

EUDAT

Towards a Pan-European Collaborative Data Infrastructure

Dr. - Ing. Morris Riedel

Head of Research Group 'High Productivity Data Processing', Juelich Supercomputing Centre, Germany

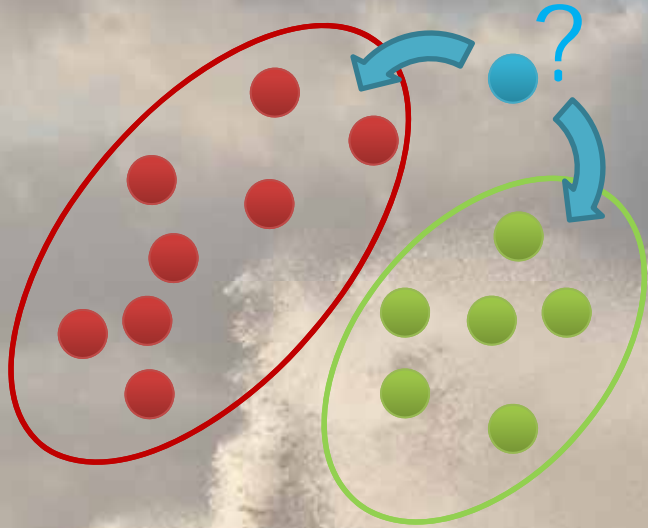
Adjunct Associated Professor, University of Iceland, Iceland



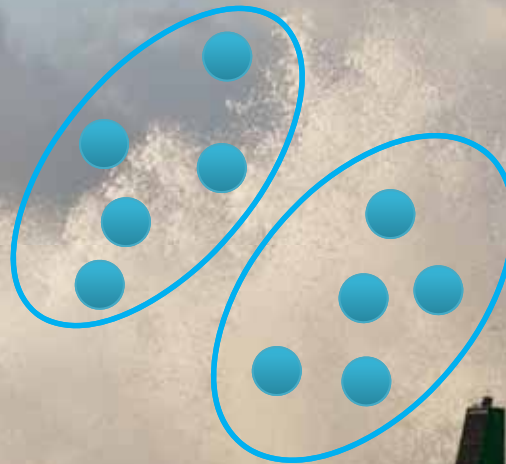
**VSC User Day, Brussels
16th January 2014**



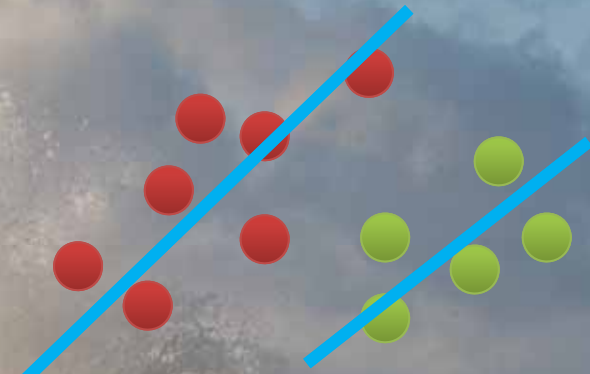
Classification



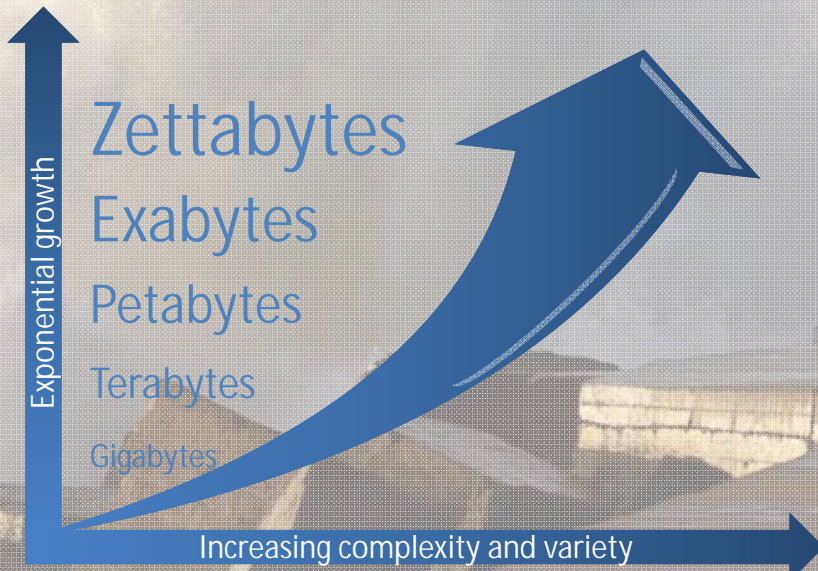
Clustering



Regression



*Data Analysis & Analytics
face 'Big Data Waves'*



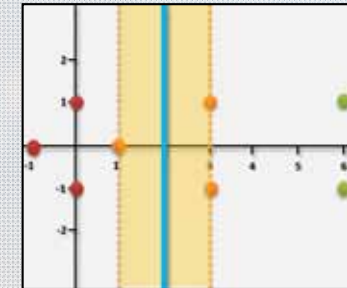
- **Where to store it?**
- **How to find it?**
- **How to make the most of it?**





HPC Simulation Pre-/Post-processing

Data results need to be analyzed and understood
Computed data must be stored and re-located
Subsets of data might be referenced in publications
Sampling vs. whole 'big data' sets (serial/parallel)
Pre-/Post-Processing & visualizations as new data

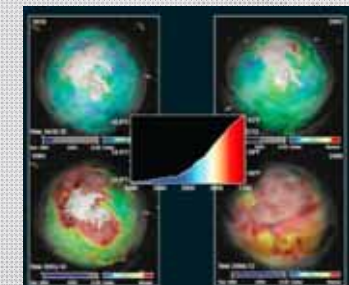


Data Analysis & Analytics in HPC facing limits & challenges



HPC Simulation & Computational Science

Increasing complexity and granularity: data $\rightarrow \infty$
How data is organized has impact on performance
Multi-physics simulations & multi-model ensemble
E.g. physical processes in climate science (land, atmosphere, ocean, sea ice) & observation validation

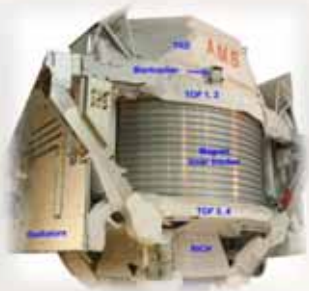


[7] DOE ASCAC report

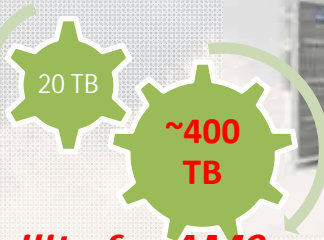


Alpha Magnetic Spectrometer (AMS) @ ISS

What are building blocks of the Universe?



Search for
Cosmic
Antimatter



JUROPA++

1-2 Pflop/s + Booster

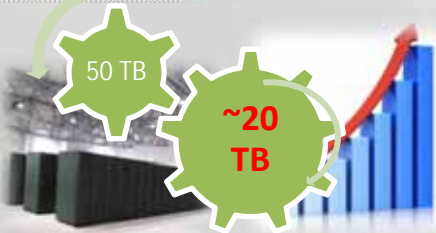
JSC is main facility for AMS computing in Germany

Selected Juelich Supercomputing Centre (JSC) Research Examples

Simulation Labs (e.g. Climate SimLab)...
...one example out of many JSC SimLabs

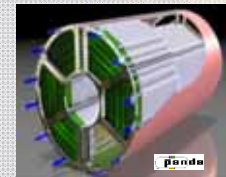
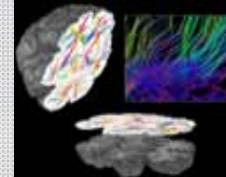


~10 years
500 HDF Files
~50 TB



Applications with
combined characteristics of
simulations and analytics...

Towards Exascale

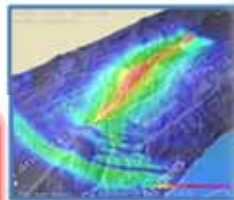


Better Prediction Accuracy... ... means 'Bigger Data'

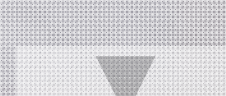
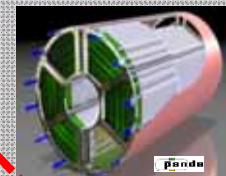
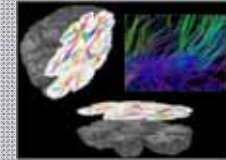
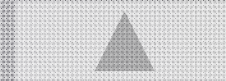
Rank	Site	System	Cores	Speed (TFlop/s)	Speed (TFlop/s)	Power (kW)
1	National Super Computer Center in Guangzhou China	Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P NUDT	3,120,000	33,862.7	54,902.4	17,808
2	DOE/SC/Oak Ridge National Laboratory United States	Titan - Cray XK7, Opteron 6274 16C 2.200GHz, Cray Gemini Interconnect, NVIDIA K20x Cray Inc.	560,640	17,590.0	27,112.5	8,209
3	DOE/NNSA/LNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1,572,864	17,173.2	20,132.7	7,890

TOP 500
HPC
Systems
11/2013

Estimated figures for simulated 240 second period, 100 hour run-time	TeraShake domain (600x300x80 km ³)	PetaShake domain (800x400x100 km ³)
Fault system interaction	NO	YES
Inner Scale	200m	25m
Resolution of terrain grid	1.8 billion mesh points	2.0 trillion mesh points
Magnitude of Earthquake	7.7	8.1
Time steps	20,000 (0.12 sec/step)	160,000 (0.015 sec/step)
Surface data	1.1 TB	1.2 PB
Volume data	43 TB	4.9 PB



'We are unable to store the output data of all computational simulations/users'



Summarizing Big Data Waves & Surfboards

How to engage in the rising tide of scientific data?

Unsolved Questions:

Scale

Heterogeneity

Stewardship

Curation

Long-Term Access and Storage

Research Challenges:

Collection, Trust, Usability

Interoperability, Diversity

Security, Smart Analytics,

Education and training

Data publication and access

Commercial exploitation

New social paradigms

Preservation and sustainability



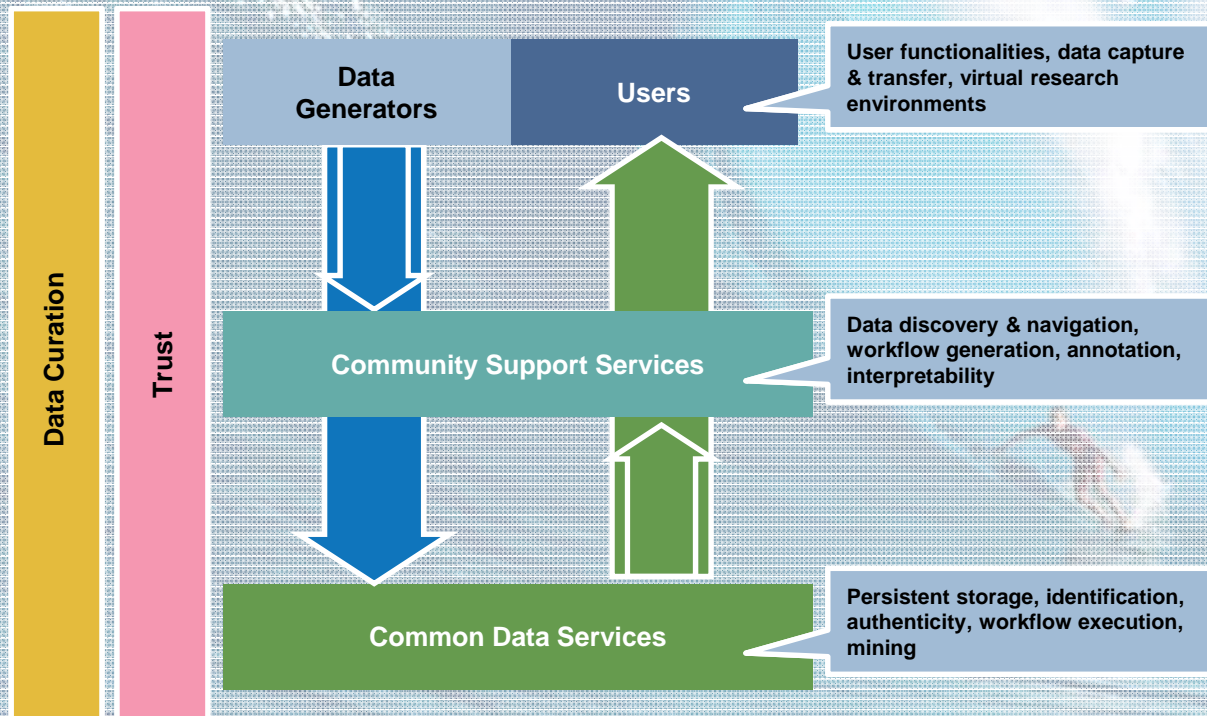
[1] HLEG Report



[2] KE Report

A framework for the future?

Collaborative Data Infrastructure



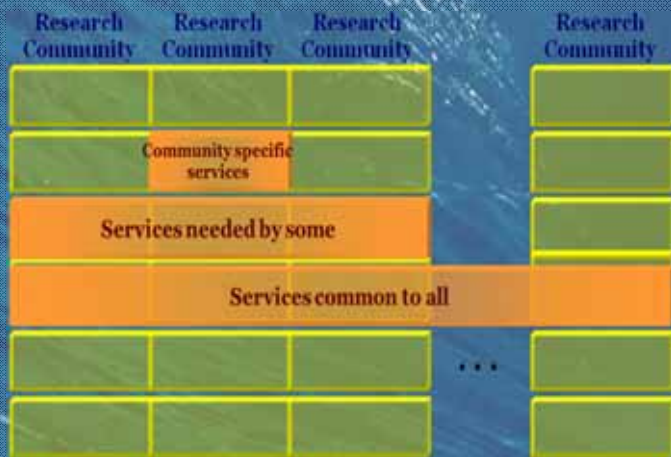
[1] HLEG Report



[3] EUDAT Web Page

Breakwaters – Offer Concrete Solutions for Researchers

Is there a common set of services often needed by scientists?



Identified Common Data Services

- Persistent Identifiers for Research Data
- Safe Replication of Scientific Data
- Transfer of Data to/from Computing
- Simple Sharing of Research Data
- Metadata Catalogue

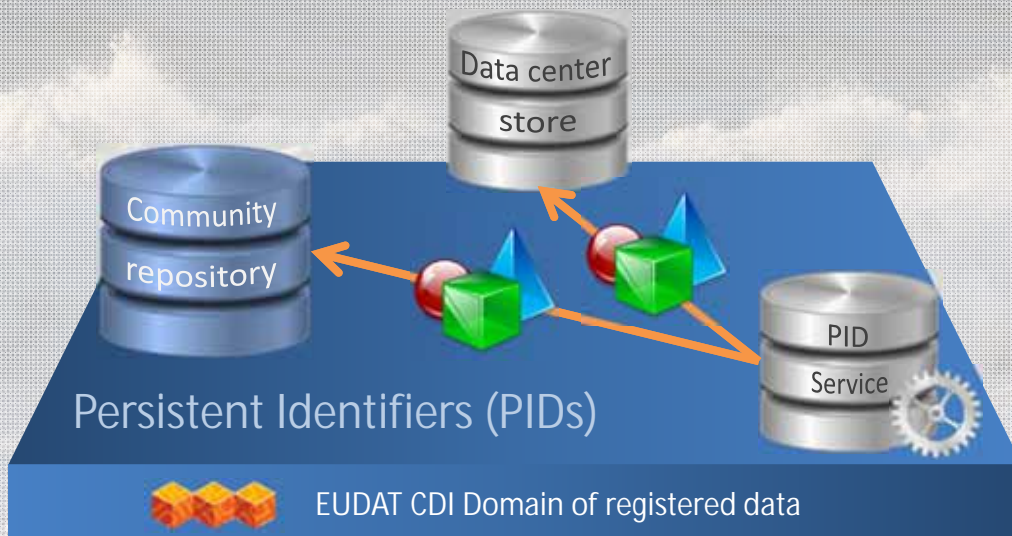
...

*'Concrete'
Next Steps →*

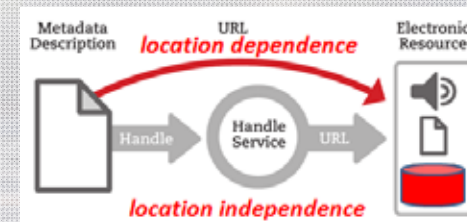


Domain of Registered Research Data

Persistent Identification of Scientific Datasets



- ✓ Provides PID for each data/digital object
- ✓ Based on the 'Handle System'

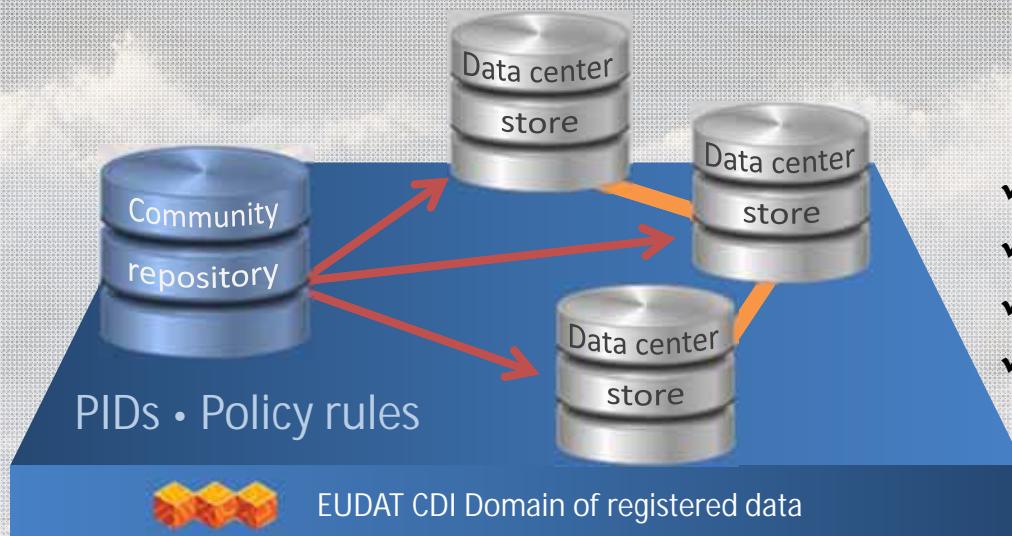


modified from [5] Handle System

Providing a robust, safe, and highly available...

Data Replication Service

...to guard against data loss in long-term archiving & preservation

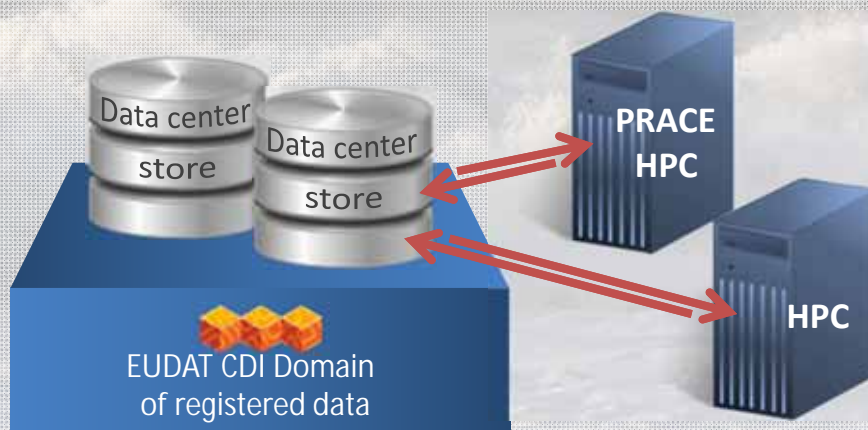


- ✓ Realized in registered data domain
- ✓ Enables reliable data curation
- ✓ Optimize data access for users
- ✓ Provides adaptable policy mechanisms

Bringing research data closer to powerful computers with a...

Data Staging Service

...for compute-intensive scientific data analysis



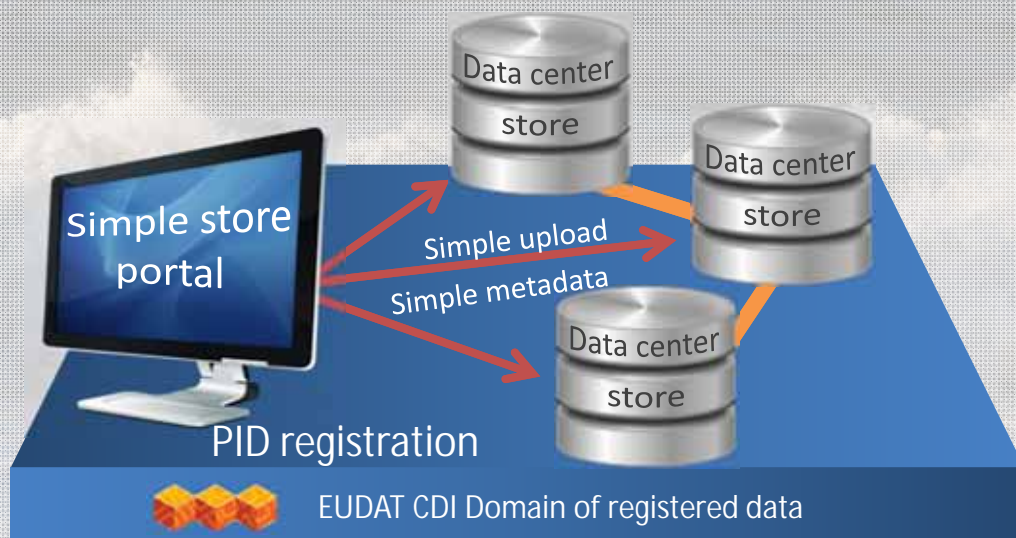
B2STAGE
Get Data to Computation

- ✓ Realized in registered data domain
- ✓ Enables easy access to execution services
- ✓ Offers CPU-intensive data transformations

Offering an easy data deposit and upload via the...

Simple Store Service

...to share data & collections with other researchers

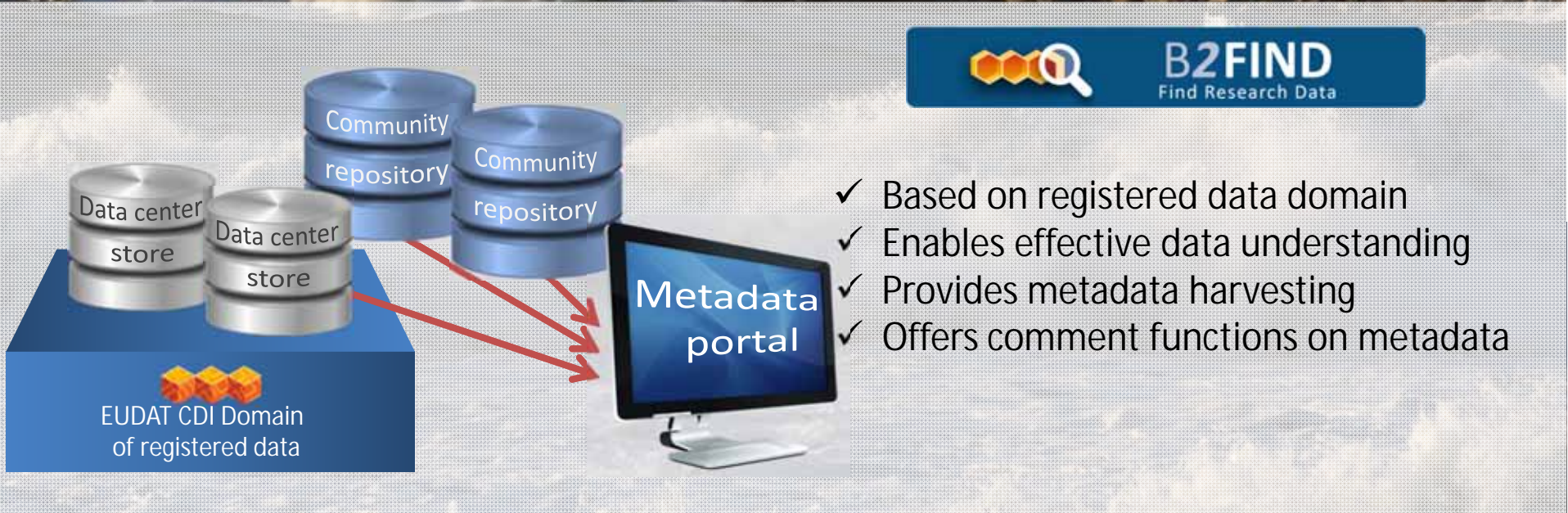


- ✓ Think as 'YouTube for Scientific Data'
- ✓ Realized in registered data domain
- ✓ Upload 'long tail' research data

Find and access research data collections via the...

MetaData Service

...in a simple and user-friendly way





[6] M. Riedel, P. Wittenburg et al., 2013



[3] EUDAT Web Page



[1] HLEG Report

Long-term Data Preservation and Curation...

bears potentials to lower 'Data Waves'

and supports data analytics & analysis



*We need to
'dive into data'*

Addressed requirements of the High Level Expert Group on Scientific Data :

- ✓ *High reliability, so data scientists can count on its availability*
- ✓ *Open deposit, allowing user-community centres to store data easily*
- ✓ *Persistent identification, allowing data centres to register a huge amount of markers to track the origins and characteristics of the information*
- ✓ *Metadata support to allow effective management, use and understanding*
- ✓ *Avoids re-creation of datasets through easy data lookups and re-use*
- ✓ *Enables easier identification of duplicates to remove them & save storage*



[1] HLEG Report

[3] EUDAT Web Page

#	Suggestions for Requirements of a Data Infrastructure	
	Long description	Short description
HLR1	Open deposit, allowing user-community centres to store data easily	Simple data storing
HLR2	Bit-stream preservation, ensuring that data authenticity will be guaranteed for a specified number of years	Bit-stream and long-term preservation
HLR3	Format and content migration, executing CPU-intensive trans-formations on large data sets at the command of the communities	CPU-intensive transformations on large data sets
HLR4	Persistent identification, allowing data centres to register a huge amount of markers to track the origins and characteristics of the Information	Persistent identification of research data
HLR5	Metadata support to allow effective management, use and understanding	Metadata services and harvesting
HLR6	Maintaining proper access rights as the basis of all trust	Proper access rights
HLR7	A variety of access and curation services that will vary between scientific disciplines and over time	Data access and curation services
HLR8	Execution services that allow a large group of researchers to operate on the stored data	Execution services for data analysis
HLR9	High reliability, so researchers can count on its availability	Reliable services
HLR10	Regular quality assessment to ensure adherence to all agreements	Quality assessment
HLR11	Distributed and collaborative authentication, authorisation and accounting	Authentication, authorization & accounting
HLR12	A high degree of interoperability at format and semantic Level	Interoperability



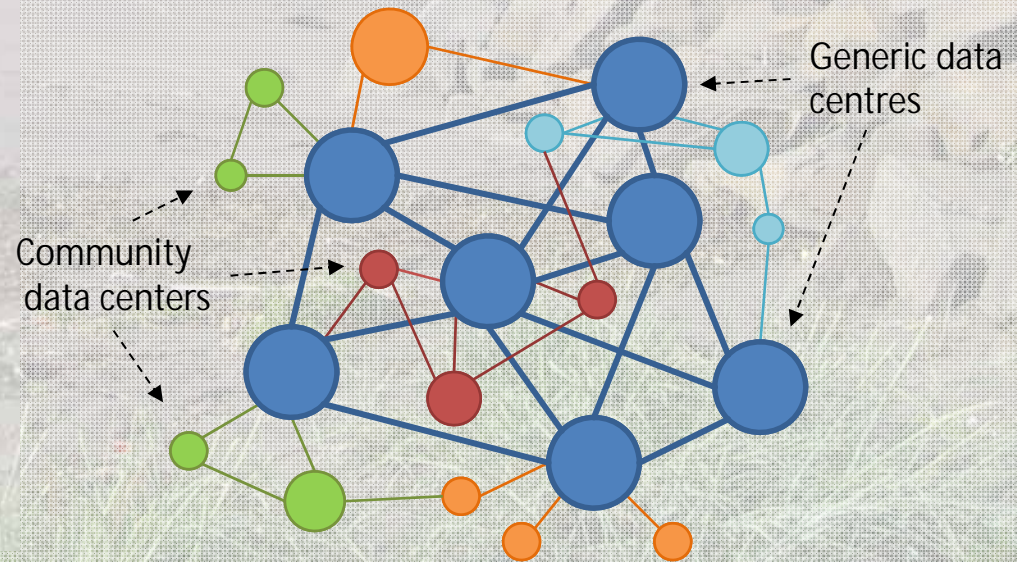
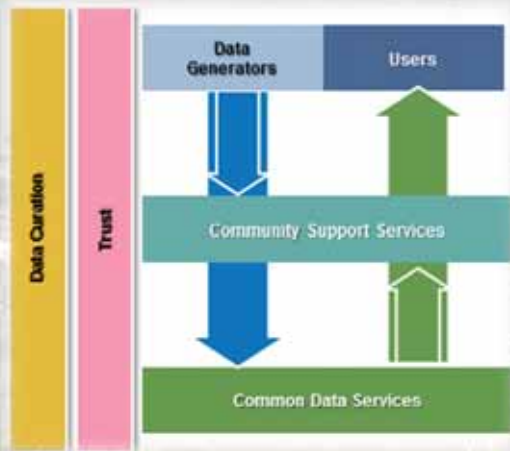
Building a trusted collaborative data infrastructure with...

Strong and Sustainable Community & Generic Data Centers

...to enable federated data services together with users



[1] HLEG Report



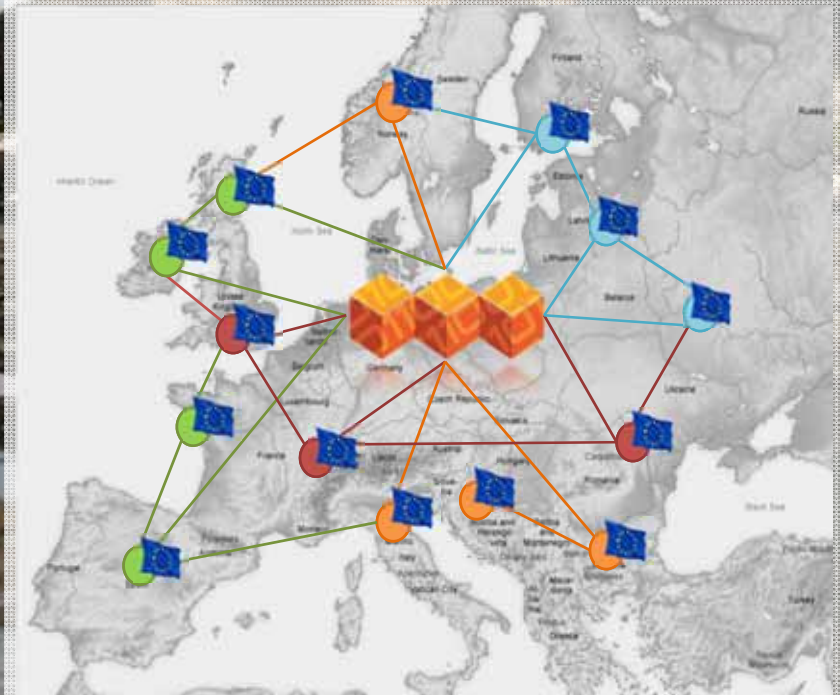


Key Approaches:

Bridging National & EU Solutions
Not 'one single data infrastructure'
Federated Network of Trusted Centers

Key Benefits for Scientific Users:

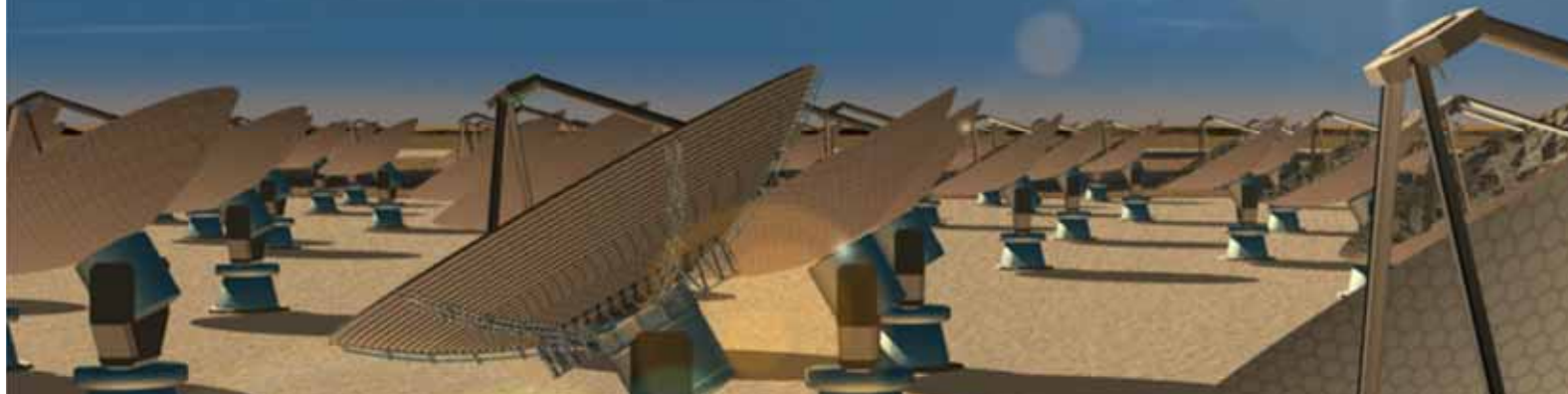
Trust, Sustainability, Interoperability,
Diversity, *Extensibility (e.g. Belgium?)*,
New Social Paradigms & Sustainability



Preparing for new data challenges on the horizon ...

The square kilometre array

... 1 PB in 20 seconds



New EUDAT Services in development with users:

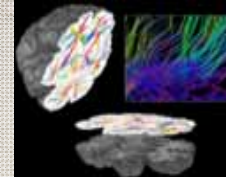
'EUDAT Box'
dropbox-like service
easy sharing
local synching



'Semantic Anno'
checking & referencing

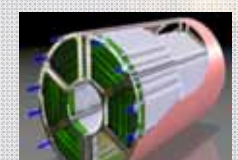
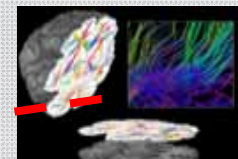
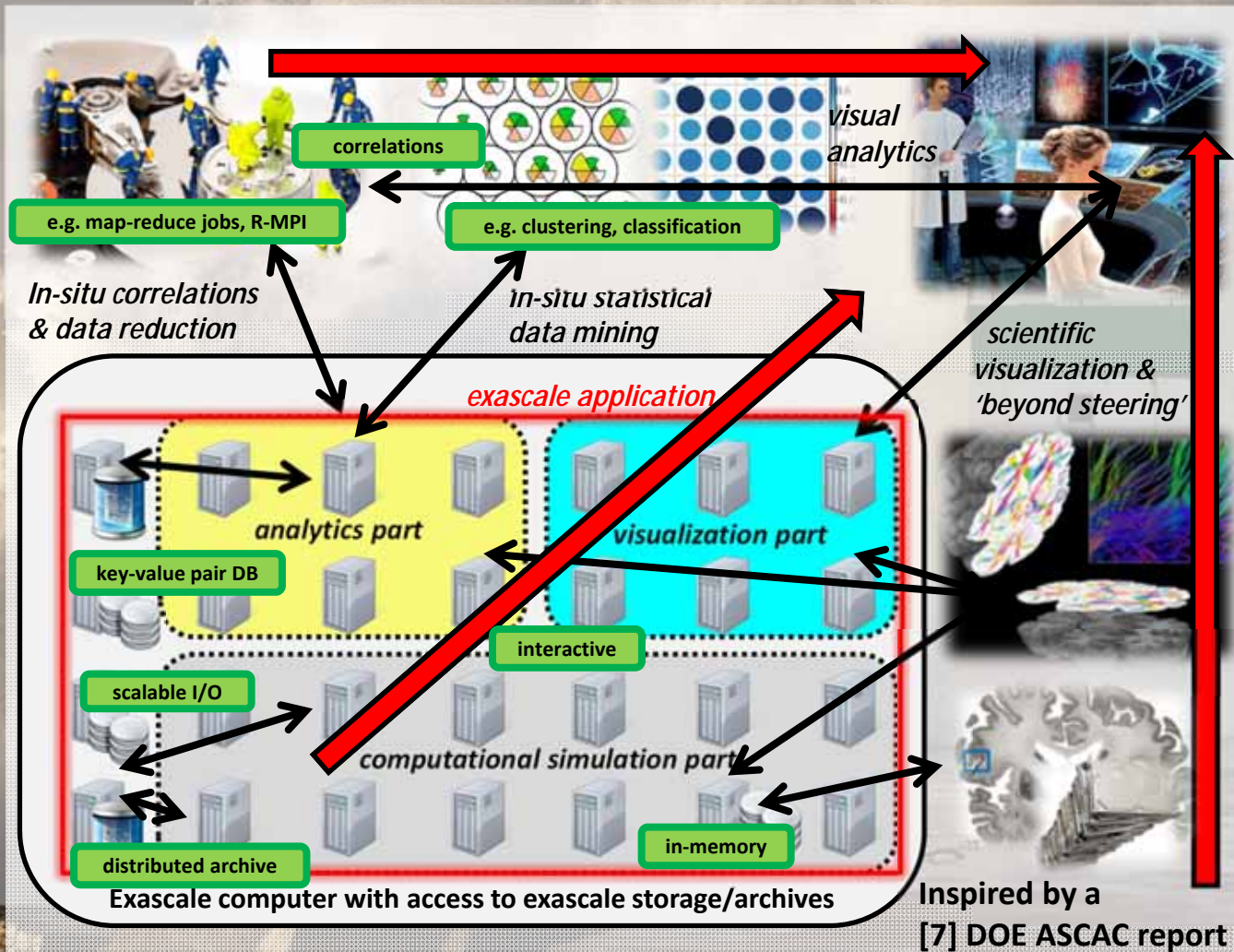


'Dynamic Data'
immediate handling




Towards Exascale: Applications with combined characteristics of simulations & data analytics

'In-Situ Analytics'



Sampling vs. Big Data




Methods & Tools

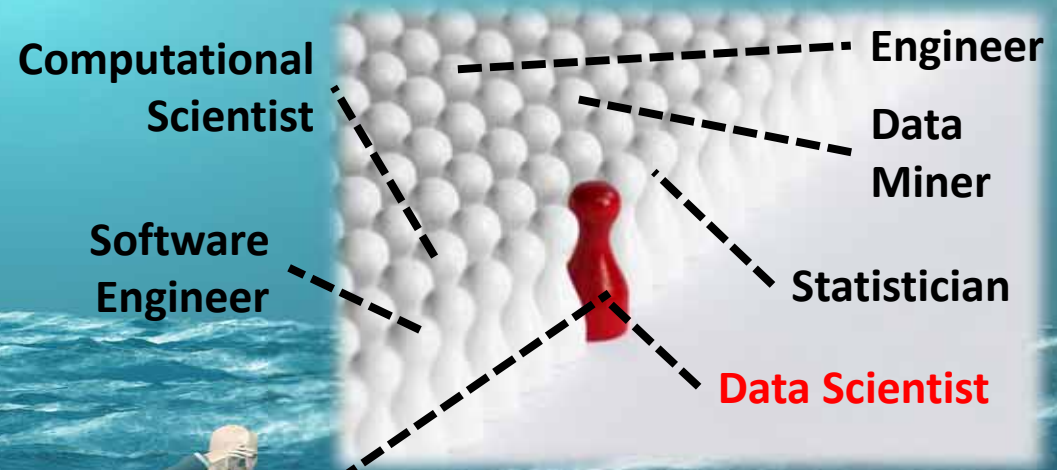
Parallelization!

Applied Statistics Data Mining

Machine Learning Scientific Computing Algorithms



new DBs **Training Data Scientists**




RESEARCH DATA ALLIANCE



Big Data Analytics

Reference Material for Data Scientists

[8] RDA BDA Webpage



EUDAT

Training Opportunity

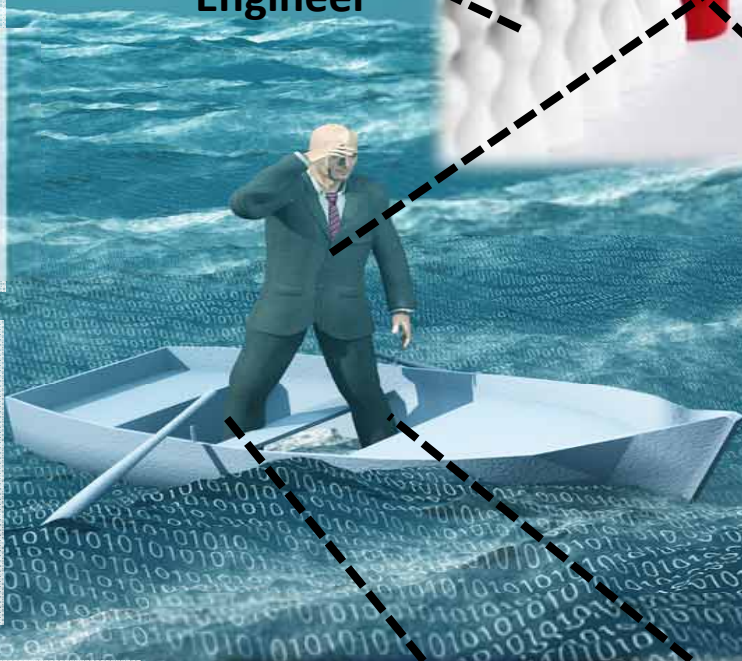
EUDAT Services, Data Management Plans, Curation, ...




UNIVERSITY OF ICELAND
SCHOOL OF ENGINEERING AND NATURAL SCIENCES
FACULTY OF INDUSTRIAL ENGINEERING,
MECHANICAL ENGINEERING AND COMPUTER SCIENCE

Statistical Data Mining Course
HPC – B(ig Data) Course

Data Scientists with skills of various fields



References

- [1] High Level Expert Group on Scientific Data, 'Riding the Wave', Report to the European Commission, October 2010, Online:
<http://cordis.europa.eu/fp7/ict/e-infrastructure/docs/hlg-sdi-report.pdf>
- [2] Knowledge Exchange Partner, 'A Surfboard for Riding the Wave', November 2012, Online:
<http://www.knowledge-exchange.info/surfboard>
- [3] EUDAT Web page, Online:
<http://www.eudat.eu>
- [4] Fran Berman, 'Maximising the Potential of Research Data', January 2013, Online:
<http://rdi2.rutgers.edu/event/rdi2-distinguished-seminar-maximizing-innovation-potential-research-data>
- [5] The Handle System, Handle.net, Online:
<http://www.handle.net/>
- [6] M. Riedel and P. Wittenburg et al. 'A Data Infrastructure Reference Model with Applications: Towards Realization of a ScienceTube Vision with a Data Replication Service', 2013, Online:
<http://www.jisajournal.com/content/4/1/1>
- [7] DOE ASCAC Data Subcommittee Report, 'Synergistic Challenges in Data-Intensive Science and Exascale Computing', 2013, Online:
http://www.sci.utah.edu/publications/chen13/ASCAC_Data_Intensive_Computing_report_final.pdf
- [8] Research Data Alliance (RDA), Big Data Analytics Interest Group, 2014, Online:
<https://www.rd-alliance.org/big-data-analytics-wiki-contents>

'Thanks'

Talk available at:

www.morrisriedel.de/talks

Contact:

m.riedel@fz-juelich.de

