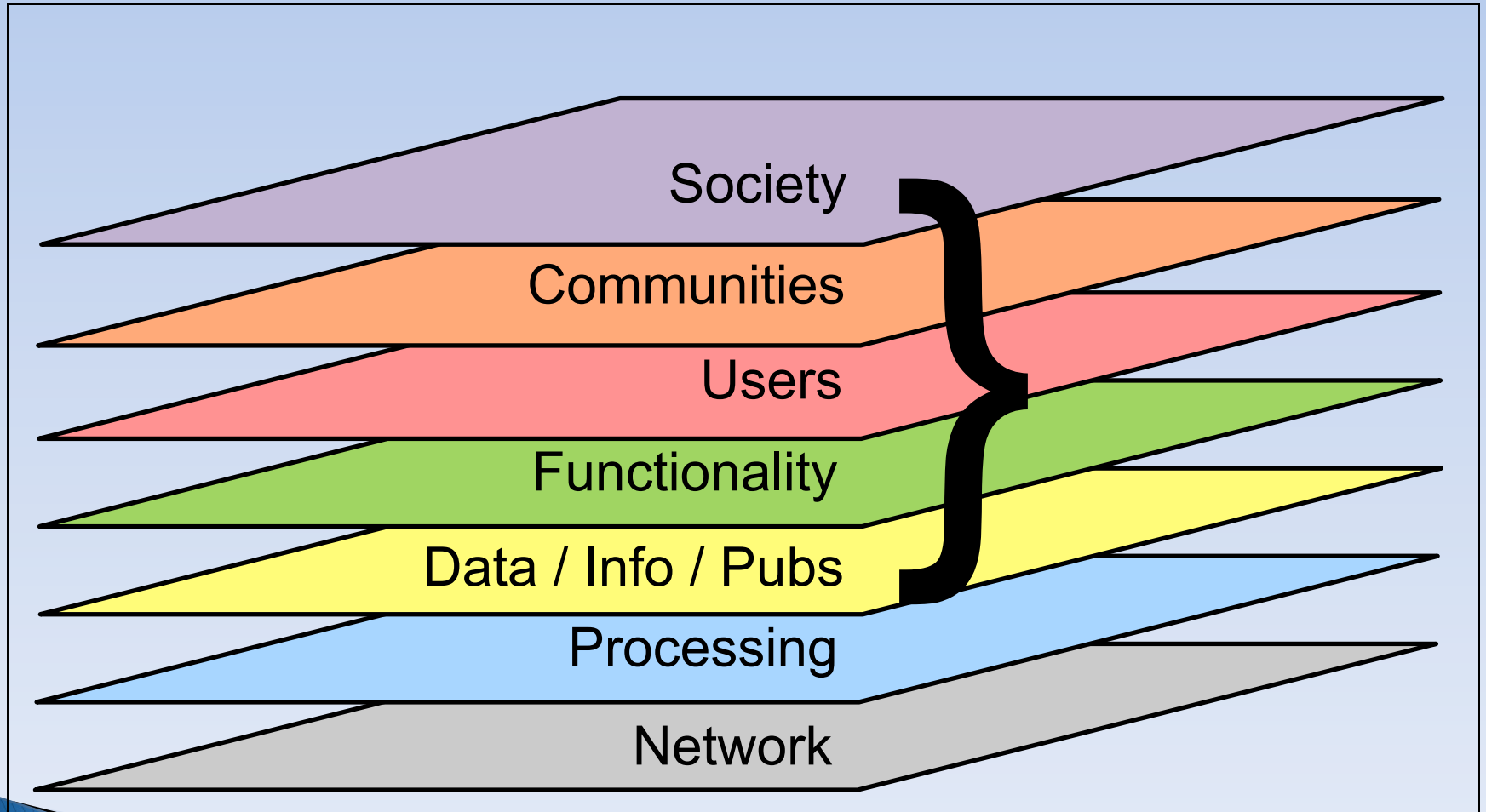


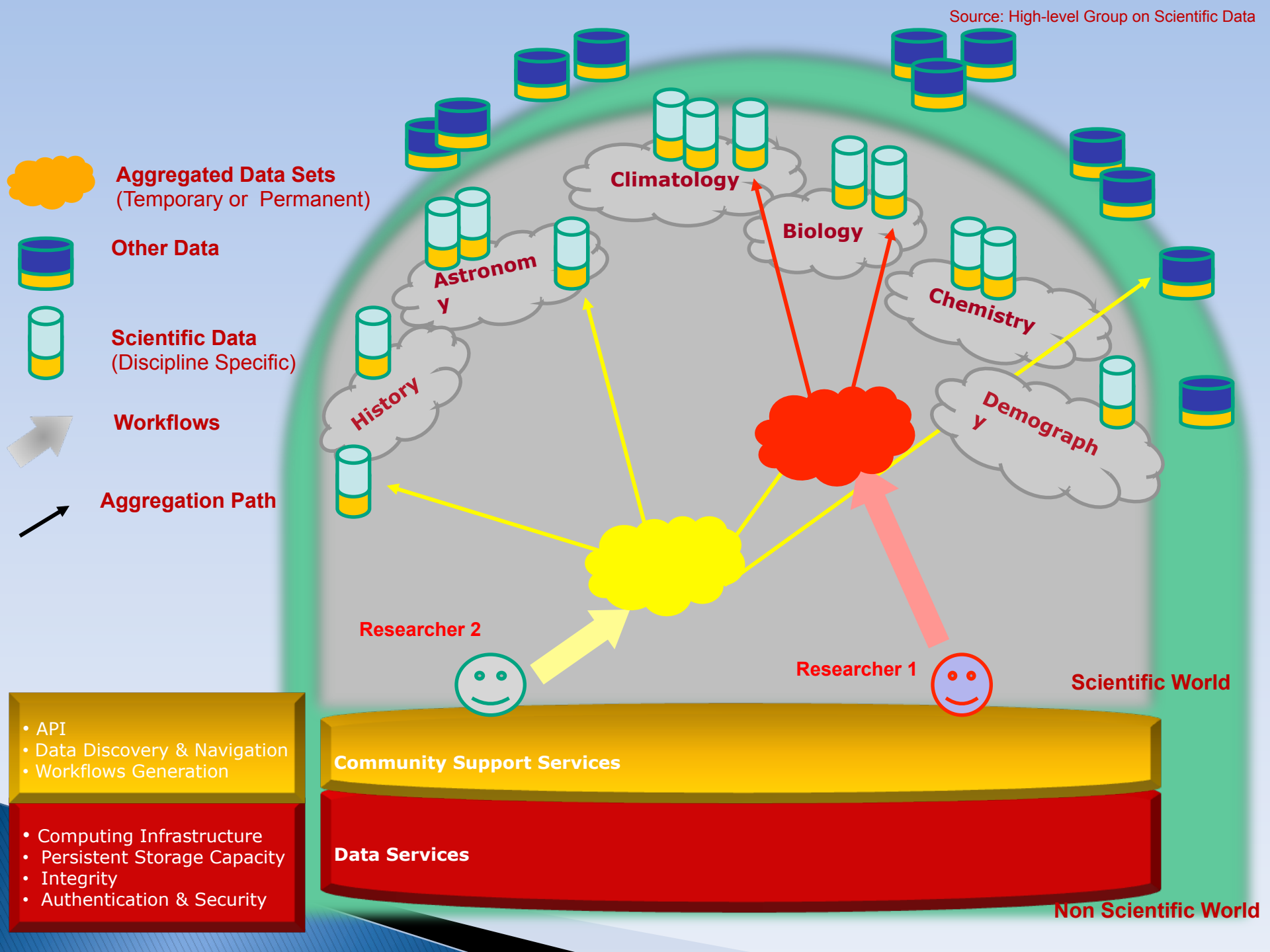
but is “probably” more than that

Wider interpretation: technologies of various kinds for **creating, collecting, annotating, manipulating, storing, finding and re-using information and services** such as those to provide **user support, training, preservation ...**

... information resources and tools such as **vocabs, ontologies, rights management and privacy protection systems, and curation**. Several of these resources depend upon manual **human** input.

eInfrastructure Layers







Riding the wave

How Europe can gain from the rising tide of scientific data

Final report of the High Level Expert Group on Scientific Data
A submission to the European Commission

October 2010

raising tide of data...

“A fundamental characteristic of our age is the **raising tide of data** – global, diverse, valuable and complex. In the realm of science, this is both an opportunity and a challenge.”

Report of the High-Level Group on Scientific Data, October 2010

“Riding the Wave: how Europe can gain from the raising tide of scientific data”

Vision 2030

“Our vision is a scientific e-Infrastructure that supports **seamless access, use, re-use and trust of data**. In a sense, the physical and technical infrastructure becomes invisible and the **data themselves become the infrastructure** – a valuable asset, on which science, technology, the economy and society can advance.”

High-Level Group on Scientific Data

“Riding the Wave: how Europe can gain from the raising tide of scientific data”

scientific Information *con·tin·u·ums* allowed in a “new world”

BETWEEN

experimental data and publications (new paradigm)

different scientific disciplines (multidisciplinary)

past, present and future (preservation)

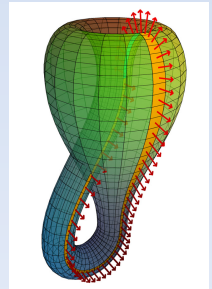
different institutions (organisation)

humans and computers (e-Infrastructure)

research and education (public mission)

current critical needs and visionary environments (evolution)

Reference: Klein Bottle with Moebius Band.
Reference to article "Imaging maths - Inside the Klein bottle" at <http://plus.maths.org/issue26/index.html>. The Klein bottle is a non-orientable surface found by Felix Klein in 1882 while working on a topological classification of surfaces.



Conclusions

- ▶ 4th Paradigm in Scientific Research
- ▶ Data size and complexity fundamental
- ▶ Data management and analysis research challenges
- ▶ Interoperability a huge challenge
 - Data, services, policies, communities, ...
- ▶ Global Research Data Infrastructures

**The name is Science!
d-Science!**

Thank you!

