# On Establishing Big Data Breakwaters

### with Analytics

### Dr. - Ing. Morris Riedel

Head of Research Group 'High Productivity Data Processing', Juelich Supercomputing Centre, Germany Adjunct Associated Professor, University of Iceland, Iceland <u>Co-chair Research Data A</u>lliance, Big Data Analytics Interest Group





### **Big Data Waves – Surfboards – Breakwaters** How to engage in the rising tide of big 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 Riding the wave References and has being the advected of the mean sector of the sector

A SURFBOARD FOR RIDING THE WAVE

[2] KE Report

[1] HLEG Report







## Volume

Variety

Velocity

Context





### 'Crowdsourcing'...

## .. increases # of Big Data Streams





Data streams with data (low trust)



Data streams with data (moderate trust)



Exabytes



Scientific/Engineering Domain Experts

Data streams with data (high trust)





The next generation radio telescope for science...

... pushing the limits of the observable universe out by billions of galaxies

## The square kilometre array

# ... 1 PB in 20 seconds









LOFAR test site Jülich

#### Large-scale Computational Massively Parallel Applications simulate Reality

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

Rank	Site		System			Cores	(TFlop/s)	(TFlop/s)	(kW)		
0	National Super Computer Center in Guangzhou China		Tianbe-2 (MilkyWay-2) - TH-IV8-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	TOP 500 HPC	
0	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	Systems	
9	DOE/NNSA/LINL United States		Sequola - B Custom IBM	ueGene/Q, Power BQC 16	C 1.60 GHz,	1,572,864	17,173.2	20,132.7	7,890	11/2015	
	Estimated figures for simulated 240 second period, 100 hour run-time	TeraShak (600x300	ke domain Ix80 km^3)	PetaShake domain (800x400x100 km^3)		A					
	Fault system interaction	NO		YES				/	We are		
	Inner Scale	20	0m	25m	-	112		nahla	to	store the	
	Resolution of terrain grid	1.8 billion mesh points		2.0 trillion mesh points							
	Magnitude of Earthquake	ke 7.7		8.1							
	Time steps	Time steps 20,0		00 160,000 c(top) (.0015 ccc/(top)			2				
1	Surface data	1.1 TB		1.2 PB		Sec.	computational				

**p**anda

2 Rensselaer

Volume data

43 TB

F. Berman: 'Maximising the Potential of Research Data' Information courtesy of the Southern California

4.9 PB

Earthquake Cente

[3] F. Berman

simulations/users'

# Making use of Big Data is important... ... but How?



#### Netflix Prize Challenge ~2009

Challenge: Improve movie recommender Prize: 1 M \$ for team with 10% improvements Open historical data to 'train a new system' Traditional Machine Learning Algorithms used Prize received with Artificial Neural Networks



It becomes a strategic differentiator in the market

#### LETTERS

#### nature

#### Detecting influenza epidemics using search engine query data

#### H1N1 Virus Made Headlines ~2009

Nature paper from Google employees Explains how Google is able to predict flus Not only national scale, but down to regions Possible via logged big data - 'search queries' Google: Seek for web using Google Coogle Seek (The beside Seek Coogle Seek (The beside Seek) Coogle Seek) Coogle Seek (The beside Seek) Coogle Seek (The bes

[4] Jeremy Ginsburg et al., 'Detecting influenza epidemics using search engine query data', Nature 457 (2009)

### Smart Analytics are Needed to Take Advantage of Big Data The challenge is to understand which analytics make sense

'Understanding climate change, finding alternative energy sources, and preserving the health of an ageing population are all cross-disciplinary problems that require high-performance data storage, **smart analytics**, transmission and mining to solve.'

*'In the data-intensive scientific world, new skills are needed for creating, handling, manipulating, analysing, and making available large amounts of data for re-use by others.'* 



[1] HLEG Report

[5] DOE ASCAC Report

*Integration of data analytics* with exascale simulations represents a new kind of workflow that will impact both data-intensive science and exascale computing.'





### Big Data Analytics Interest Group

RESEARCH DATA ALLIANCE

#### Can we establish a

'neutral group gathering clear evidence' which analytics make sense in Research

- <u>Develops community based recommendations on feasible data analytics</u> approaches to address community needs of utilizing large quantities of data
- <u>Analyzes different scientific research domain applications</u> and their potential use of various big data analytics techniques
- <u>A systematic classification of feasible combinations</u> of analysis algorithms, analytical tools, data and resource characteristics and scientific queries will be covered in these recommendations



### Big Data Analytics Interest Group

#### RESEARCH DATA ALLIANCE

- Agricultural Data Interoperability IG
- Big Data Analytics IG
- Brokering IG
- Certification of Digital Repositories IG
- Community Capability Model WG
- Data Citation WG
- Data Foundation and Terminology WG
- Data in Context IG
- Data Type Registries WG
- Defining Urban Data Exchange for Science IG
- Digital Practices in History and Ethnography IG
- Engagement Group IG
- Legal Interoperability IG
- Long tail of research data IG

- Marine Data Harmonization IG
- Metadata IG
- Metadata Standards Directory WG
- PID Information Types WG
- Practical Policy WG
- Preservation e-Infrastructure IG
- Publishing Data IG
- Standardization of Data Categories and Codes IG
- Structural Biology IG
- Toxicogenomics Interoperability IG
- UPC Code for Data IG
- Wheat Data Interoperability WG

[6] Research Data Alliance Big Data Analytics IG Web page





### **Big Data Analytics Interest Group**

RESEARCH DATA ALLIANCE

### New Technology Impact As Challenge for Users 'NoSQL Databases Example'

#### **Selected Features**

Simplicity of design and deployment Horizontal scaling Less constrained consistency models Finer control over availability Simple retrieval and appending

#### Types

Key-Value-based (e.g. Cassandra) Column-based (e.g. Apache Hbase) Document-based (e.g. MongoDB) Graph-based (e.g. Neo4J)









#### Big Data Analytics Techniques require their Parallelization



[10] Gesmundo et al.



[11] Zhanquan Sun et al.



✓ Parallel K-Means
Clustering Algorithms

 ✓ Parallel Perceptron Using Map-Reduce  ✓ Parallel Support Vector Machines



## 'Underutilized Unique Big Data'

'Big Data' not always means 'Big Volume', but also Context (uniqueness)

In-flight Measurements ( $\rightarrow$ 'Context')

MOZAIC 1994 – today (~20 years):

Ozone, water vepor, CO, NOy, + p, T, Wind

IAGOS 2011 – today ('better data'):

✓ + "Cloud droplets" + CO<sub>2</sub>, CH<sub>4</sub>, Aerosole





**Challenge:** 















**Big Data** 

Analytics

Large 'underutilized' data collections exist



















## **Automating Quality Control**

But 'Bigger Data' not always means 'Better data'



Challenge: Data items ~7.9 billions Large data collection exist with many measurements

#### Approach:

 Initial general data analytics techniques
→ E.g. automated outlier detection for quality control (Previously only manual sampling, but better?)

Domain-specific focussed data validation by experts



*With thanks to Robert Huber, Marum, Bremen* 















Long-term Data Preservation and Curation... bears potentials to lower 'Data Waves'

and supports big data analytics process



Selected Benefits of open data infrastructures for science & engineering:

- ✓ *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







## Talk available at:

#### www.morrisriedel.de/talks

## **Contact:**

### m.riedel@fz-juelich.de