

# Big Data Architectures and Technologies

Marcel Kunze

Research Group Cloud Computing - Steinbuch Centre for Computing

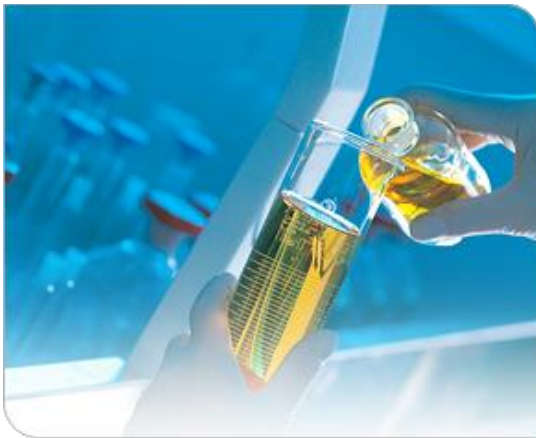


# Karlsruhe Institute of Technology (KIT)

Largest European scientific Institution

Main Topics: Energy, Nanotechnology, Astrophysics, Engineering

Mission:



■ Research



■ Education



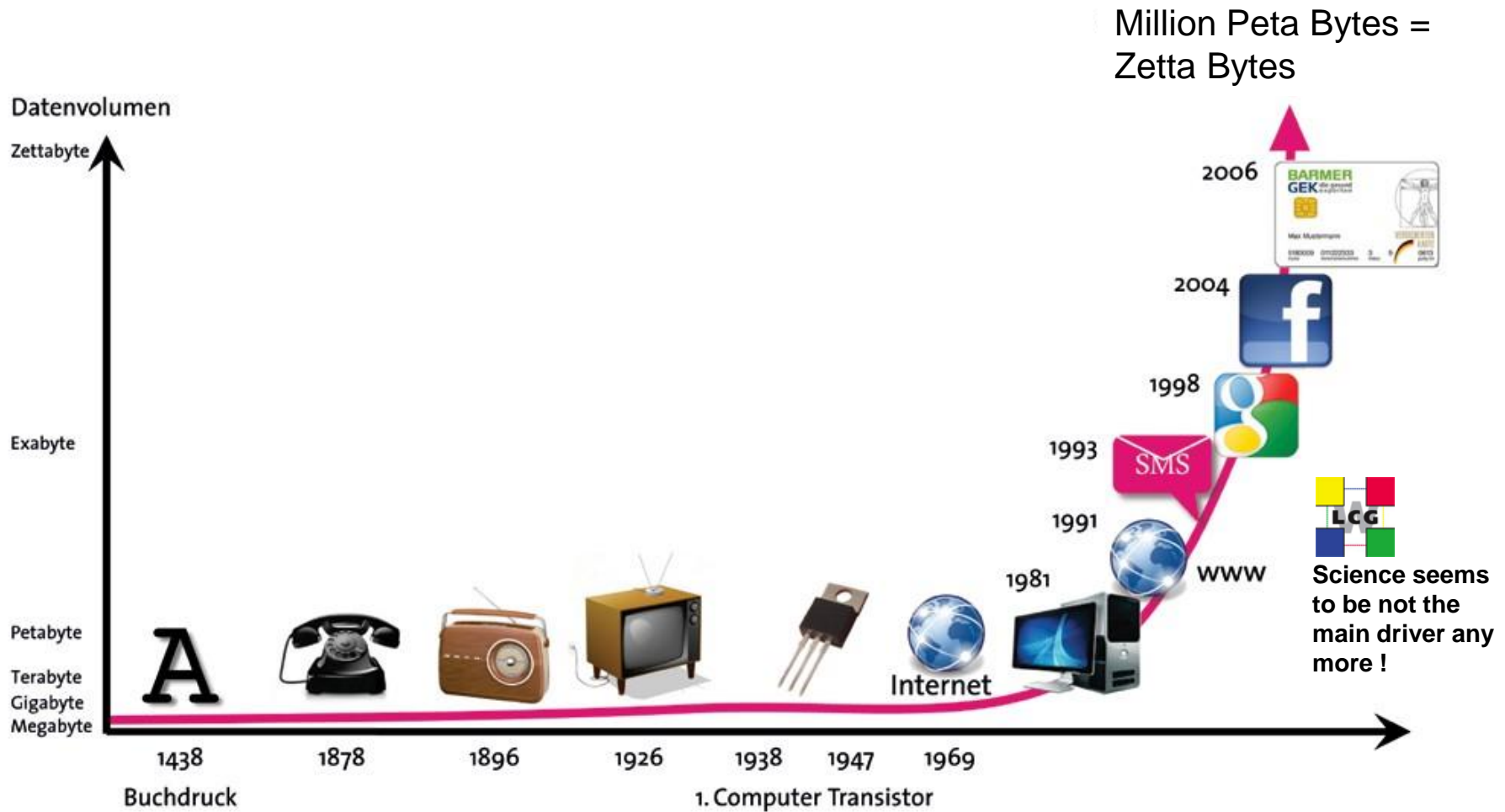
■ Innovation

# Agenda

- What is “Big Data” ?
- Big Data Management
- Big Data Toolbox
- R&D Projects

# What is “Big Data” ?

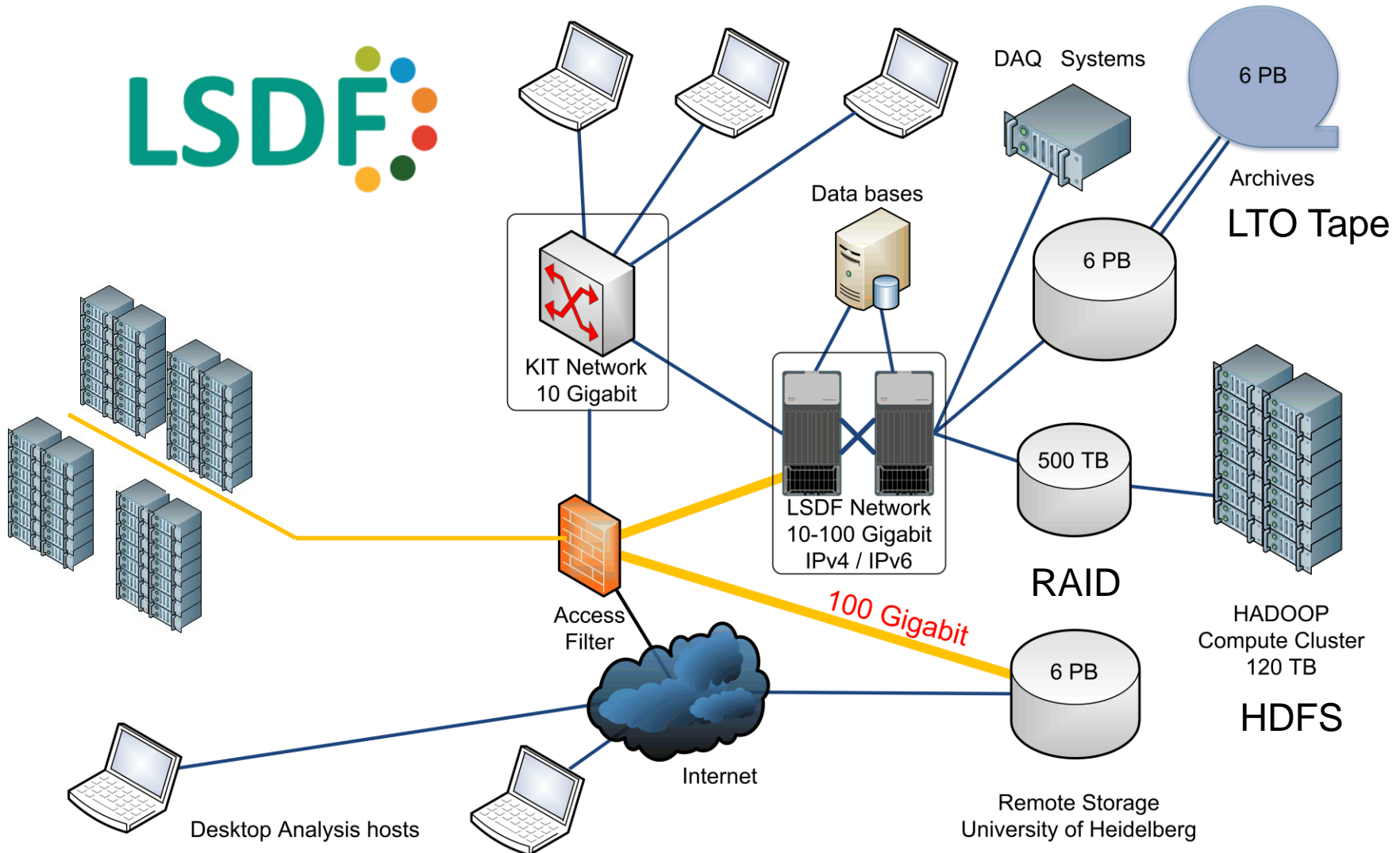
# Big Data Volume



Source: „Big-Data im Praxiseinsatz, Leitfaden“, BITKOM 2012

# Large Scale Data Facility (LSDF)

KIT Institutes



# The Data Center as a Computer

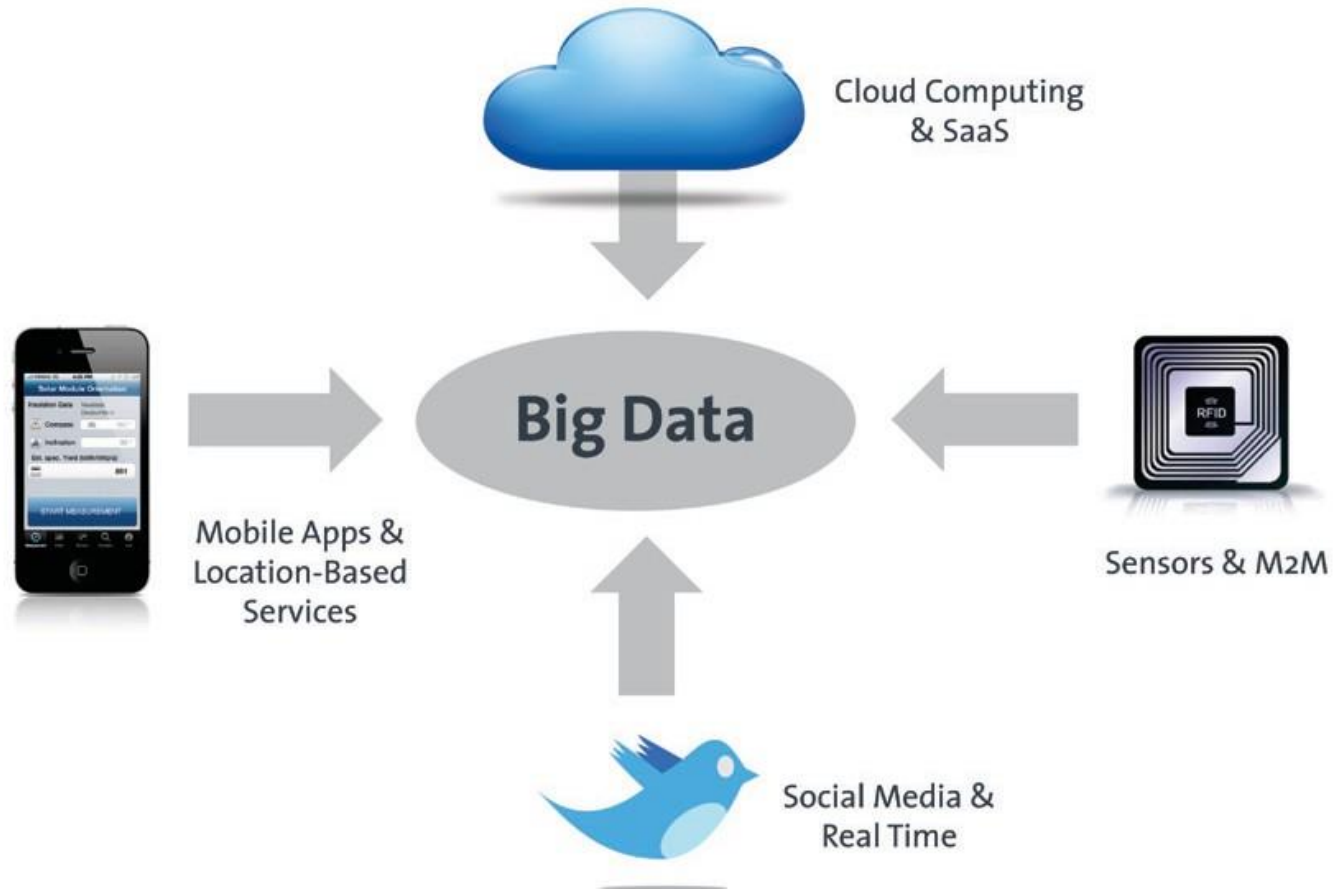


■ Apple data center in Maiden, NW Carolina: Exa-Scale data center (iCloud)



Comparison:  
SCC@KIT,  
LSDf, LHC Tier-1

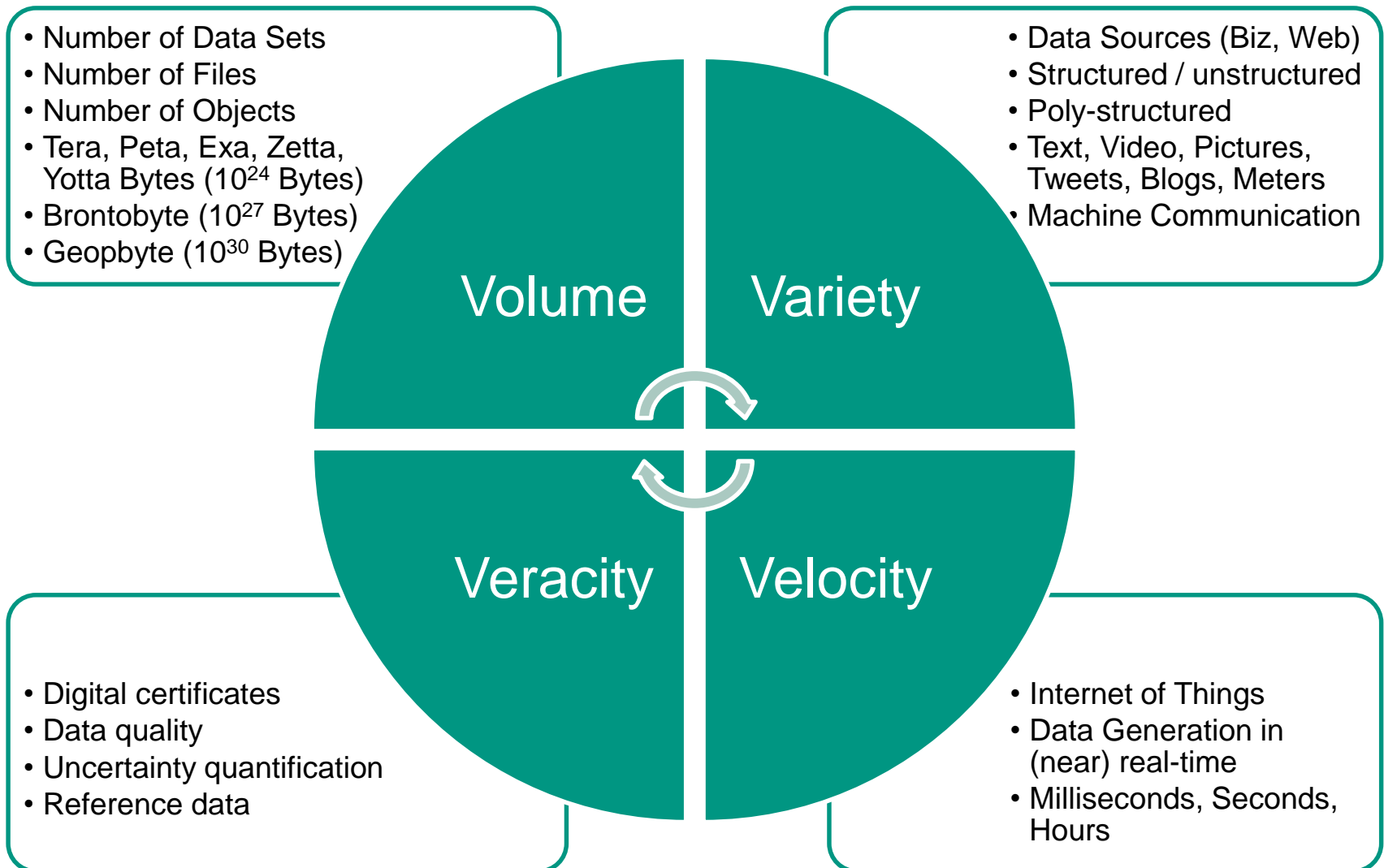
# Big Data Drivers



Source: „Big-Data im Praxiseinsatz, Leitfaden“, BITKOM 2012



# Big Data Dimensions (4V)



# The 5<sup>th</sup> Dimension: Value

- A major new trend in information processing will be the trading of original and enriched data, effectively creating an information economy
  - Data mining
  - Descriptive analytics (Past)
  - Predictive analytics (Future)
  - Prescriptive analytics (Actionable insight)
    - Correlation of data
    - Intelligence of patterns, relations, etc.
    - ...

*„When hardware became commoditized, software was valuable. Now that software is being commoditized, data is valuable.“ (TIM O'REILLY)*

*„The important question isn't who owns the data. Ultimately, we all do. A better question is, who owns the means of analysis?“ (A. CROLL, MASHABLE, 2011)*

# Added Value in the Data Economy



Big Data Technologies and Cloud Infrastructure

- Example: Weather data - heterogeneous data/simulation sources – forecast/probability – map/app

# Big Data Management

# Data Management

- How can we manage billions of files/objects?
- How can we manage exabytes of data?
  - Data security
  - Data usage control
  - Data movement
  - Data availability
  - Data preservation
  - Data publication

# Data Security

- There is no absolute security, neither in the cloud nor on your own premises. Spies are everywhere...
- Risk management:
  - Risk = probability of disaster \* cost of disaster
  - Given a lower cost, a higher probability of data loss could be tolerable
  - A web server and clickstream analyses might be safely migrated into a public cloud
- Strong encryption helps to treat compliance and legal issues. Control of:
  - Keys
  - Algorithms
  - Data



**“Encryption helps...”**

# Data Usage Control

- Serious cloud providers are offering corresponding services that even comply with the highest legal standards
  - Regional concepts to define the data location
  - Trusted data stores with encryption
  - Identity and access management
  - Data usage control via transactional trails logging

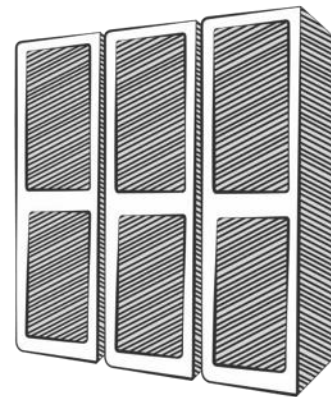


Who



had What

access to



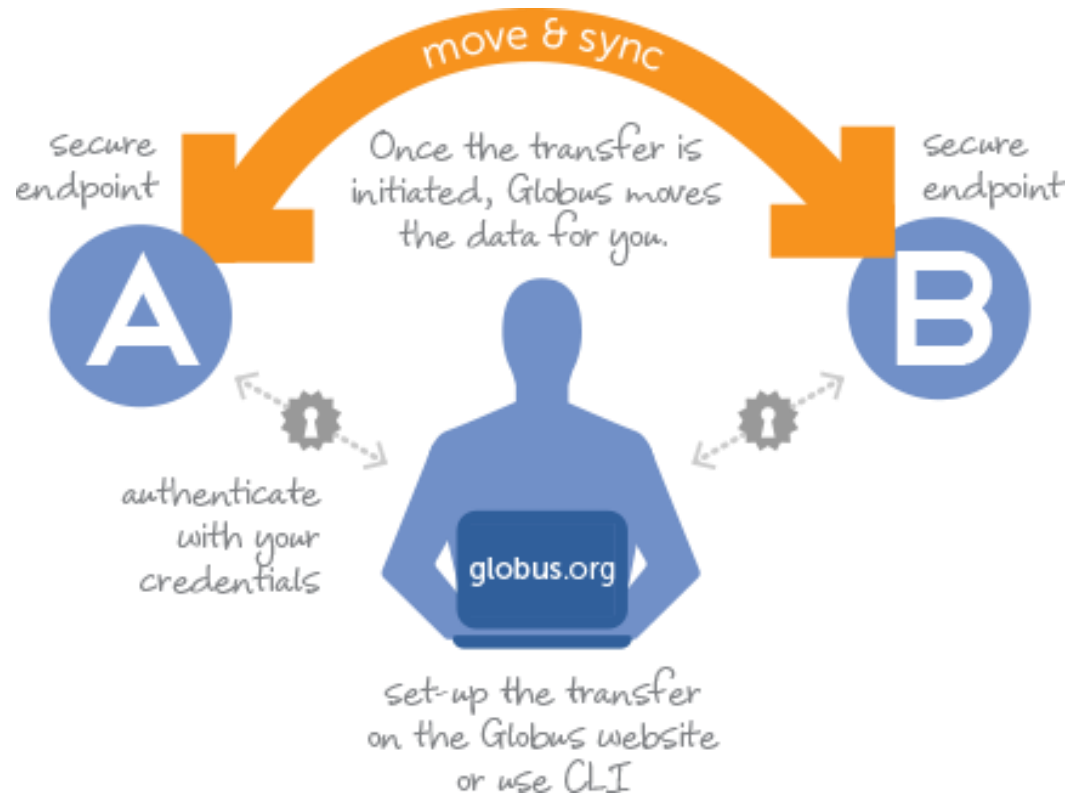
Which resources



When

Source: Amazon CloudTrail

# Data Movement



Source: <https://www.globus.org/>

- Globus Online: Fully managed and automated Big Data file transfer
- SaaS: Separation of control plane and data plane (Third party transfer)



# Data Publication

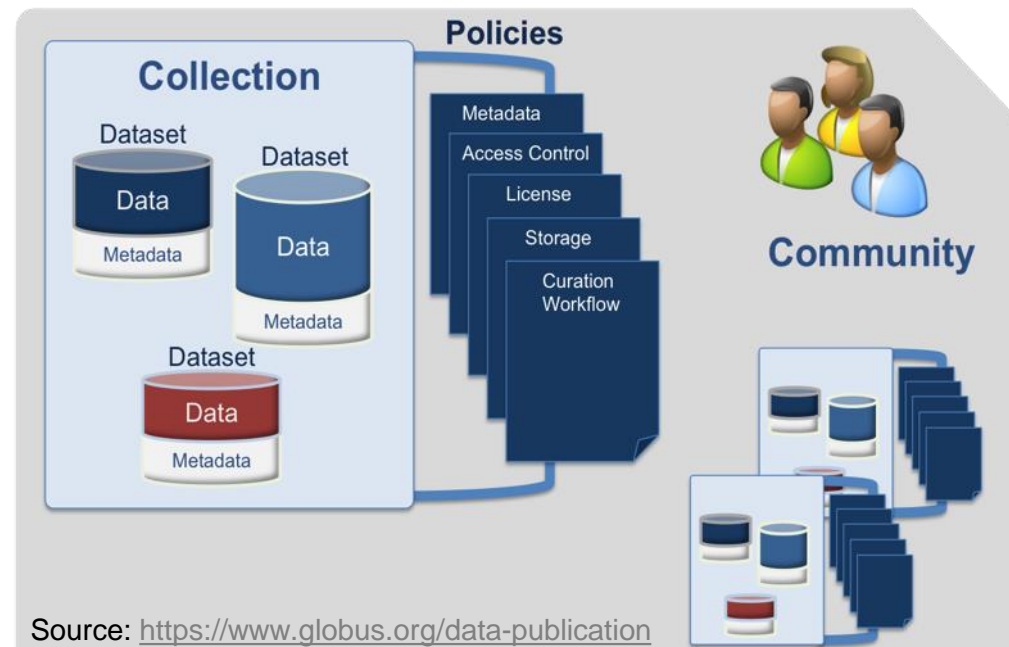
- What does it mean to publish?
- Data is:
  - Identified
  - Described
  - Curated
  - Verifiable
  - Accessible
  - Preserved

# Data Discovery

- What does it mean to discover?
- Data can be:
  - Searched
  - Browsed
  - Accessed

# Globus Data Publication Services

- SaaS for publishing Big Data
- Bring your own storage
- Extensible metadata
- Publication and curation workflows
- Public and restricted collections
- Rich discovery model
  
- Architecture:
  - Communities create collections of datasets
  - Describe the metadata
  - Define policies for data access
  - Define processes for publication



# Supply Domain specific Meta Data

g globus
blaiszik

Describe
Describe
Upload
Verify
License
Complete

## Submit: Describe this Item ?

Please fill further information about this submission below.

Enter appropriate subject keywords or phrases below.

**Subject Keywords**

self-healing	Remove Entry	circuit	Remove Entry
microcapsules	Remove Entry		+ Add More

Enter the names of any sponsors and/or funding codes in the box below.

**Sponsors**

This material is based upon work supported as part of the Center for Electrical Energy Storage - Tailored Interfaces, an Energy Frontier Research Center funded by the U.S. Department of Energy, Office of Science, Office of Basic Energy Sciences under Award Number (919 DOE ANL 9F-31921 NS).

Enter a description for this item in the box below.

**Description**

Thermomechanical failure of conductive pathways in highly integrated circuits results in loss of function that is often impossible to repair and remains a long-standing problem hindering advanced electronic packaging. Prior approaches to restoration of conductivity rely on external intervention in the form of heating or manual delivery of relatively low conductivity materials. Here, we demonstrate autonomic healing of an electrical circuit with nearly full recovery of conductance (ca. 99%) less than one millisecond after damage. The rapid restorative

Enter the name of experiment for this item below.

**Experiment**

self-healing-10vtpcr

Enter the names of materials used in this experiment below.

**Material**

Gallium	Remove Entry	Gold	Remove Entry
Indium		circuitboard	+ Add More

Enter the energy density used in this experiment.

**Energy Density (mAh/g)**

2000

Enter the Argonne GUP that this experiment was conducted under.

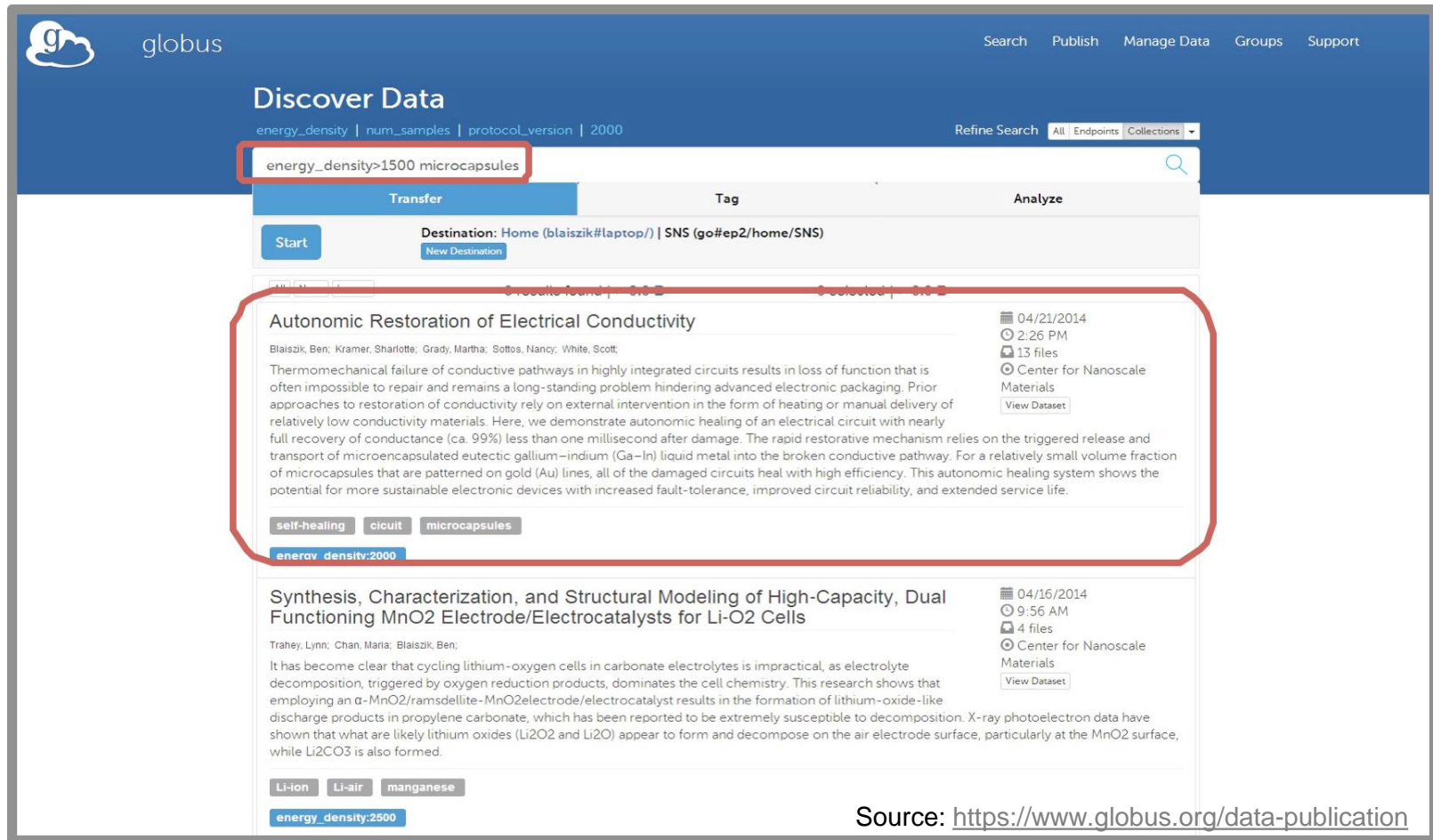
**GUP**

345-455-2543

< Previous
Cancel/Save
Next >

Source: <https://www.globus.org/data-publication>

# Search within and across Collections



The screenshot shows the Globus Discover Data interface. At the top, there's a search bar with the query "energy\_density>1500 microcapsules" highlighted by a red box. Below the search bar, there are tabs for "Transfer", "Tag", and "Analyze". The "Transfer" tab is active, showing a "Start" button and a "Destination" field set to "Home (blaiszik#laptop/) | SNS (go#ep2/home/SNS)".

The search results are displayed in a list. The first result is titled "Autonomic Restoration of Electrical Conductivity" and is highlighted with a red rounded rectangle. It includes the following information:

- Title:** Autonomic Restoration of Electrical Conductivity
- Authors:** Blaiszik, Ben; Kramer, Shariotte; Grady, Martha; Sottos, Nancy; White, Scott
- Date:** 04/21/2014
- Time:** 2:26 PM
- Files:** 13 files
- Organization:** Center for Nanoscale Materials
- View Dataset:** [button]
- Tags:** self-healing, circuit, microcapsules
- Filter:** energy\_density:2000

The second result is titled "Synthesis, Characterization, and Structural Modeling of High-Capacity, Dual Functioning MnO2 Electrode/Electrocatalysts for Li-O2 Cells". It includes the following information:

- Title:** Synthesis, Characterization, and Structural Modeling of High-Capacity, Dual Functioning MnO2 Electrode/Electrocatalysts for Li-O2 Cells
- Authors:** Trahey, Lynn; Chan, Maria; Blaiszik, Ben;
- Date:** 04/16/2014
- Time:** 9:56 AM
- Files:** 4 files
- Organization:** Center for Nanoscale Materials
- View Dataset:** [button]
- Tags:** Li-ion, Li-air, manganese
- Filter:** energy\_density:2500

Source: <https://www.globus.org/data-publication>

■ Globus may be a pathfinder project to create open data markets

End-to-end sequencing analysis.  
Flexible, scalable, simplified.



## A solution as cutting-edge as your research.

Globus Genomics combines state-of-the-art algorithms, data management tools, a graphical workflow environment, and an elastic computing infrastructure. We take care of IT complexity so you can focus on your research.

### Researchers

Easily leverage advanced tools to scale your work.

[learn more](#)

### Core Labs

Provide customers with cost-effective services.

[learn more](#)

## Scalability and flexibility for big genomics data.

Genome sequencing is notoriously data-intensive. We make it easy to manipulate, store, and share your data, no matter how big it gets.

## Your workflows, simplified.

We help you easily connect and configure all the tools you need to automate your pipelines.

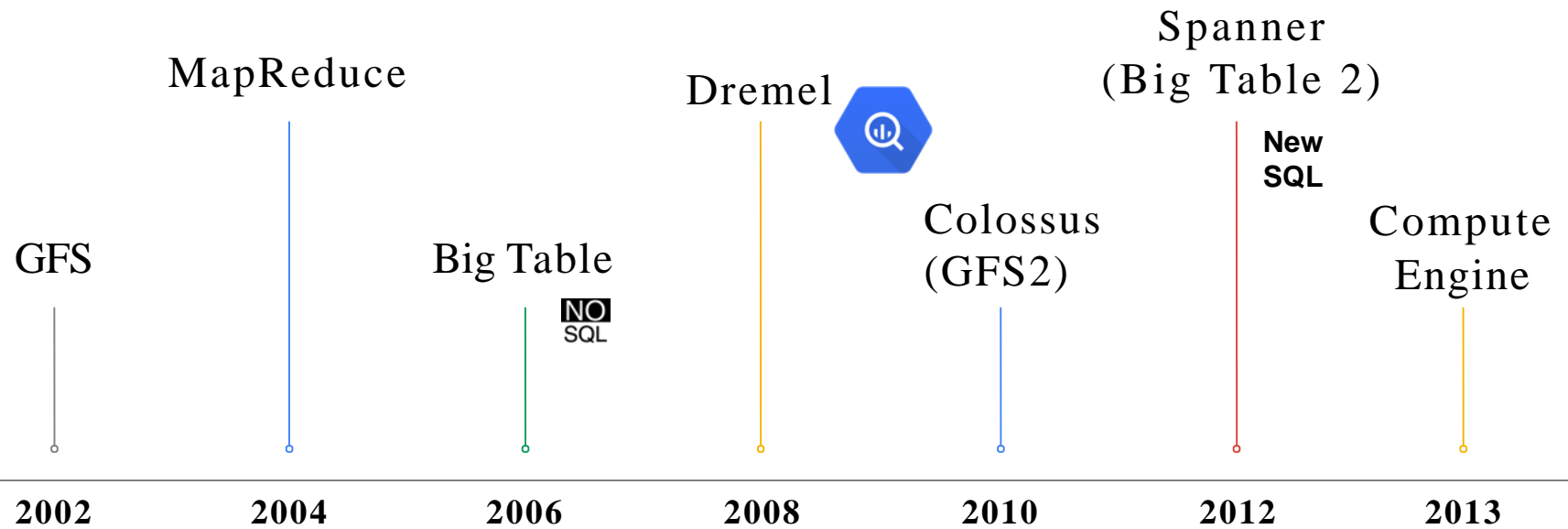
## Let the cloud power your analysis.

Globus provides the data management solution using Amazon's compute cloud. The efficiencies are all yours.

Source: <https://www.globus.org/genomics/>

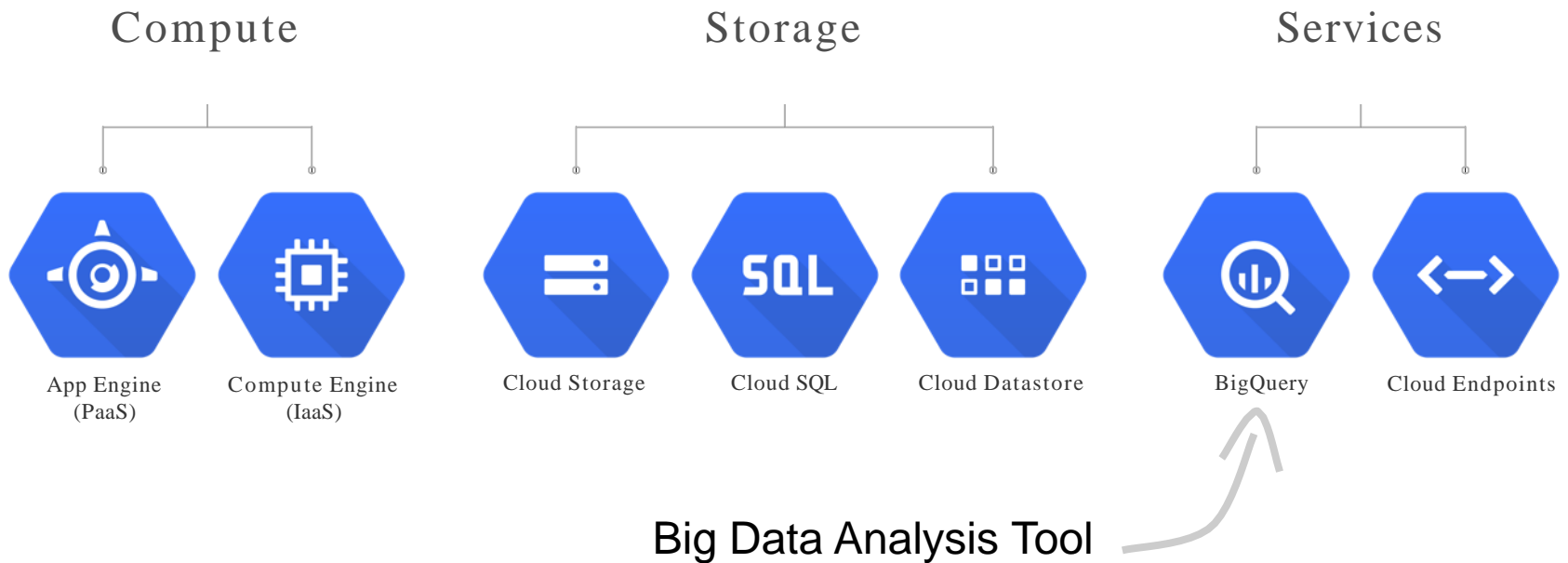
# Big Data Toolbox

# Google Innovations over the last twelve Years





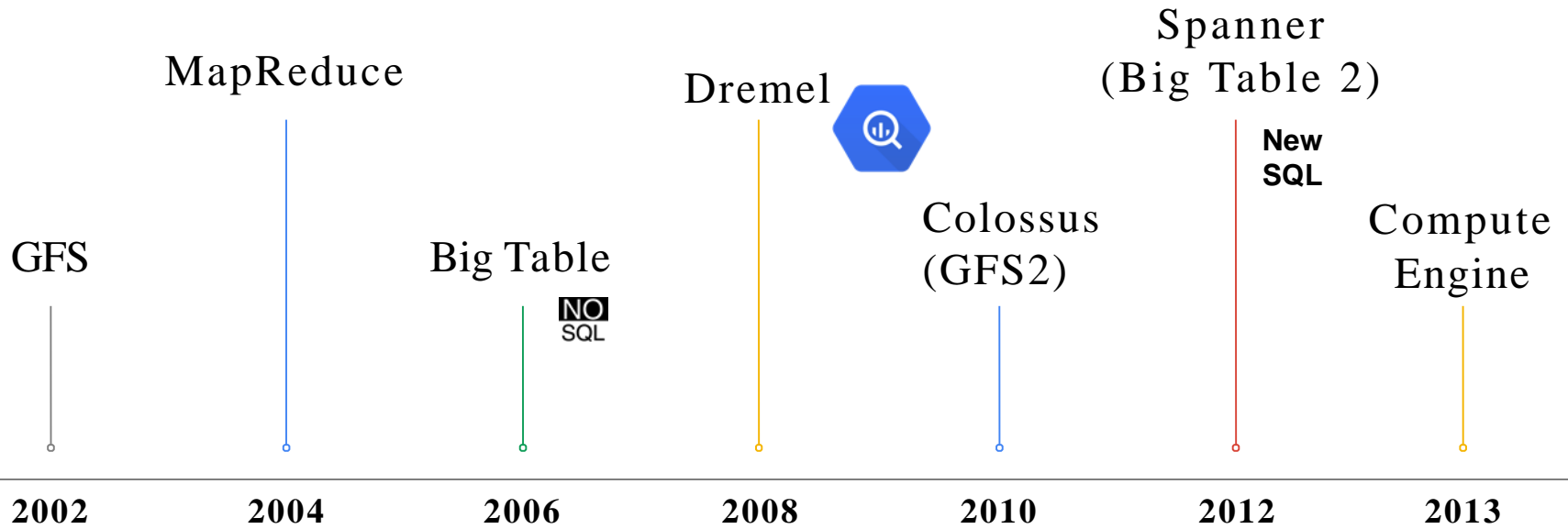
# Google Cloud Platform



- Analyze terabytes of data in seconds
- Data imported in bulk or using streaming
- Supports CSV and JSON
- Browser tools, command line tool, or REST-API



# Google Innovations over the last twelve Years



Opensource:





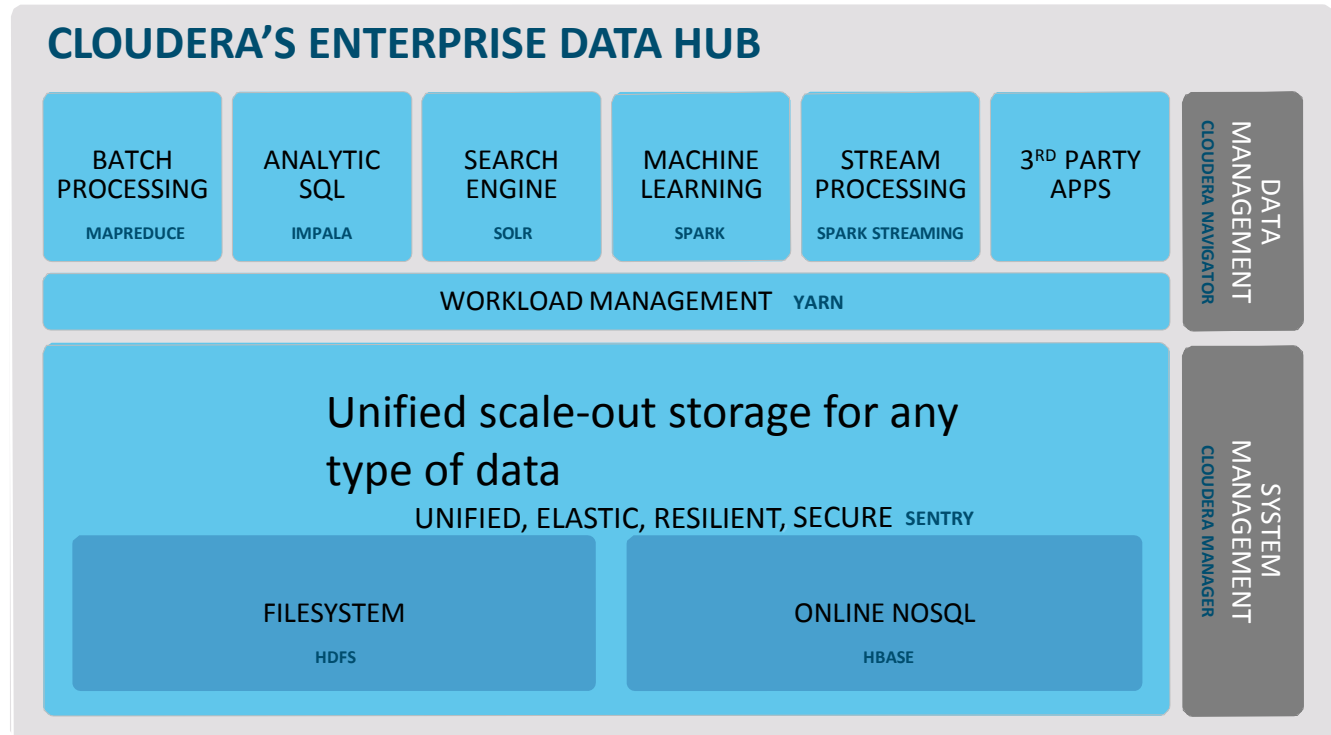
# Hadoop

- Hadoop is a Big Data ecosystem that implements
  - Hadoop core utilities
  - HBase: A scalable, distributed database for large tables.
  - HDFS: A distributed file system.
  - Hive: A data warehouse, data summarization and ad hoc querying.
  - MapReduce: distributed processing on compute clusters.
  - Oozie: Workflow management.
  - Pig: A high-level data-flow language for parallel computation.
  - Spark: Ultra-fast in-memory computing.
  - ZooKeeper: coordination service for distributed applications.
  - And much more ...



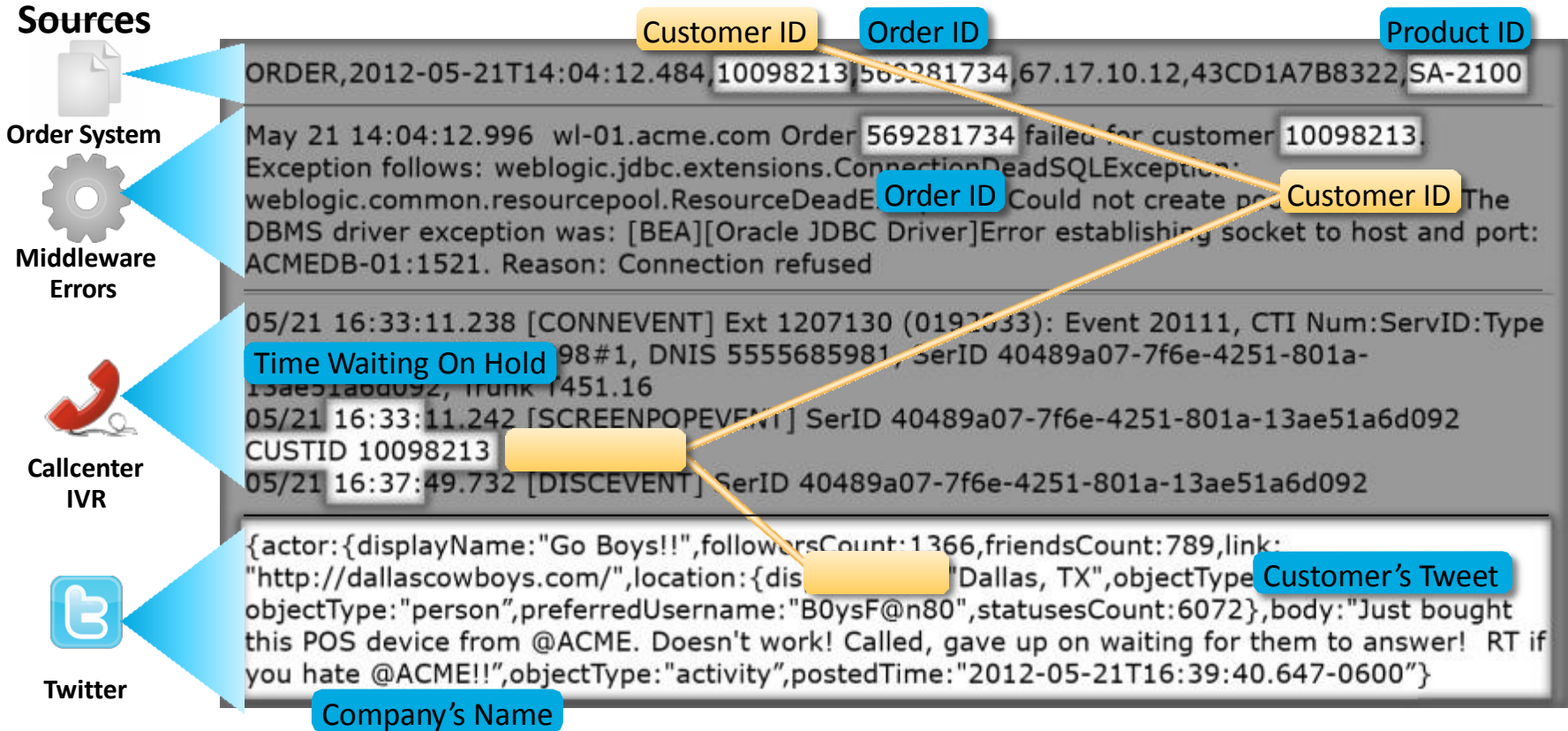
# From Hadoop to the Enterprise Data Hub

<b>Opensource</b> <b>Scalable</b> <b>Flexible</b> <b>Cost effective</b>	✓
<b>Managed</b>	✓
<b>Open Architecture</b>	✓
<b>Secure</b>	✓



<http://www.cloudera.com/content/cloudera/en/products-and-services/cloudera-enterprise.html>

# How do we conveniently treat Machine Data?

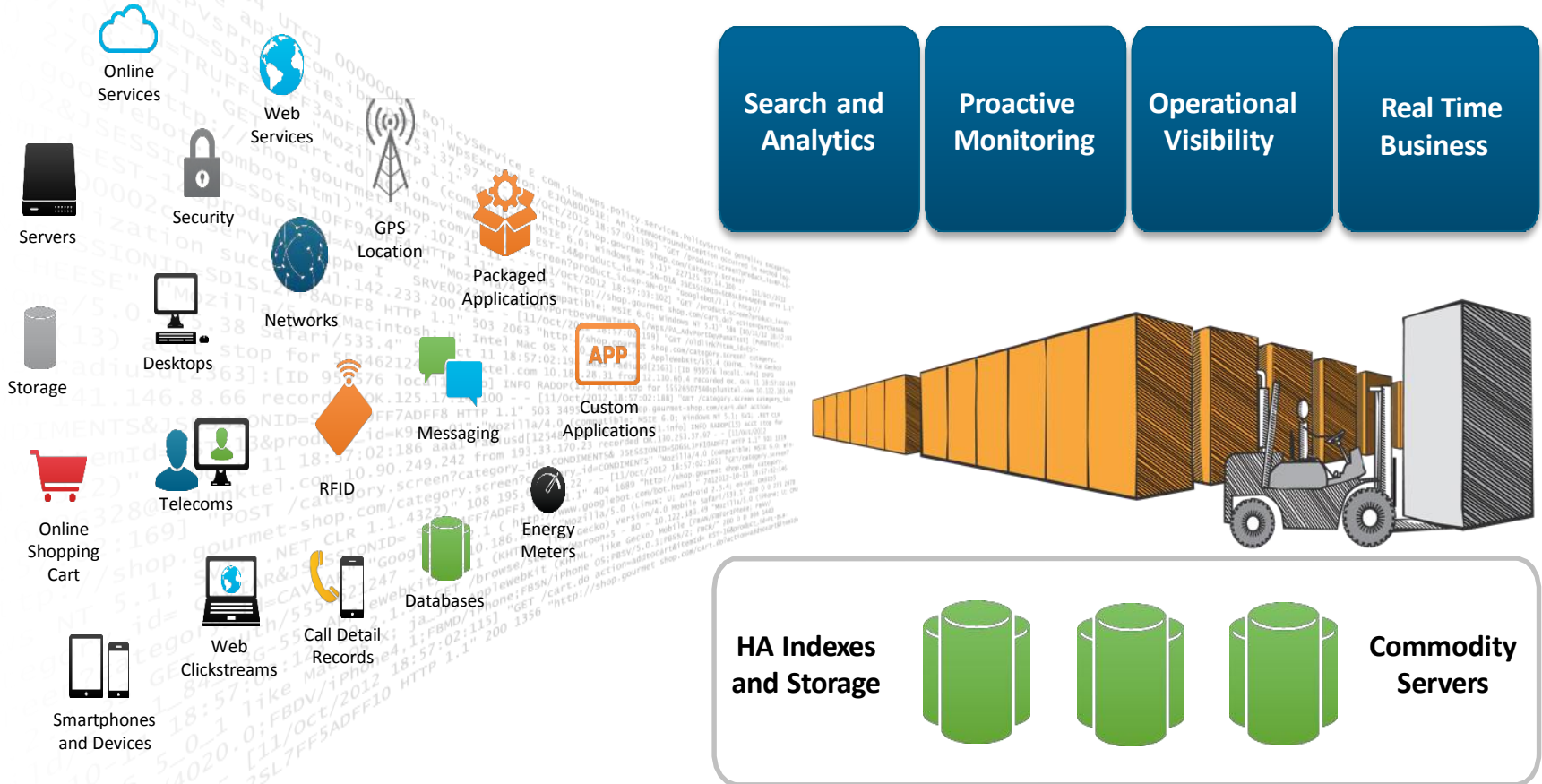


■ Variety: Data are “poly-structured” and live in lots of formats

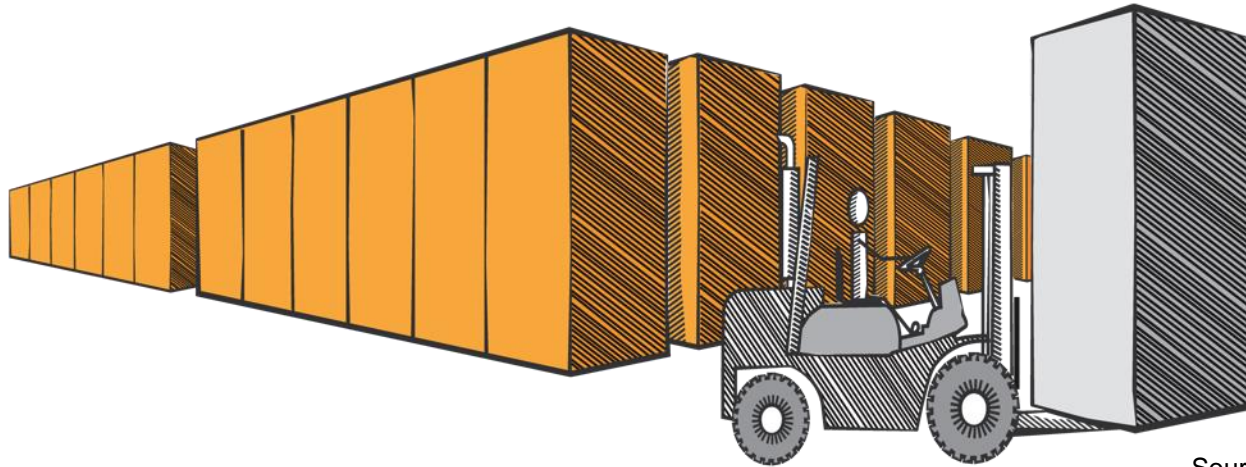
# The Data Warehouse

## Machine Data

## Operational Intelligence



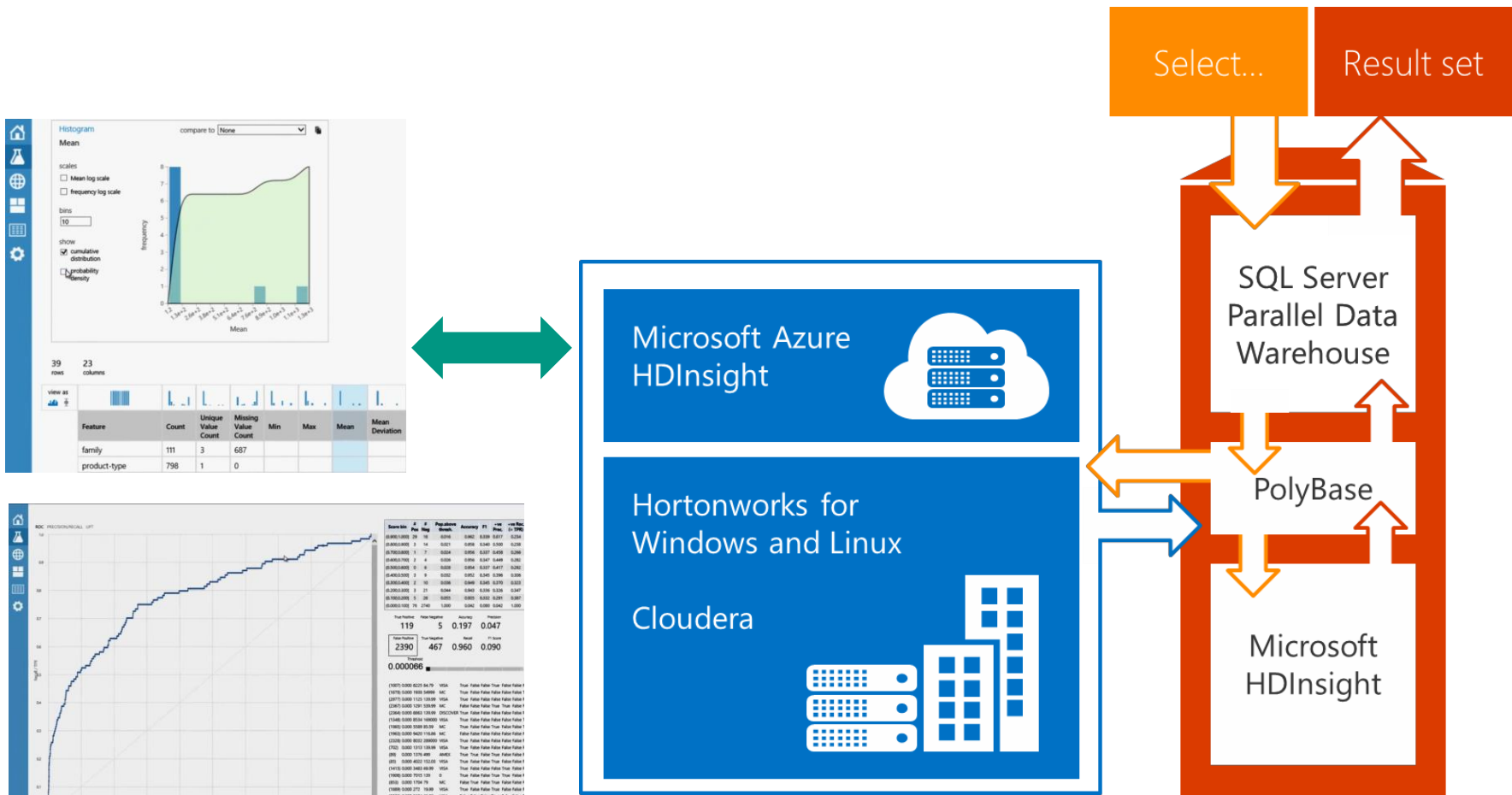
# Big Data Analytics in the Cloud: Amazon Redshift



Source: [Amazon Redshift and DynamoDB](#)

- Fully managed, massively parallel relational data warehouse
- Takes care of cluster management and distribution of data
- Optimized for complex queries across many large tables
- Use standard SQL & standard BI tools
- Can be combined with Hadoop on-demand (Elastic MapReduce)

# Microsoft Analytics Platform System



- Bringing together Hadoop with the data warehouse (Windows Azure)
- New: Microsoft machine learning service (Azure ML)

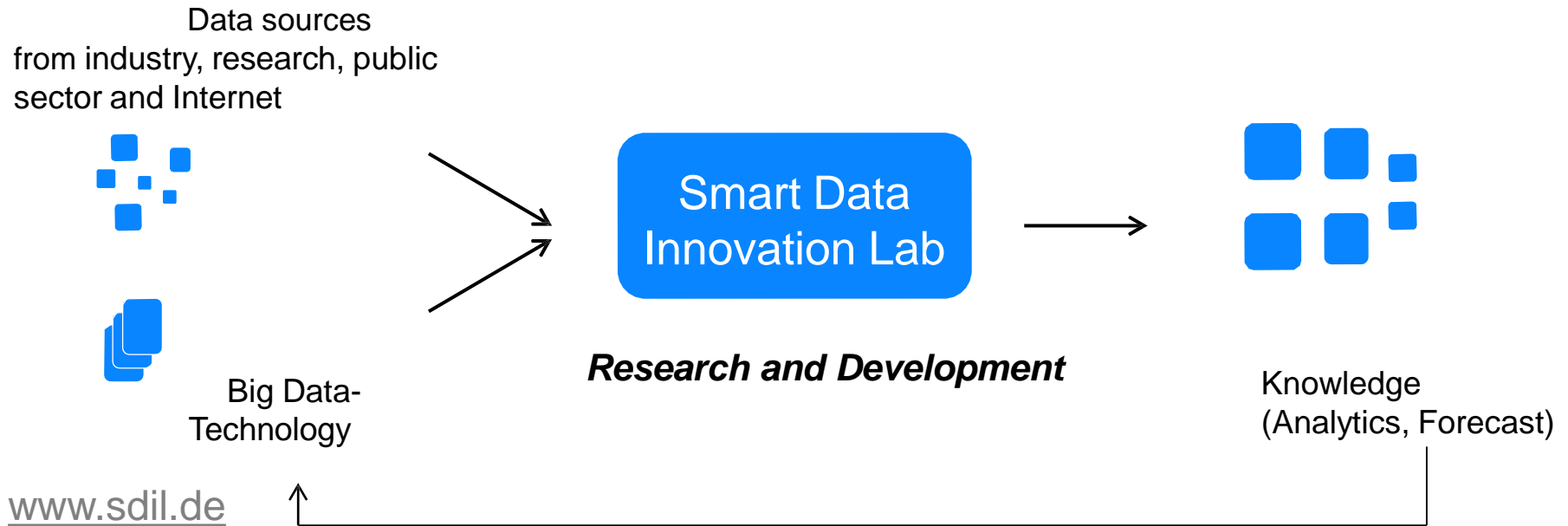


# Research and Development

# Ingredients of a Big Data Project

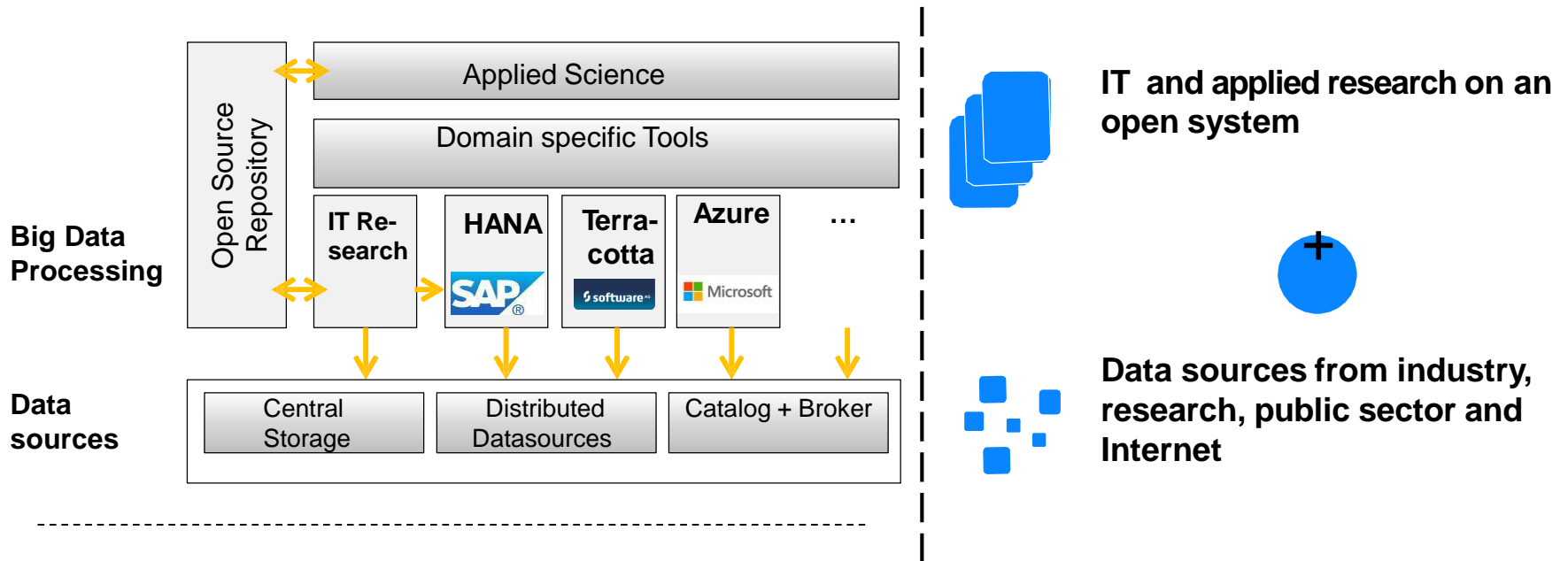
- Technology
    - Data preparation
    - Scalable processing
    - Scalable platform
- } **Cloud Computing**
- Mathematical analysis methods
    - Machine learning
    - Statistics
    - Optimization
    - ...
  - Toolset
    - Natural Language Processing
    - Image processing
    - Visualization
  - Application
    - Real-world analysis problem

# Smart Data Innovation Lab (SDIL)

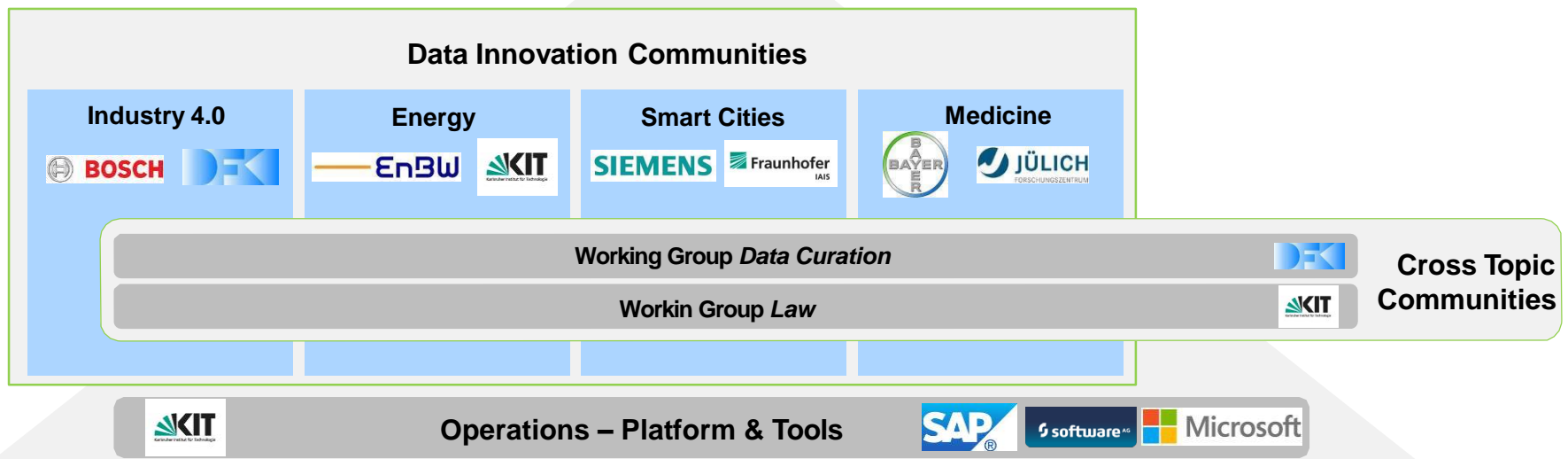


- Cooperation between industry and science to spur innovation
- Pilot R&D projects on dedicated Big Data infrastructure

# Architecture



# Application Area



# Partners



# Research and Development Areas

## Applications

- Industry 4.0
- Logistics
- Smart Grids
- Smart City
- Personalized Medicine

## Methods

- Data mining
- Machine Learning
- Statistics Analysis
- Predictive Analytics
- Tools

## Storage

- Data Warehouses
- NoSQL Databases
- Column Stores
- In memory DBs

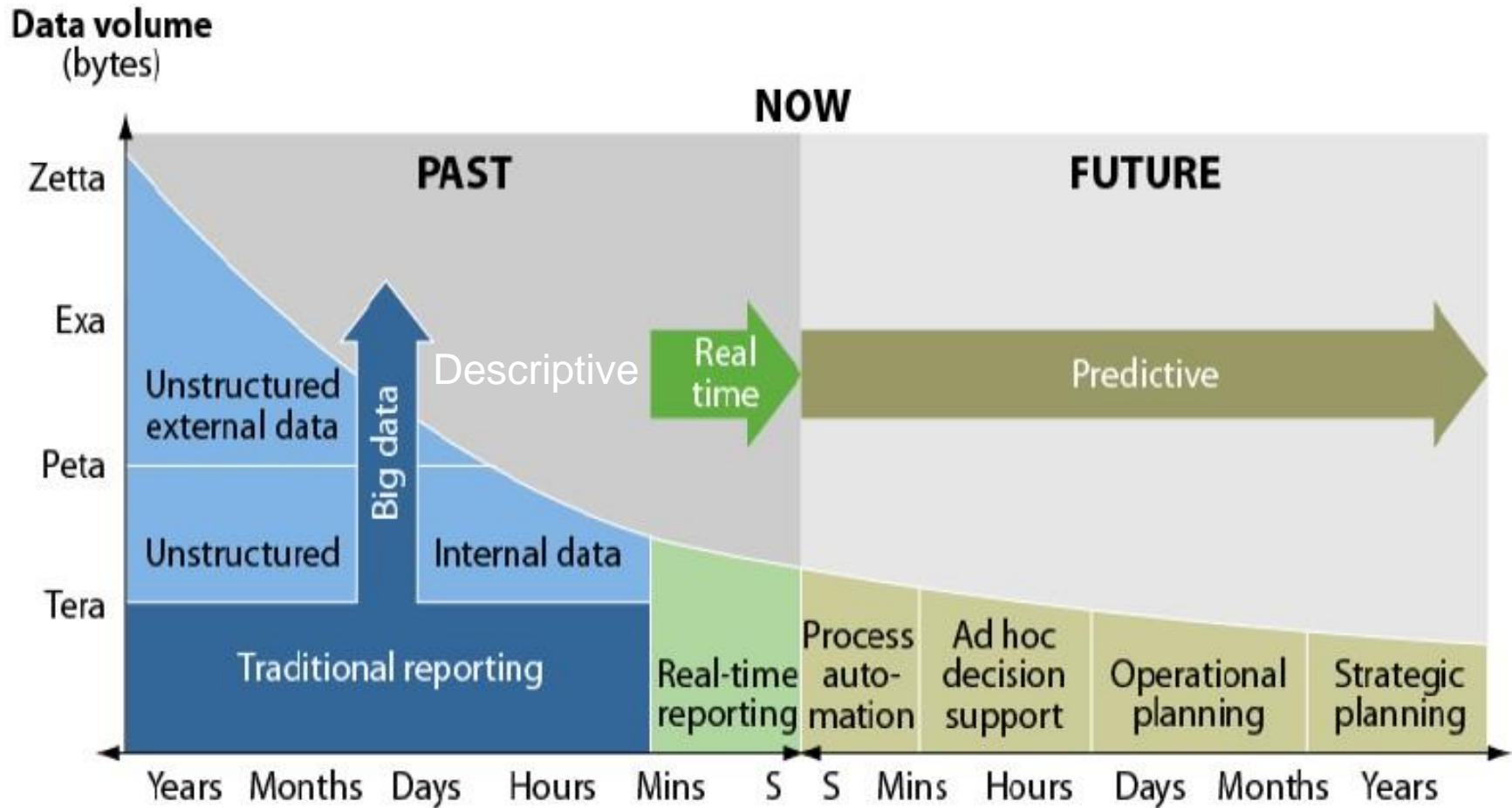
## Processing

- Hadoop Engines
- Real time Analytics
- Software Defined Data Center

## Representation

- Dashboards
- Visualization
- Rich Clients
- Collaboration Platforms

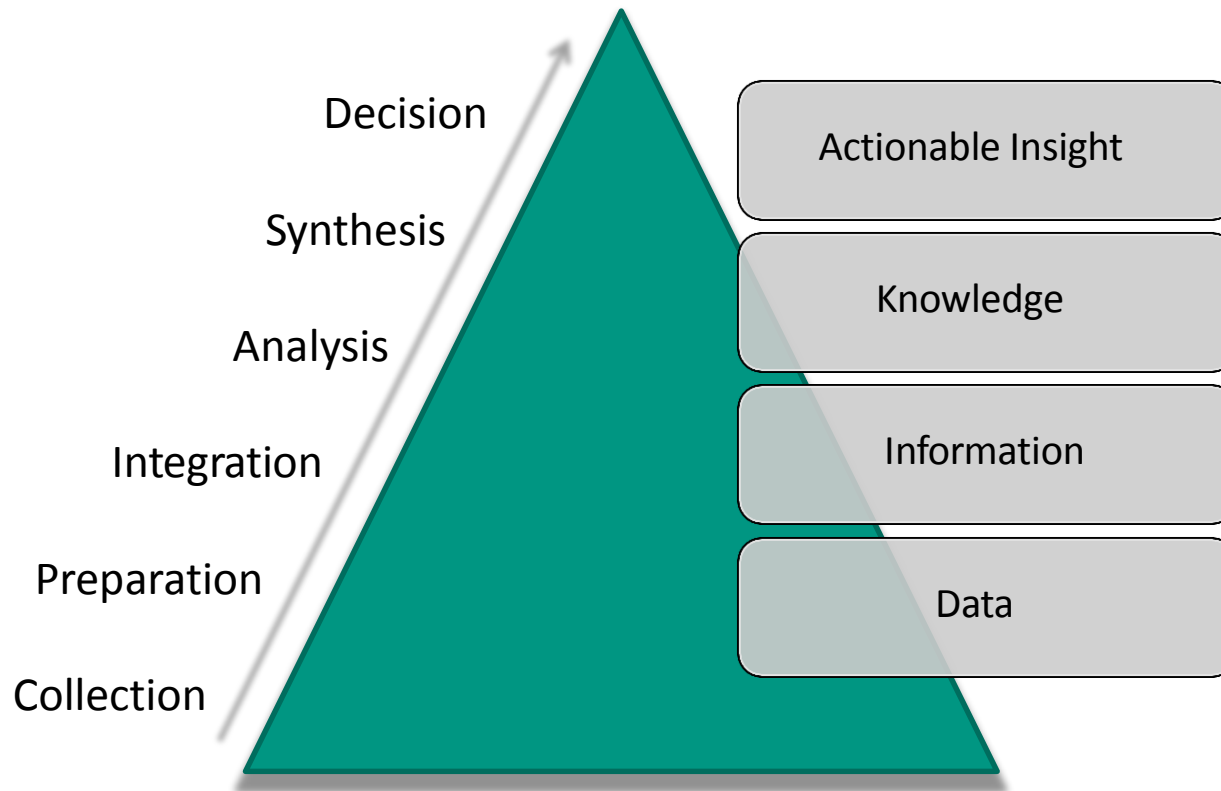
# Predictive Analytics



Source: [blue yonder](#)

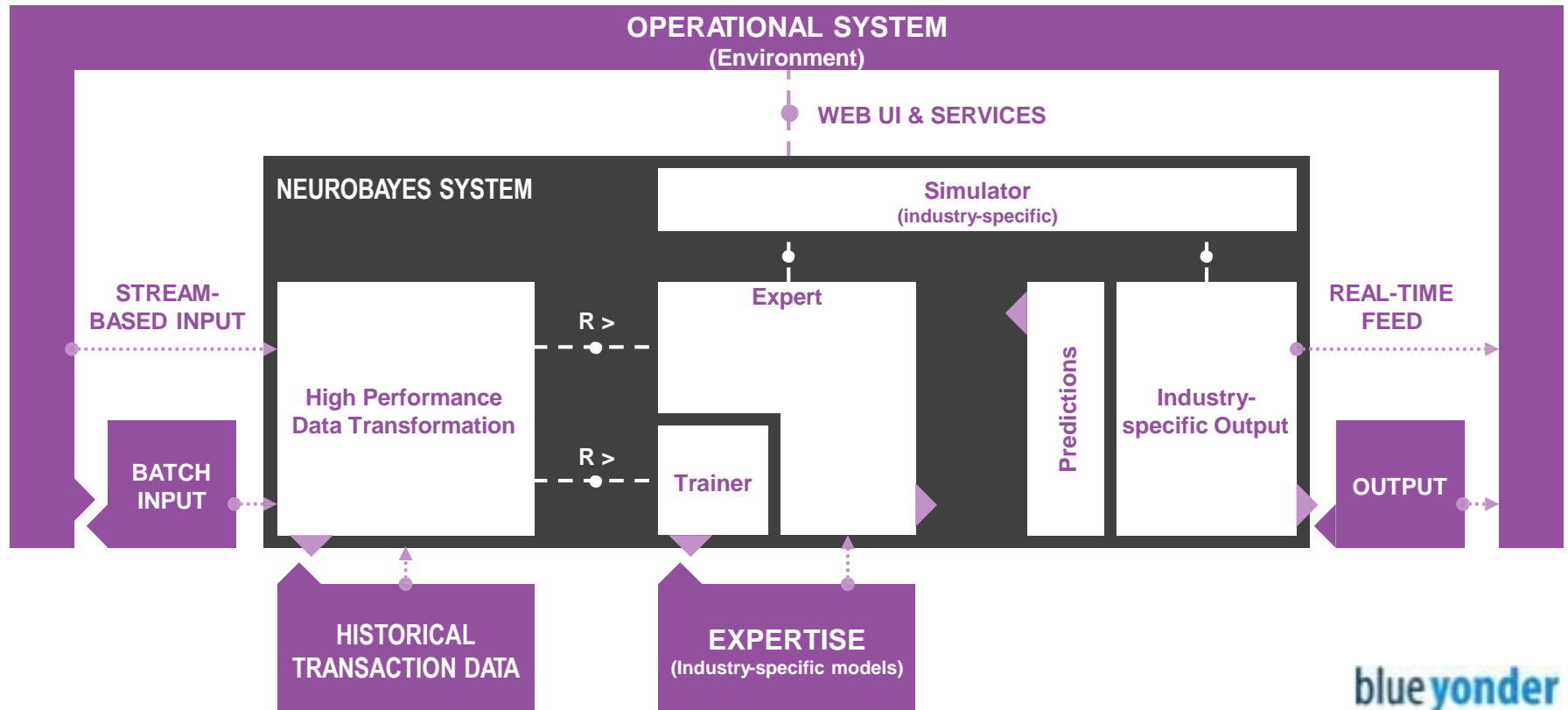


# Prescriptive Analytics



- From raw data to decision processes:
  - Data analysis takes time and often works with offline data only
  - Is it possible to improve the process and react in real time ?

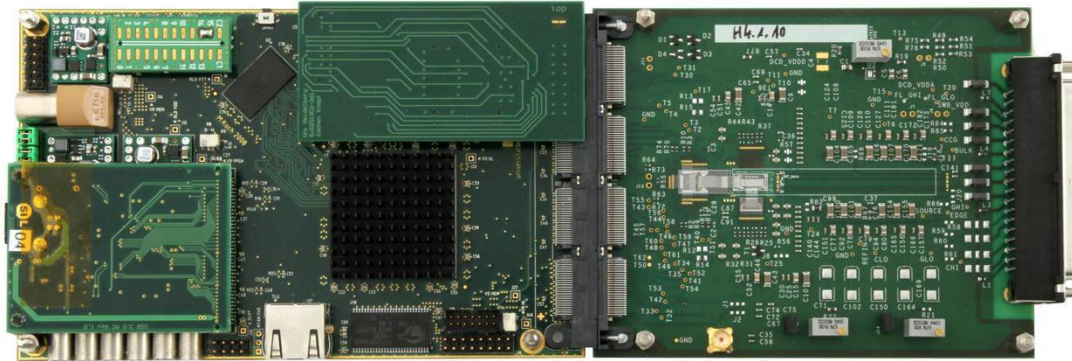
# Blue Yonder forward demand Architecture



blueyonder

- Machine learning utilizing modern in-memory database technology
- Direct integration into business processes (not just simple data-mining)

# Future: Algorithm in Hardware



- NeuroBayes machine learning algorithm on FPGA
- Field Programmable Gate Array: (XILINX Virtex6 VLX75T)
- Clock frequency: 250 MHz
- Approx. 1 decision per clock cycle (fully pipelined architecture)
- 250 million decisions per second
- Throughput: 100 Gbit/s
- Interesting for real-time investigation of online streaming data

# Summary

- Big Data depends on **scalable dynamic models** (Cloud is essential)
- Big Data is **interdisciplinary**: Computer Science, Mathematics, ...
- **Hadoop** may assume the role of the **data hub**
- Big data is more about **value** than about volume

## Contact:

[marcel.kunze@kit.edu](mailto:marcel.kunze@kit.edu)

Dr. Marcel Kunze  
Karlsruhe Institute of Technology (KIT)  
Forschungsgruppe Cloud Computing  
Steinbuch Centre for Computing  
Hermann-von-Helmholtz-Platz 1  
D-76344 Eggenstein-Leopoldshafen

Research Group Cloud Computing - Steinbuch Centre for Computing

