
**Information technology — Big data
reference architecture —**

**Part 2:
Use cases and derived requirements**

*Technologies de l'information — Architecture de référence des big
data —*

Partie 2: Cas pratiques et exigences dérivées



Copyrighted document, no reproduction or circulation

STANDARDSISO.COM: Click to view the full PDF of ISO/IEC TR 20547 WG - 2:2018

Oct 2024



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2018

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
CP 401 • Ch. de Blandonnet 8
CH-1214 Vernier, Geneva, Switzerland
Tel. +41 22 749 01 11
Fax +41 22 749 09 47
copyright@iso.org
www.iso.org

Published in Switzerland

Contents

Page

Foreword	v
Introduction	vi
1 Scope	1
2 Normative references	1
3 Terms and definitions	1
3.1 Terms defined elsewhere.....	1
3.2 Terms defined in this document.....	1
3.3 Abbreviated terms.....	1
4 Use case properties for survey	6
4.1 Overall description.....	6
4.2 Current solution.....	7
4.3 Big data characteristics.....	7
4.4 Big data science.....	7
4.5 Overall big data issues.....	8
4.6 Big data use case Template.....	8
5 Use cases summaries	9
5.1 Use case development process.....	9
5.2 Government operation.....	10
5.2.1 Use case 1: Census 2010 and 2000 — Title 13 big data.....	10
5.2.2 Use case 2: NARA Accession, Search, Retrieve, Preservation.....	10
5.2.3 Use case 3: Statistical survey response improvement.....	11
5.2.4 Use case 4: Non-Traditional Data in Statistical Survey Response Improvement (Adaptive Design).....	11
5.3 Commercial.....	12
5.3.1 Use case 5: Cloud Eco-System for Financial Industries.....	12
5.3.2 Use case 6: Mendeley — An International Network of Research.....	12
5.3.3 Use case 7: Multi-media streaming service.....	13
5.3.4 Use case 8: Web Search.....	13
5.3.5 Use case 9: Big data Business Continuity and Disaster Recovery Within a Cloud Eco-System.....	14
5.3.6 Use case 10: Cargo Shipping.....	14
5.3.7 Use case 11: Materials Data for Manufacturing.....	14
5.3.8 Use case 12: Simulation-Driven Materials Genomics.....	15
5.4 Defense.....	16
5.4.1 Use case 13: Cloud Large-Scale Geospatial Analysis and Visualization.....	16
5.4.2 Use case 14: Object Identification and Tracking from Wide-Area Large Format Imagery or Full Motion Video—Persistent Surveillance.....	16
5.4.3 Use case 15: Intelligence Data Processing and Analysis.....	17
5.5 Health care and life sciences.....	17
5.5.1 Use case 16: Electronic Medical Record Data.....	17
5.5.2 Use case 17: Pathology Imaging/Digital Pathology.....	18
5.5.3 Use case 18: Computational Bioimaging.....	18
5.5.4 Use case 19: Genomic Measurements.....	19
5.5.5 Use case 20: Comparative Analysis for Metagenomes and Genomes.....	19
5.5.6 Use case 21: Individualized Diabetes Management.....	19
5.5.7 Use case 22: Statistical Relational Artificial Intelligence for Health Care.....	20
5.5.8 Use case 23: World Population-Scale Epidemiological Study.....	20
5.5.9 Use case 24: Social Contagion Modeling for Planning, Public Health, and Disaster Management.....	21
5.5.10 Use case 25: Biodiversity and LifeWatch.....	21
5.6 Deep Learning and Social Media.....	22
5.6.1 Use case 26: Large-Scale Deep Learning.....	22

5.6.2	Use case 27: Organizing Large-Scale, Unstructured Collections of Consumer Photos	22
5.6.3	Use case 28: Truthy—Information Diffusion Research from Twitter Data	23
5.6.4	Use case 29: Crowd Sourcing in the Humanities as Source for Big and Dynamic Data	23
5.6.5	Use case 30: CINET—Cyberinfrastructure for Network (Graph) Science and Analytics	23
5.6.6	Use case 31: NIST Information Access Division — Analytic Technology Performance Measurements, Evaluations, and Standards	24
5.7	The Ecosystem for research	24
5.7.1	Use case 32: DataNet Federation Consortium	24
5.7.2	Use case 33: The Discinnet Process	25
5.7.3	Use case 34: Semantic Graph Search on Scientific Chemical and Text-Based Data	25
5.7.4	Use case 35: Light Source Beamlines	26
5.8	Astronomy and physics	26
5.8.1	Use case 36: Catalina Real-Time Transient Survey: A Digital, Panoramic, Synoptic Sky Survey	26
5.8.2	Use case 37: DOE Extreme Data from Cosmological Sky Survey and Simulations	27
5.8.3	Use case 38: Large Survey Data for Cosmology	27
5.8.4	Use case 39: Particle Physics—Analysis of Large Hadron Collider Data: Discovery of Higgs Particle	28
5.8.5	Use case 40: Belle II High Energy Physics Experiment	29
5.9	Earth, environmental, and polar science	29
5.9.1	Use case 41: European Incoherent Scatter Scientific Association 3D Incoherent Scatter Radar System	29
5.9.2	Use case 42: Common Operations of Environmental Research Infrastructure	30
5.9.3	Use case 43: Radar Data Analysis for the Center for Remote Sensing of Ice Sheets	31
5.9.4	Use case 44: Unmanned Air Vehicle Synthetic Aperture Radar (UAVSAR) Data Processing, Data Product Delivery, and Data Services	31
5.9.5	Use case 45: NASA Langley Research Center/ Goddard Space Flight Center iRODS Federation Test Bed	32
5.9.6	Use case 46: MERRA Analytic Services (MERRA/AS)	32
5.9.7	Use case 47: Atmospheric Turbulence – Event Discovery and Predictive Analytics	32
5.9.8	Use case 48: Climate Studies Using the Community Earth System Model at the U.S. Department of Energy (DOE) NERSC Center	33
5.9.9	Use case 49: DOE Biological and Environmental Research (BER) Subsurface Biogeochemistry Scientific Focus Area	33
5.9.10	Use case 50: DOE BER AmeriFlux and FLUXNET Networks	34
5.10	Energy	34
5.10.1	Use case 51: Consumption Forecasting in Smart Grids	34
5.10.2	Use case 52: Home Energy Management System	34
6	Use cases derived technical considerations	35
6.1	Use case specific technical considerations	35
6.2	Summary of requirements analysis	35
6.3	Features of use cases	37
Annex A Submitted use case studies		40
Annex B Summary of Key Properties		197
Annex C Use case technical considerations summary		207
Annex D Use case detail technical considerations		225
Bibliography		252

Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work. In the field of information technology, ISO and IEC have established a joint technical committee, ISO/IEC JTC 1.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights. Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see www.iso.org/patents).

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation on the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see the following URL: www.iso.org/iso/foreword.html.

This document was prepared by Technical Committee ISO/IEC JTC 1, *Information Technology*.

A list of all parts in the ISO/IEC 20547-series can be found on the ISO website.

Introduction

This document is focuses on forming a community of interest from industry, academia, and government, with the goal of developing a consensus list of big data technical considerations across all stakeholders. This included gathering and understanding various examples of use cases from diversified areas (i.e., application domains). To achieve this goal, the following tasks were done:

- gathered input from all stakeholders regarding big data technical considerations;
- analyzed and prioritized a list of challenging use case specific technical considerations that may delay or prevent adoption of big data deployment;
- developed a comprehensive list of generalized big data technical considerations for ISO/IEC 20547-3, *Information technology – Big data reference architecture - Part 3: Reference architecture*; and
- documented the findings in this document.

Information technology — Big data reference architecture —

Part 2: Use cases and derived requirements

1 Scope

This document provides examples of big data use cases with application domains and technical considerations derived from the contributed use cases.

2 Normative references

The following documents, in whole or in part, are normatively referenced in this document and are indispensable for its application. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 20546 *Information technology — Big data — Definition and vocabulary*

3 Terms and definitions

For the purposes of this document, the terms and definitions given in ISO/IEC 20546 and the following apply.

ISO and IEC maintain terminological databases for use in standardization at the following addresses:

- IEC Electropedia: available at <http://www.electropedia.org/>
- ISO Online browsing platform: available at <http://www.iso.org/obp>

3.1 Terms defined elsewhere

None.

3.2 Terms defined in this document

3.2.1 use case

typical application stated at a high level for the purposes of extracting technical considerations or comparing usages across fields

3.3 Abbreviated terms

2D	two-Dimensional
3D	three-Dimensional
6D	six-Dimensional
AOD	Analysis Object Data

API	Application Programming Interface
ASDC	Atmospheric Science Data Center
ASTM	American Society for Testing and Materials
AWS	Amazon Web Services
BC/DR	Business Continuity and Disaster Recovery
BD	Big data
BER	Biological and Environmental Research
BNL	Brookhaven National Laboratory
CAaaS	Climate Analytics as a Service
CADRG	Compressed ARC Digitized Raster Graphics
CBSP	Cloud Brokerage Service Provider
CERES	Clouds and Earth's Radiant Energy System
CERN	European Organization for Nuclear Research
CESM	Community Earth System Model
CFTC	U.S. Commodity Futures Trading Commission
CIA	Confidentiality, Integrity, and Availability
CINET	Cyberinfrastructure for Network (Graph) Science and Analytics
CMIP	Coupled Model Intercomparison Project
CMIP5	Climate Model Intercomparison Project
CMS	Compact Muon Solenoid
COSO	Committee of Sponsoring Organizations
CPU	Central Processing Unit
CReSIS	Center for Remote Sensing of Ice Sheets
CRTS	Catalina Real-Time Transient Survey
CSP	Cloud Service Provider
CSS	Catalina Sky Survey proper
CV	Controlled Vocabulary
DFC	DataNet Federation Consortium
DHTC	Distributed High Throughput Computing
DNA	DeoxyriboNucleic Acid
DOE	U.S. Department of Energy

DOJ	U.S. Department of Justice
DPO	Data Products Online
EBAF–TOA	Energy Balanced and Filled–Top of Atmosphere
EC2	Elastic Compute Cloud
EDT	Enterprise Data Trust
EHR	Electronic Health Record
EMR	Electronic Medical Record
EMSO	European Multidisciplinary Seafloor and Water Column Observatory
ENVRI	Common Operations of Environmental Research Infrastructures
ENVRI RM	ENVRI Reference Model
EPOS	European Plate Observing System
ESFRI	European Strategy Forum on Research Infrastructures
ESG	Earth System Grid
ESGF	Earth System Grid Federation
FDIC	U.S. Federal Deposit Insurance Corporation
FI	Financial Industries
FLUXNET	Flux Tower Network
FMV	Full Motion Video
FNAL	Fermi National Accelerator Laboratory
GAAP	U.S. Generally Accepted Accounting Principles
GB	Giga Byte
GCM	General Circulation Model
GEOS-5	Goddard Earth Observing System version 5
GeoTiff	Geo Tagged Image File Format
GEWaSC	Genome-Enabled Watershed Simulation Capability
GHG	Green House Gas
GMAO	Global Modeling and Assimilation Office
GPFS	General Parallel File System
GPS	Global Positioning System
GPU	Graphics Processing Unit
GRC	Governance, Risk management, and Compliance

GSFC	Goddard Space Flight Center
HDF5	Hierarchical Data Format
HDFS	Hadoop Distributed File System
HPC	High-Performance Computing
HTC	High-Throughput Computing
HVS	Hosted Virtual Server
I/O	Input Output
IaaS	Infrastructure as a Service
IAGOS	In-service Aircraft for a Global Observing System
ICD	International Classification of Diseases
ICOS	Integrated Carbon Observation System
IMG	Integrated Microbial Genomes
INPC	Indiana Network for Patient Care
IPCC	Intergovernmental Panel on Climate Change
iRODS	Integrated Rule-Oriented Data System
ISACA	International Society of Auditors and Computer Analysts
isc2	International Security Computer and Systems Auditors
ISO	International Organization for Standardization
ITIL	Information Technology Infrastructure Library
JGI	Joint Genome Institute
KML	Keyhole Markup Language
kWh	kilowatt-hour
LaRC	Langley Research Center
LBNL	Lawrence Berkeley National Laboratory
LDA	latent Dirichlet allocation
LHC	Large Hadron Collider
LPL	Lunar and Planetary Laboratory
LSST	Large Synoptic Survey Telescope
MERRA	Modern Era Retrospective Analysis for Research and Applications
MERRA/AS	MERRA Analytic Services
MPI	Message Passing Interface

MRI	Magnetic Resonance Imaging
NARA	National Archives and Records Administration
NARR	North American Regional Reanalysis
NaaS	Network as a Service
NASA	National Aeronautics and Space Administration
NCAR	National Center for Atmospheric Research
NCBI	National Center for Biotechnology Information
NCCS	NASA Center for Climate Simulation
NERSC	National Energy Research Scientific Computing Center
NetCDF	Network Common Data Form
NEX	NASA Earth Exchange
NFS	Network File System
NIKE	NIST Integrated Knowledge Editorial Net
NIST	National Institute of Standards and Technology
NITF	National Imagery Transmission Format
NLP	Natural Language Processing
NRT	Near Real Time
NSF	National Science Foundation
ODP	Open Distributed Processing
OGC	Open Geospatial Consortium
PB	PetaByte
PCA	Principal Component Analysis
PCAOB	Public Company Accounting and Oversight Board
PID	persistent identification
PII	Personally Identifiable Information
PNNL	Pacific Northwest National Laboratory
RDBMS	relational database management system
RDF	Resource Description Framework
RECOVER	Rehabilitation Capability Convergence for Ecosystem Recovery
ROI	return on investment
RPI	Repeat Pass Interferometry

RPO	Recovery Point Objective
RTO	Response Time Objective
SAN	storage area network
SAR	Synthetic Aperture Radar
SDN	software-defined networking
SIOS	Svalbard Integrated Arctic Earth Observing System
SPADE	Support for Provenance Auditing in Distributed Environments
SSH	Secure Shell
SSO	Single Sign-On
TB	TeraByte
tf-idf	term frequency–inverse document frequency
UA	University of Arizona
UAVSAR	Unmanned Air Vehicle Synthetic Aperture Radar
UC	Use Case
UI	User Interface
UPS	United Parcel Service
UQ	Uncertainty Quantification
VASP	Vienna Ab initio Simulation Package
vCDS	virtual Climate Data Server
VO	Virtual Observatory
VOIP	Voice over IP
WALF	Wide Area Large Format Imagery
WLCG	Worldwide LHC Computing Grid
XBRL	extensible Business Related Markup Language
XML	Extensible Markup Language
ZTF	Zwicky Transient Factory

4 Use case properties for survey

4.1 Overall description

- **Use case title:** Title provided by the use case author
- **Vertical (area):** Intended to categorize the use cases. However, an ontology was not created prior to the use case submissions so this field was not used in the use case compilation.

- **Author/company/email:** Name, company, and email (if provided) of the person(s) submitting the use case
- **Actors/ stakeholders and their roles and responsibilities:** Description of the players and their roles in the use case
- **Goals:** Objectives of the use case
- **Use case description:** Brief description of the use case

4.2 Current solution

Current solutions describe current approach to processing big data at the hardware and software infrastructure and analytics level.

- **Compute (System):** Computing component of the data analysis system
- **Storage:** Storage component of the data analysis system
- **Networking:** Networking component of the data analysis system
- **Software:** Software component of the data analysis system

4.3 Big data characteristics

Big data Characteristics describe the properties of the (raw) data including the four major 'V's' of big data.

- **Data source:** The origin of data, which could be from instruments, Internet of Things, Web, Surveys, Commercial activity, or from simulations. The source(s) can be distributed, centralized, local, or remote.
- **Data destination:** If data transformed in use case, where the final results end up.
- **Volume:** The characteristic of datasets that is most associated with big data. Volume represents the extensive amount of data available for analysis to extract valuable information. The assumption that you can extract the most value by analysing as much of the volume of data as possible was one of the primary drivers for the creation of the new scaling technologies.
- **Velocity:** The rate of flow at which the data is created, stored, analysed, or visualized. Big data velocity means a large quantity of data needs to be processed in a short amount of time. Dealing with high velocity data is commonly referred to as techniques for streaming data.
- **Variety:** The need to analyse data from a number of domains and a number of data types. The variety of data was handled through transformations or pre-analytics to extract features that would allow integration with other data. The wider range of data formats, logical models, timescales, and semantics, which is desirable to be used in analytics, complicates the integration of the variety of data. Metadata is increasingly used to aid in the integration.
- **Variability:** Changes in data rate, format/structure, semantics, and/or quality that impact the supported application, analytic, or problem. Impacts can include the need to refactor architectures, interfaces, processing/algorithms, integration/fusion, storage, applicability, or use of the data.

4.4 Big data science

Big data science describes the high level aspects of the data analysis process.

- **Veracity and data quality:** This covers the completeness and accuracy of the data with respect to semantic content as well as syntactical quality of data (such as presence of missing fields or incorrect values).

- **Visualization:** Refers to the way data is viewed by an analyst making decisions based on the data. Typically, visualization is the final stage of a technical data analysis pipeline and follows the data analytics stage.
- **Data types:** Refers to the style of data such as structured, unstructured, images (e.g., pixels), text (e.g., characters), gene sequences, and numerical.
- **Metadata:** Comments on quality and richness of metadata.
- **Curation and governance:** Comment on process to ensure good data quality and who is responsible.

NOTE The use case template has a separate item to describe security and privacy issues.

- **Data analytics:** Refers broadly to tools and algorithms used in processing the data at any stage including the data to information or knowledge to wisdom stages, as well as the information to knowledge stage.

4.5 Overall big data issues

- **Other big data issues:** Did we miss something important that your use case highlights? Your chance to address questions which we should have asked.
- **User Interface and mobile access issues:** Refers to issues in accessing or generating big data from clients including smart phones and tablets.
- **List key features and related use cases:** Put use case in context of related use cases. What features generalize and what are idiosyncratic to this use case.
- **Project future:** How do you expect application, and approach (hardware, software, analytics) to change in future?
- **More project information (URLs):** Put a collection of useful links.

4.6 Big data use case Template

This clause provides one blank use case template. The below blank use case template was used for the purpose of capturing use cases to derived technical consideration.

NOTE The terms used in this template may or may not match with ISO/IEC 20546 and other parts of the ISO/IEC 20547-series.

Use case title		
Vertical (area)		
Author/company/email		
Actors/stakeholders and their roles and responsibilities		
Goals		
Use case description		
Current solutions	Compute(System)	
	Storage	
	Networking	
	Software	

Big data characteristics	Data source (distributed/centralized)	
	Volume (size)	
	Velocity (e.g. real time)	
	Variety (multiple datasets, mashup)	
	Variability (rate of change)	
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	
	Visualization	
	Data quality (syntax)	
	Data types	
	Data analytics	
Big data specific challenges (Gaps)		
Big data specific challenges in mobility		
Security and privacy technical considerations		
Highlight issues for generalizing this Use case (e.g. for ref. architecture)		
More information (URLs)		
NOTE <additional comments>		

5 Use cases summaries

5.1 Use case development process

A use case is a typical application stated at a high level for the purposes of extracting technical considerations or comparing usages across fields. In order to develop a consensus list of big data technical considerations across all stakeholders, publicly available information was collected for various big data architectures. After collection of use cases, application domains were identified to better organize the collection of use cases.

NOTE 1 The list of application domains reflects the use cases submitted and is not intended to be exhaustive.

The nine application domains were as follows:

- **Government operation** (4): National Archives and Records Administration, Census Bureau;
- **Commercial** (8): Finance in Cloud, Cloud Backup, Citations, Multi-media streaming, Web Search, Digital Materials, Cargo Shipping;
- **Defense** (3): Sensors, Image Surveillance, Situation Assessment ;
- **Healthcare and life sciences** (10): Medical Records, Graph and Probabilistic Analysis, Pathology, Bioimaging, Genomics, Epidemiology, People Activity Models, Biodiversity;
- **Deep learning and social media** (6) Self-driving cars, Geolocate Images, SNS, Crowd Sourcing, Network Science, Benchmark Datasets;

- **Ecosystem for research** (4): Metadata, Collaboration, Language Translation, Light Source Experiments;
- **Astronomy and physics** (5): Sky Surveys (and comparisons to simulation), LHC at CERN, Belle Accelerator II;
- **Earth, environmental, and polar science** (10): Radar Scattering in Atmosphere, Earthquake, Ocean, Earth Observation, Ice Sheet Radar Scattering, Earth Radar Mapping, Climate Simulation Datasets, Atmospheric Turbulence Identification, Subsurface Biogeochemistry (microbes to watersheds), Gas Sensors;
- **Energy** (2): Smart Grid, Home energy management.

NOTE 2 The template was valuable for gathering consistent information to develop supporting analysis and comparison of the use cases. However, varied levels of detail and quantitative or qualitative information was received for each use case template section. For some application domains, several similar big data use cases are presented, providing a more complete view of big data technical considerations within that application domain.

The examples of use cases are presented in this clause with the information originally submitted. The original content (See [Annex A](#)) has not been modified.

NOTE 3 Specific vendor solutions and technologies are mentioned in the use cases. However, the listing of these solutions and technologies does not constitute endorsement from the JTC 1 WG 9.

The use cases are numbered sequentially to facilitate cross-referencing between the use case summaries presented in this clause, the original use cases ([Annex A](#)), and the use case summary tables ([Annexes B, C, and D](#)).

5.2 Government operation

5.2.1 Use case 1: Census 2010 and 2000 — Title 13 big data

Application:

Census 2010 and 2000—Title 13 data must be preserved for several decades so they can be accessed and analyzed after 75 years. Data must be maintained 'as-is' with no access and no data analytics for 75 years, preserved at the bit level, and curated, which may include format transformation. Access and analytics must be provided after 75 years. Title 13 of the U.S. Code authorizes the U.S. Census Bureau to collect and preserve census related data and guarantees that individual and industry-specific data are protected.

Current approach:

The dataset contains 380 TB of scanned documents.

Future:

Future data scenarios and applications were not expressed for this use case.

5.2.2 Use case 2: NARA Accession, Search, Retrieve, Preservation

Application:

This area comprises accession, search, retrieval, and long-term preservation of government data.

Current approach:

The data are currently handled as follows.

- Get physical and legal custody of the data.
- Pre-process data for conducting virus scans, identifying file format, and removing empty files.

- Index the data.
- Categorize records (e.g., sensitive, non-sensitive, privacy data).
- Transform old file formats to modern formats.
- Conduct e-discovery.
- Search and retrieve to respond to special requests..
- Search and retrieve public records by public users.

Hundreds of TBs are stored centrally in commercial databases supported by custom software and commercial search products.

Future:

Federal agencies possess many distributed data sources, which currently must be transferred to centralized storage. In the future, those data sources may reside in multiple cloud environments. In this case, physical custody should avoid transferring big data from cloud to cloud or from cloud to data center.

5.2.3 Use case 3: Statistical survey response improvement

Application:

Survey costs are increasing as survey responses decline. The goal of this work is to increase the quality — and reduce the cost — of field surveys by using advanced ‘recommendation system techniques.’ These techniques are open and scientifically objective, using data mashed up from several sources and also historical survey para-data (i.e., administrative data about the survey.)

Current approach:

This use case handles about a PB of data coming from surveys and other government administrative sources. Data can be streamed. During the decennial census, approximately 150 million records transmitted as field data are streamed continuously. All data must be both confidential and secure. All processes must be auditable for security and confidentiality as required by various legal statutes. Data quality should be high and statistically checked for accuracy and reliability throughout the collection process. Solution information is described in [Clause A.1.3](#)

Future:

Improved recommendation systems are needed similar to those used in e-commerce (e.g., similar to the use case [5.3.3](#) that reduce costs and improve quality, while providing confidentiality safeguards that are reliable and publicly auditable. Data visualization is useful for data review, operational activity, and general analysis. The system continues to evolve and incorporate important features such as mobile access.

5.2.4 Use case 4: Non-Traditional Data in Statistical Survey Response Improvement (Adaptive Design)

Application:

Survey costs are increasing as survey response declines. This use case has goals similar to those of the Statistical Survey Response Improvement use case (see [Clause 5.2.3](#)). However, this case involves non-traditional commercial and public data sources from the web, wireless communication, and electronic transactions mashed up analytically with traditional surveys. The purpose of the mashup is to improve statistics for small area geographies and new measures, as well as the timeliness of released statistics.

Current approach:

Data from a range of sources are integrated including survey data, other government administrative data, web scrapped data, wireless data, e-transaction data, possibly social media data, and positioning

data from various sources. Software, visualization, and data characteristics are similar to those in the Statistical Survey Response Improvement use case.

Future:

Analytics need to be developed that give more detailed statistical estimations, on a more near real-time basis, for less cost. The reliability of estimated statistics from such mashed up sources still must be evaluated.

5.3 Commercial

5.3.1 Use case 5: Cloud Eco-System for Financial Industries

Application:

Use of cloud (e.g., big data) technologies needs to be extended in financial industries (i.e., banking, securities and investments, insurance) transacting business within the U.S.

Current approach:

The financial industry is already using big data for fraud detection, risk analysis, assessments, as well as improving their knowledge and understanding of customers. At the same time, the industry is still using traditional client/server/data warehouse/relational database management system (RDBMS) for the handling, processing, storage, and archival of financial data. Real-time data and analysis are important in these applications.

Future:

Security, privacy, and regulation must be addressed. For example, the financial industry must examine SEC-mandated use of XBRL (extensible business-related markup language) and use of other cloud functions.

5.3.2 Use case 6: Mendeley — An International Network of Research

Application:

Mendeley has built a database of research documents and facilitates the creation of shared bibliographies. Mendeley collects and uses the information about research reading patterns and other activities conducted via their software to build more efficient literature discovery and analysis tools. Text mining and classification systems enable automatic recommendation of relevant research, improving research teams' performance and cost-efficiency, particularly those engaged in curation of literature on a particular subject.

Current approach:

Data size is presently 15 TB and growing at a rate of about 1 TB per month. Solution information is described in [Clause A.2.2](#). The database uses standard libraries for machine learning and analytics, latent Dirichlet allocation (LDA, a generative probabilistic model for discrete data collection), and custom-built reporting tools for aggregating readership and social activities for each document.

Future:

Currently big data storage batch jobs are scheduled daily, but work has begun on real-time recommendation. The database contains approximately 400 million documents and roughly 80 million unique documents, and receives 500,000 to 700,000 new uploads on a weekday. Thus a major challenge is clustering matching documents together in a computationally efficient way (i.e., scalable and parallelized) when they are uploaded from different sources and have been slightly modified via third-party annotation tools or publisher watermarks and cover pages.

5.3.3 Use case 7: Multi-media streaming service

Application:

This allows streaming of user-selected movies to satisfy multiple objectives (for different stakeholders)—but with a focus on retaining subscribers. The company needs to find the best possible ordering of a set of videos for a user (e.g., household) within a given context in real time, with the objective of maximizing movie consumption. Recommendation systems and streaming video delivery are core Netflix technologies. Recommendation systems are always personalized and use logistic/linear regression, elastic nets, matrix factorization, clustering, LDA, association rules, gradient-boosted decision trees, and other tools. Digital movies are stored in the cloud with metadata, along with individual user profiles and rankings for small fraction of movies. The current system uses multiple criteria: a content-based recommendation system, a user-based recommendation system, and diversity. Algorithms are continuously refined with A/B testing (i.e., two-variable randomized experiments used in online marketing).

Current approach:

It held a competition for the best collaborative filtering algorithm to predict user ratings for films—the purpose of which was to improve ratings by 10 %. The winning system combined over 100 different algorithms. Solution information is described in [Clause A.2.3](#). Business initiatives have been used to increase viewership.

Future:

Streaming video is a very competitive business. It needs to be aware of other companies and trends in both content (e.g., which movies are popular) and big data technology.

5.3.4 Use case 8: Web Search

Application:

A web search function returns results in $\sim 0,1$ s based on search terms with an average of three words. It is important to maximize quantities such as “precision@10” for the number of highly accurate/appropriate responses in the top 10 ranked results.

Current approach:

The current approach uses the following steps.

- Crawl the web.
- Pre-process data to identify what is searchable (words, positions).
- Form an inverted index, which maps words to their locations in documents
- Rank the relevance of documents using the PageRank algorithm.
- Employ advertising technology, e.g., using reverse engineering to identify ranking models—or preventing reverse engineering.
- Cluster documents into topics (as in Google News).
- Update results efficiently.

Modern clouds and technologies such as Map/Reduce have been heavily influenced by this application, which now comprises ~ 45 billion web pages total.

Future:

Web search is a very competitive field, so continuous innovation is needed. Two important innovation areas are addressing the growing segment of mobile clients, and increasing sophistication of responses

and layout to maximize the total benefit of clients, advertisers, and the search company. The “deep web” (content not indexed by standard search engines, buried behind user interfaces to databases, etc.) and multimedia searches are also of increasing importance. Each day, 500 million photos are uploaded, and each minute, 100 hours of video are uploaded to YouTube.

5.3.5 Use case 9: Big data Business Continuity and Disaster Recovery Within a Cloud Eco-System

Application:

Business Continuity and Disaster Recovery (BC/DR) needs to consider the role that four overlaying and interdependent forces will play in ensuring a workable solution to an entity's business continuity plan and requisite disaster recovery strategy. The four areas are people (i.e., resources), processes (e.g., time/cost/return on investment [ROI]), technology (e.g., various operating systems, platforms, and footprints), and governance (e.g., subject to various and multiple regulatory agencies).

Current approach:

Data replication services are provided through cloud ecosystems, incorporating IaaS and supported by Tier 3 data centers. Replication is different from backup and only moves the changes that took place since the previous replication, including block-level changes. The replication can be done quickly—with a five-second window—while the data are replicated every four hours. This data snapshot is retained for seven business days, or longer if necessary. Replicated data can be moved to a failover center (i.e., a backup system) to satisfy an organization's recovery point objectives (RPO) and recovery time objectives (RTO). Related solution information is described in [Annex A](#). Data sizes range from terabytes to petabytes.

Future:

Migrating from a primary site to either a replication site or a backup site is not yet fully automated. The goal is to enable the user to automatically initiate the failover sequence. Both organizations must know which servers have to be restored and what the dependencies and inter-dependencies are between the primary site servers and replication and/or backup site servers. This knowledge requires continuous monitoring of both.

5.3.6 Use case 10: Cargo Shipping

Application:

Delivery companies need optimal means of monitoring and tracking cargo.

Current approach:

Information is updated only when items are checked with a bar code scanner, which sends data to the central server. An item's location is not currently displayed in real time.

Future:

Tracking items in real time is feasible through the Internet of Things application, in which objects are given unique identifiers and capability to transfer data automatically, i.e., without human interaction. A new aspect will be the item's status condition, including sensor information, global positioning system (GPS) coordinates, and a unique identification schema based upon standards (ISO/IEC 29161:2016 *Information technology – Data structure – Unique identification for the Internet of Things*).

5.3.7 Use case 11: Materials Data for Manufacturing

Application:

Every physical product is made from a material that has been selected for its properties, cost, and availability. This translates into hundreds of billions of dollars of material decisions made every year. However, the adoption of new materials normally takes decades (usually two to three decades) rather

than a small number of years, in part because data on new materials are not easily available. To speed adoption time, accessibility, quality, and usability must be broadened, and proprietary barriers to sharing materials data must be overcome. Sufficiently large repositories of materials data are needed to support discovery.

Current approach:

Decisions about materials usage are currently unnecessarily conservative, are often based on older rather than newer materials research and development data, and do not take advantage of advances in modeling and simulation.

Future:

Materials informatics is an area in which the new tools of data science can have a major impact by predicting the performance of real materials (in gram to ton quantities) starting at the atomistic, nanometer, and/or micrometer levels of description. The following efforts are needed to support this area.

- Establish materials data repositories, beyond the existing ones, that focus on fundamental data.
- Develop internationally accepted data recording standards that can be used by a very diverse materials community, including developers of materials test standards (e.g., American Society for Testing and Materials [ASTM] International and ISO), testing companies, materials producers, and research and development labs.
- Develop tools and procedures to help organizations that need to deposit proprietary materials in data repositories to mask proprietary information while maintaining the data's usability.
- Develop multi-variable materials data visualization tools in which the number of variables can be quite high.

5.3.8 Use case 12: Simulation-Driven Materials Genomics

Application:

Massive simulations spanning wide spaces of possible design lead to innovative battery technologies. Systematic computational studies are being conducted to examine innovation possibilities in photovoltaics. Search and simulation is the basis for rational design of materials. All these require management of simulation results contributing to the materials genome.

Current approach:

Survey results are produced using PyMatGen, FireWorks, Vienna Ab initio Simulation Package (VASP), ABINIT, NWChem, BerkeleyGW, and varied materials community codes running on large supercomputers, such as the Hopper at the National Energy Research Scientific Computing Center (NERSC), a 150,000-core machine that produces high-resolution simulations.

Future:

Large-scale computing and flexible data methods at scale for messy data are needed for simulation science. The advancement of goal-driven thinking in materials design requires machine learning and knowledge systems that integrate data from publications, experiments, and simulations. Other needs include scalable key-value and object store databases; the current 100 TB of data will grow to 500 TB over the next five years.

5.4 Defense

5.4.1 Use case 13: Cloud Large-Scale Geospatial Analysis and Visualization

Application:

Large-scale geospatial data analysis and visualization must be supported. As the number of geospatially aware sensors and geospatially tagged data sources increase, the volume of geospatial data requiring complex analysis and visualization is growing exponentially.

Current approach:

Traditional geographic information systems (GISs) are generally capable of analyzing millions of objects and visualizing thousands. Data types include imagery (various formats such as National Imagery Transmission Format (NITF), Geo Tagged Image File Format (GeoTiff), and Compressed ARC Digitized Raster Graphics (CADRG)) and vector (various formats such as shape files, KML [Keyhole Markup Language], and text streams). Object types include points, lines, areas, polylines, circles, and ellipses. Image registration—transforming various data into one system—requires data and sensor accuracy. Analytics include principal component analysis (PCA) and independent component analysis (ICA) and consider closest point of approach, deviation from route, and point density over time. Software information is described in [A.3.1](#).

Future:

Today's intelligence systems often contain trillions of geospatial objects and must visualize and interact with millions of objects. Critical issues are indexing, retrieval and distributed analysis (note that geospatial data requires unique approaches to indexing and distributed analysis); visualization generation and transmission; and visualization of data at the end of low-bandwidth wireless connections. Data are sensitive and must be completely secure in transit and at rest (particularly on handhelds).

5.4.2 Use case 14: Object Identification and Tracking from Wide-Area Large Format Imagery or Full Motion Video—Persistent Surveillance

Application:

Persistent surveillance sensors can easily collect PB of imagery data in the space of a few hours. The data should be reduced to a set of geospatial objects (e.g., points, tracks) that can be easily integrated with other data to form a common operational picture. Typical processing involves extracting and tracking entities (e.g., vehicles, people, and packages) over time from the raw image data.

Current approach:

It is not feasible for humans to process these data for either alerting or tracking purposes. The data need to be processed close to the sensor, which is likely forward-deployed since it is too large to be easily transmitted. Typical object extraction systems are currently small (e.g., 1 to 20 nodes) graphics processing unit (GPU)-enhanced clusters. There are a wide range of custom software and tools, including traditional RDBMSs and display tools. Real-time data are obtained at Full Motion Video (FMV)—30 to 60 frames per second at full-color 1 080p resolution (i.e., 1 920 × 1 080 pixels, a high-definition progressive scan) or Wide-Area Large Format Imagery (WALF)—1 to 10 frames per second at 10,000 pixels × 10,000 pixels and full-color resolution. Visualization of extracted outputs will typically be as overlays on a geospatial (i.e., GIS) display. Analytics are basic object detection analytics and integration with sophisticated situation awareness tools with data fusion. Significant security issues must be considered; sources and methods cannot be compromised (i.e., “the enemy” should not know what we see).

Future:

A typical problem is integration of this processing into a large GPU cluster capable of processing data from several sensors in parallel and in near real time. Transmission of data from sensor to system is also a major challenge.

5.4.3 Use case 15: Intelligence Data Processing and Analysis

Application:

Intelligence analysts need the following capabilities:

- Identify relationships between entities (e.g., people, organizations, places, equipment).
- Spot trends in sentiment or intent for either the general population or a leadership group such as state and non-state actors.
- Identify the locations and possibly timing of hostile actions including implantation of improvised explosive devices.
- Track the location and actions of potentially hostile actors.
- Reason against and derive knowledge from diverse, disconnected, and frequently unstructured (e.g., text) data sources.
- Process data close to the point of collection, and allow for easy sharing of data to/from individual soldiers, forward-deployed units, and senior leadership in garrisons.

Current approach:

Data size ranges from tens of terabytes to hundreds of petabytes, with imagery intelligence devices gathering a petabyte in a few hours. Dismounted warfighters typically have at most one to hundreds of gigabytes (GBs), which is typically handheld data storage. Software information is described in Annex [A.3.3](#).

Future:

Data currently exist in disparate silos. These data must be accessible through a semantically integrated data space. A wide variety of data types, sources, structures, and quality will span domains and require integrated search and reasoning. Most critical data are either unstructured or maintained as imagery or video, which requires significant processing to extract entities and information. Network quality, provenance, and security are essential.

5.5 Health care and life sciences

5.5.1 Use case 16: Electronic Medical Record Data

Application:

Large national initiatives around health data are emerging. These include developing a digital learning health care system to support increasingly evidence-based clinical decisions with timely, accurate, and up-to-date patient-centered clinical information; using electronic observational clinical data to efficiently and rapidly translate scientific discoveries into effective clinical treatments; and electronically sharing integrated health data to improve healthcare process efficiency and outcomes. These key initiatives all rely on high-quality, large-scale, standardized, and aggregate health data. Advanced methods are needed for normalizing patient, provider, facility, and clinical concept identification within and among separate health care organizations. With these methods in place, feature selection, information retrieval, and enhanced machine learning decision-models can be used to define and extract clinical phenotypes from non-standard, discrete, and free-text clinical data. Clinical phenotype data must be leveraged to support cohort selection, clinical outcomes research, and clinical decision support.

Current approach:

The Indiana Network for Patient Care (INPC), the nation's largest and longest-running health information exchange, houses clinical data from more than 1,100 discrete logical operational healthcare sources. More than 20 TB of raw data, these data describe over 12 million patients and over 4 billion

discrete clinical observations. Between 500,000 and 1,5 million new real-time clinical transactions are added every day.

Future:

Running on an Indiana University supercomputer, Teradata, PostgreSQL, and MongoDB will support information retrieval methods to identify relevant clinical features (e.g., term frequency-inverse document frequency [tf-idf], latent semantic analysis, and mutual information). NLP techniques will extract relevant clinical features. Validated features will be used to parameterize clinical phenotype decision models based on maximum likelihood estimators and Bayesian networks. Decision models will be used to identify a variety of clinical phenotypes such as diabetes, congestive heart failure, and pancreatic cancer.

5.5.2 Use case 17: Pathology Imaging/Digital Pathology

Application:

Digital pathology imaging is an emerging field in which examination of high-resolution images of tissue specimens enables novel and more effective ways to diagnose diseases. Pathology image analysis segments massive spatial objects (e.g., millions of objects per image) such as nuclei and blood vessels, represented with their boundaries, along with many extracted image features from these objects. The derived information is used for many complex queries and analytics to support biomedical research and clinical diagnosis

Current approach:

Each 2D image comprises 1 GB of raw image data and entails 1,5 GB of analytical results. Message Passing Interface (MPI) is used for image analysis. Solution information is described in [A.4.2](#).

Future:

Recently, 3D pathology imaging has been made possible using 3D laser technologies or serially sectioning hundreds of tissue sections onto slides and scanning them into digital images. Segmenting 3D microanatomic objects from registered serial images could produce tens of millions of 3D objects from a single image. This provides a deep 'map' of human tissues for next-generation diagnosis. 3D images can comprise 1 TB of raw image data and entail 1 TB of analytical results. A moderated hospital would generate 1 PB of data per year.

5.5.3 Use case 18: Computational Bioimaging

Application:

Data delivered from bioimaging are increasingly automated, higher resolution, and multi-modal. This has created a data analysis bottleneck that, if resolved, can advance bioscience discovery through big data techniques.

Current approach:

The current piecemeal analysis approach does not scale to situations in which a single scan on emerging machines is 32 TB and medical diagnostic imaging is annually around 70 PB, excluding cardiology. A web-based, one-stop shop is needed for high-performance, high-throughput image processing for producers and consumers of models built on bio-imaging data.

Future:

The goal is to resolve that bottleneck with extreme-scale computing and community-focused science gateways, both of which apply massive data analysis toward massive imaging data sets. Workflow components include data acquisition, storage, enhancement, noise minimization, segmentation of regions of interest, crowd-based selection and extraction of features, and object classification, as well as organization and search. Possible software packages are described [A.4.3](#).

5.5.4 Use case 19: Genomic Measurements

Application:

The NIST Genome in a Bottle Consortium integrates data from multiple sequencing technologies and methods to develop highly confident characterization of whole human genomes as reference materials. The consortium also develops methods to use these reference materials to assess performance of any genome sequencing run.

Current approach:

NIST's approximately 40 TB network file system (NFS) is full. The National Institutes of Health (NIH) and the National Center for Biotechnology Information (NCBI) are also currently storing PBs of data. NIST is also storing data using open-source sequencing bioinformatics software from academic groups (UNIX-based) on a 72-core cluster, supplemented by larger systems at collaborators.

Future:

DNA sequencers can generate ~300 GB of compressed data per day, and this volume has increased much faster than Moore's Law gives for increase in computer processing power. Future data could include other 'omics' (e.g., genomics) measurements, which will be even larger than DNA sequencing. Clouds have been explored as a cost effective scalable approach.

5.5.5 Use case 20: Comparative Analysis for Metagenomes and Genomes

Application:

Given a metagenomic sample this use case aims to do the following:

- determine the community composition in terms of other reference isolate genomes;
- characterize the function of its genes;
- begin to infer possible functional pathways;
- characterize similarity or dissimilarity with other metagenomic samples;
- begin to characterize changes in community composition and function due to changes in environmental pressures;
- isolate subsections of data based on quality measures and community composition.

Current approach:

The current integrated comparative analysis system for metagenomes and genomes is front-ended by an interactive web user interface (UI) with core data. The system involves backend precomputations and batch job computation submission from the UI. The system provides an interface to standard bioinformatics tools (e.g., BLAST, HMMER, multiple alignment and phylogenetic tools, gene callers, sequence feature predictors).

Future:

Management of heterogeneity of biological data is currently performed by a RDBMS (i.e., Oracle). Unfortunately, it does not scale for even the current volume, 50 TB of data. NoSQL solutions aim at providing an alternative, but unfortunately they do not always lend themselves to real-time interactive use or rapid and parallel bulk loading, and sometimes they have issues regarding robustness.

5.5.6 Use case 21: Individualized Diabetes Management

Application:

Diabetes is a growing illness in the world population, affecting both developing and developed countries. Current management strategies do not adequately take into account individual patient profiles, such as co-morbidities and medications, which are common in patients with chronic illnesses. Advanced graph-based data mining techniques must be applied to electronic health records (EHRs), converting them into RDF (Resource Description Framework) graphs. These advanced techniques would facilitate searches for diabetes patients and allow for extraction of their EHR data for outcome evaluation.

Current approach:

Typical patient data records are composed of 100 controlled vocabulary values and 1,000 continuous values. Most values have a timestamp. The traditional paradigm of relational row-column lookup needs to be updated to semantic graph traversal.

Future:

The first step is to compare patient records to identify similar patients from a large EHR database (i.e., an individualized cohort.) Each patient's management outcome should be evaluated to formulate the most appropriate solution for a given patient with diabetes. The time-dependent properties need to be processed before query to allow matching based on derivatives and other derived properties. Software information is described in [A.4.6](#).

5.5.7 Use case 22: Statistical Relational Artificial Intelligence for Health Care

Application:

The goal of the project is to analyze large, multi-modal medical data, including different data types such as imaging, EHR, and genetic and natural language. This approach employs relational probabilistic models that have the capability of handling rich relational data and modeling uncertainty using probability theory. The software learns models from multiple data types, and can possibly integrate information and reason about complex queries. Users can provide a set of descriptions, for instance: magnetic resonance imaging (MRI) images and demographic data about a particular subject. They can then query for the onset of a particular disease (e.g., Alzheimer's), and the system will provide a probability distribution over the possible occurrence of this disease.

Current approach:

A single server can handle a test cohort of a few hundred patients with associated data of hundreds of GBs.

Future:

A cohort of millions of patients can involve PB size datasets. A major issue is the availability of too much data (e.g., images, genetic sequences), which can make the analysis complicated. Sometimes, large amounts of data about a single subject are available, but the number of subjects is not very high (i.e., data imbalance). This can result in learning algorithms picking up random correlations between the multiple data types as important features in analysis. Another challenge lies in aligning the data and merging from multiple sources in a form that will be useful for a combined analysis.

5.5.8 Use case 23: World Population-Scale Epidemiological Study

Application:

There is a need for reliable, real-time prediction and control of pandemics similar to the 2009 H1N1 influenza. Addressing various kinds of contagion diffusion may involve modeling and computing information, diseases, and social unrest. Agent-based models can utilize the underlying interaction network (i.e., a network defined by a model of people, vehicles, and their activities) to study the evolution of the desired phenomena.

Current approach:

There is a two-step approach: (1) build a synthetic global population; and (2) run simulations over the global population to reason about outbreaks and various intervention strategies. The current 100 TB

dataset was generated centrally with an MPI-based simulation system written in Charm++. Parallelism is achieved by exploiting the disease residence time period.

Future:

Large social contagion models can be used to study complex global-scale issues, greatly increasing the size of systems used.

5.5.9 Use case 24: Social Contagion Modeling for Planning, Public Health, and Disaster Management

Application:

Social behavior models are applicable to national security, public health, viral marketing, city planning, and disaster preparedness. In a social unrest application, people take to the streets to voice either unhappiness with or support for government leadership. Models would help quantify the degree to which normal business and activities are disrupted because of fear and anger, the possibility of peaceful demonstrations and/or violent protests, and the potential for government responses ranging from appeasement, to allowing protests, to issuing threats against protestors, to taking actions to thwart protests. Addressing these issues would require fine-resolution models (at the level of individual people, vehicles, and buildings) and datasets.

Current approach:

The social contagion model infrastructure simulates different types of human-to-human interactions (e.g., face-to-face versus online media), and also interactions between people, services (e.g., transportation), and infrastructure (e.g. Internet, electric power). These activity models are generated from averages such as census data.

Future:

One significant concern is data fusion (i.e., how to combine data from different sources and how to deal with missing or incomplete data.). A valid modeling process must take into account heterogeneous features of hundreds of millions or billions of individuals, as well as cultural variations across countries. For such large and complex models, the validation process itself is also a challenge.

5.5.10 Use case 25: Biodiversity and LifeWatch

Application:

Research and monitor different ecosystems, biological species, their dynamics, and their migration with a mix of custom sensors and data access/processing, and a federation with relevant projects in the area. Particular case studies include monitoring alien species, migrating birds, and wetlands. One of many efforts from the consortium titled Common Operations for Environmental Research Infrastructures (ENVRI) is investigating integration of LifeWatch with other environmental e-infrastructures.

Current approach:

At this time, this project is in the preliminary planning phases and, therefore, the current approach is not fully developed.

Future:

The LifeWatch initiative will provide integrated access to a variety of data, analytical, and modeling tools as served by a variety of collaborating initiatives. It will also offer data and tools in selected workflows for specific scientific communities. In addition, LifeWatch will provide opportunities to construct personalized “virtual labs,” allowing participants to enter and access new data and analytical tools. New data will be shared with the data facilities cooperating with LifeWatch, including both the Global Biodiversity Information Facility and the Biodiversity Catalogue, also known as the Biodiversity Science Web Services Registry. Data include ‘omics’, species information, ecological information (e.g.,

biomass, population density), and ecosystem data (e.g., carbon dioxide [CO₂] fluxes, algal blooming, water and soil characteristics.)

5.6 Deep Learning and Social Media

5.6.1 Use case 26: Large-Scale Deep Learning

Application:

There is a need to increase the size of datasets and models that can be tackled with deep learning algorithms. Large models (e.g., neural networks with more neurons and connections) combined with large datasets are increasingly the top performers in benchmark tasks for vision, speech, and NLP. It will be necessary to train a deep neural network from a large (e.g., much greater than 1 TB) corpus of data, which is typically comprised of imagery, video, audio, or text. Such training procedures often require customization of the neural network architecture, learning criteria, and dataset pre-processing. In addition to the computational expense demanded by the learning algorithms, the need for rapid prototyping and ease of development is extremely high.

Current approach:

The largest applications so far are to image recognition and scientific studies of unsupervised learning with 10 million images and up to 11 billion parameters on a 64 GPU HPC Infiniband cluster. Both supervised (i.e., using existing classified images) and unsupervised applications are being investigated.

Future:

Large datasets of 100 TB or more may be necessary to exploit the representational power of the larger models. Training a self-driving car could take 100 million images at megapixel resolution. Deep learning shares many characteristics with the broader field of machine learning. The paramount requirements are high computational throughput for mostly dense linear algebra operations, and extremely high productivity for researcher exploration. High-performance libraries must be integrated with high-level prototyping environments.

5.6.2 Use case 27: Organizing Large-Scale, Unstructured Collections of Consumer Photos

Application:

Collections of millions to billions of consumer images are used to produce 3D reconstructions of scenes—with no a priori knowledge of either the scene structure or the camera positions. The resulting 3D models allow efficient and effective browsing of large-scale photo collections by geographic position. New images can be geolocated by matching them to 3D models, and object recognition can be performed on each image. The 3D reconstruction can be posed as a robust, non-linear, least squares optimization problem: observed or noisy correspondences between images are constraints, and unknowns are six-dimensional (6D) camera poses of each image and 3D positions of each point in the scene.

Current approach:

The current system information is described in [A.5.2](#). Over 505 billion images are currently on SNS, with over 500 million images added to social media sites each day.

Future:

Necessary maintenance and upgrades require many analytics including feature extraction, feature matching, and large-scale probabilistic inference. These analytics appear in many or most computer vision and image processing problems, including recognition, stereo resolution, and image denoising. Other needs are visualizing large-scale, 3D reconstructions and navigating large-scale collections of images that have been aligned to maps.

5.6.3 Use case 28: Truthy—Information Diffusion Research from Twitter Data

Application:

How communication spreads on socio-technical networks must be better understood, and methods are needed to detect potentially harmful information spread at early stages (e.g., deceiving messages, orchestrated campaigns, and untrustworthy information).

Current approach:

Twitter generates a large volume of continuous streaming data—about 30 TB a year, compressed—through circulation of ~100 million messages per day. The increase over time is roughly 500 GB data per day. All these data must be acquired and stored. Additional needs include near real-time analysis of such data for anomaly detection, stream clustering, signal classification, and online-learning; and data retrieval, big data visualization, data-interactive web interfaces, and public application programming interfaces (APIs) for data querying. Software information is described in [A.5.4](#). Information diffusion, clustering, and dynamic network visualization capabilities already exist.

Future:

It plans to expand and so needs to move toward advanced distributed storage programs, and in-memory database, described in [A.5.4](#) for real-time analysis. Solutions will need to incorporate streaming clustering, anomaly detection, and online learning.

5.6.4 Use case 29: Crowd Sourcing in the Humanities as Source for Big and Dynamic Data

Application:

Information is captured from many individuals and their devices using a range of sources: manually entered, recorded multimedia, reaction times, pictures, sensor information. These data are used to characterize wide-ranging individual, social, cultural, and linguistic variations among several dimensions (e.g., space, social space, time).

Current approach:

At this point, typical systems used are Extensible Markup Language (XML) technology and traditional relational databases. Other than pictures, not much multi-media is employed yet.

Future:

Crowd sourcing is beginning to be used on a larger scale. However, the availability of sensors in mobile devices provides a huge potential for collecting large amount of data from numerous individuals. This possibility has not been explored on a large scale so far; existing crowd sourcing projects are usually of a limited scale and web-based. Privacy issues may be involved because of access to individuals' audio-visual files; anonymization may be necessary but not always possible. Data management and curation are critical. With multimedia, the size could be hundreds of terabytes.

5.6.5 Use case 30: CINET—Cyberinfrastructure for Network (Graph) Science and Analytics

Application:

CINET provides a common web-based platform that allows the end user seamless access to the following:

- network and graph analysis tools such as SNAP, NetworkX, and Galib;
- real-world and synthetic networks;
- computing resources;
- data management systems.

Current approach:

CINET uses an InfiniBand-connected HPC cluster with 720 cores to provide HPC as a service. The platform is being used for research and education. CINET is used in classes and to support research by social science and social networking communities

Future:

Rapid repository growth is expected to lead to at least 1,000 to 5,000 networks and methods in about a year. As more fields use graphs of increasing size, parallel algorithms will be important. Two critical challenges are data manipulation and bookkeeping of the derived data, as there are no well-defined and effective models and tools for unified management of various graph data.

5.6.6 Use case 31: NIST Information Access Division — Analytic Technology Performance Measurements, Evaluations, and Standards

Application:

Performance metrics, measurement methods, and community evaluations are needed to ground and accelerate development of advanced analytic technologies in the areas of speech and language processing, video and multimedia processing, biometric image processing, and heterogeneous data processing, as well as the interaction of analytics with users. Typically one of two processing models are employed: (1) push test data out to test participants, and analyze the output of participant systems, and (2) push algorithm test harness interfaces out to participants, bring in their algorithms, and test them on internal computing clusters.

Current approach:

There is a large annotated corpora of unstructured/semi-structured text, audio, video, images, multimedia, and heterogeneous collections of the above, including ground truth annotations for training, developmental testing, and summative evaluations. The test corpora exceed 900 million web pages occupying 30 TB of storage, 100 million tweets, 100 million ground-truthed biometric images, several hundred thousand partially ground-truthed video clips, and terabytes of smaller fully ground-truthed test collections.

Future:

Even larger data collections are being planned for future evaluations of analytics involving multiple data streams and very heterogeneous data. In addition to larger datasets, the future includes testing of streaming algorithms with multiple heterogeneous data. The use of clouds is being explored.

5.7 The Ecosystem for research

5.7.1 Use case 32: DataNet Federation Consortium

Application:

The DataNet Federation Consortium (DFC) promotes collaborative and interdisciplinary research through a federation of data management systems across federal repositories, national academic research initiatives, institutional repositories, and international collaborations. The collaboration environment runs at scale and includes petabytes of data, hundreds of millions of files, hundreds of millions of metadata attributes, tens of thousands of users, and a thousand storage resources.

Current approach:

Currently, 25 science and engineering domains have projects that rely on the iRODS (Integrated Rule-Oriented Data System) policy-based data management system. Active organizations include the National Science Foundation, with major projects such as the Ocean Observatories Initiative (sensor archiving); Temporal Dynamics of Learning Center (cognitive science data grid); iPlant Collaborative (plant genomics); Drexel's engineering digital library; and H. W. Odum Institute for Research in Social

Science (data grid federation with Dataverse). iRODS currently manages PB of data, hundreds of millions of files, hundreds of millions of metadata attributes, tens of thousands of users, and a thousand storage resources. It interoperates with workflow systems (e.g., National Center for Computing Applications' [NCSA's] Cyberintegrator, Kepler, Taverna), cloud, and more traditional storage models, as well as different transport protocols.

Future:

Future data scenarios and applications were not expressed for this use case.

5.7.2 Use case 33: The Discinnet Process

Application:

Discinnet has developed a Web 2.0 collaborative platform and research prototype as a pilot installation, which is now being deployed and tested by researchers from a growing number of diverse research fields. The goal is to reach a wide enough sample of active research fields, represented as clusters (i.e., researchers projected and aggregating within a manifold of mostly shared experimental dimensions) to test general, hence potentially interdisciplinary, epistemological models throughout the present decade.

Current approach:

Currently, 35 clusters have been started, with close to 100 awaiting more resources. There is potential for many more to be created, administered, and animated by research communities. Examples of clusters include optics, cosmology, materials, microalgae, health care, applied math, computation, rubber, and other chemical products/issues.

Future:

Discinnet itself would not be big data but rather will generate metadata when applied to a cluster that involves big data. In interdisciplinary integration of several fields, the process would reconcile metadata from many complexity levels.

5.7.3 Use case 34: Semantic Graph Search on Scientific Chemical and Text-Based Data

Application:

Social media-based infrastructure, terminology and semantic data-graphs are established to annotate and present technology information. The process uses root- and rule-based methods currently associated primarily with certain Indo-European languages, such as Sanskrit and Latin.

Current approach:

Many reports, including a recent one on the Material Genome Project, find that exclusive top-down solutions to facilitate data sharing and integration are not desirable for multi-disciplinary efforts. However, a bottom-up approach can be chaotic. For this reason, there is need for a balanced blend of the two approaches to support easy-to-use techniques to metadata creation, integration, and sharing. This challenge is very similar to the challenge faced by language developers, so a recently developed method is based on these ideas. There are ongoing efforts to extend this method to publications of interest to the Material Genome Initiative, the Open Government movement, and the NIST Integrated Knowledge Editorial Net (NIKE) a NIST-wide publication archive. These efforts are a component of the Research Data Alliance Metadata Standards Directory Working Group.

Future:

A cloud infrastructure should be created for social media of scientific information. Scientists from across the world could use this infrastructure to participate and deposit results of their experiments. Prior to establishing a scientific social medium, some issues must be resolved including the following:

- Minimize challenges related to establishing re-usable, interdisciplinary, scalable, on-demand, use-case, and user-friendly vocabulary.

- Adopt an existing or create new on-demand 'data-graph' to place information in an intuitive way, such that it would easily integrate with existing data-graphs in a federated environment, independently of details of data management.
- Find relevant scientific data without spending too much time on the Internet.

Start with resources such as the Open Government movement, Material Genome Initiative, and Protein Databank. This effort includes many local and networked resources. Developing an infrastructure to automatically integrate information from all these resources using data-graphs is a challenge, but steps are being taken to solve it. Strong database tools and servers for data-graph manipulation are needed.

5.7.4 Use case 35: Light Source Beamlines

Application:

Samples are exposed to X-rays from light sources in a variety of configurations, depending on the experiment. Detectors, essentially high-speed digital cameras, collect the data. The data are then analyzed to reconstruct a view of the sample or process being studied.

Current approach:

A variety of commercial and open source software is used for data analysis. Data transfer is accomplished using physical transport of portable media, which severely limits performance, high-performance GridFTP, managed by Globus Online, or workflow systems such as SPADE (Support for Provenance Auditing in Distributed Environments, an open source software infrastructure).

Future:

Camera resolution is continually increasing. Data transfer to large-scale computing facilities is becoming necessary because of the computational power required to conduct the analysis on timescales useful to the experiment. Because of the large number of beamlines (e.g., 39 at the LBNL Advanced Light Source), aggregate data load is likely to increase significantly over the coming years, as will the need for a generalized infrastructure for analyzing GB per second of data from many beamline detectors at multiple facilities.

5.8 Astronomy and physics

5.8.1 Use case 36: Catalina Real-Time Transient Survey: A Digital, Panoramic, Synoptic Sky Survey

Application:

Catalina Real-Time Transient Survey (CRTS) explores the variable universe in the visible light regime, on timescales ranging from minutes to years, by searching for variable and transient sources. It discovers a broad variety of astrophysical objects and phenomena, including various types of cosmic explosions (e.g., supernovae), variable stars, phenomena associated with accretion to massive black holes (e.g., active galactic nuclei) and their relativistic jets, and high proper motion stars. The data are collected from three telescopes (two in Arizona and one in Australia), with additional ones expected in the near future in Chile.

Current approach:

The survey generates up to approximately 0,1 TB on a clear night with a total of approximately 100 TB in current data holdings. The data are pre-processed at the telescope and then transferred to the University of Arizona and Caltech for further analysis, distribution, and archiving. The data are processed in real time, and detected transient events are published electronically through a variety of dissemination mechanisms, with no proprietary withholding period (CRTS has a completely open data policy). Further data analysis includes classification of the detected transient events, additional observations using other telescopes, scientific interpretation, and publishing. This process makes

heavy use of the archival data (several PBs) from a wide variety of geographically distributed resources connected through the virtual observatory (VO) framework.

Future:

CRTS is a scientific and methodological test bed and precursor of larger surveys to come, notably the Large Synoptic Survey Telescope (LSST), expected to operate in the 2020s and selected as the highest-priority ground-based instrument in the 2010 Astronomy and Astrophysics Decadal Survey. LSST will gather about 30 TB per night.

Survey pipelines from telescopes (on the ground or in space) produce transient event data streams, and the events, along with their observational descriptions, are ingested by one or more depositories, from which the event data can be disseminated electronically to human astronomers or robotic telescopes. Each event is assigned an evolving portfolio of information, which includes all available data on that celestial position. The data are gathered from a wide variety of data archives unified under the Virtual Observatory framework, expert annotations, etc.

Representations of such federated information can be both human-readable and machine-readable. The data are fed into one or more automated event characterization, classification, and prioritization engines that deploy a variety of machine learning tools for these tasks.

The engines' output, which evolves dynamically as new information arrives and is processed, informs the follow-up observations of the selected events, and the resulting data are communicated back to the event portfolios for the next iteration.

Users, either human or robotic, can tap into the system at multiple points, both for information retrieval and to contribute new information, through a standardized set of formats and protocols. This could be done in (near) real-time or in archival (i.e., not time-critical) modes.

5.8.2 Use case 37: DOE Extreme Data from Cosmological Sky Survey and Simulations

Application:

A cosmology discovery tool integrates simulations and observation to clarify the nature of dark matter, dark energy, and inflation—some of the most exciting, perplexing, and challenging questions facing modern physics, including the properties of fundamental particles affecting the early universe. The simulations will generate data sizes comparable to observation.

Current approach:

At this time, this project is in the preliminary planning phases and, therefore, the current approach is not fully developed.

Future:

These systems will use huge amounts of supercomputer time — over 200 million hours. Associated data sizes are as follows:

- Dark Energy Survey (DES): 4 PB per year in 2015;
- Zwicky Transient Factory (ZTF): 1 PB per year in 2015;
- LSST (see CRTS discussion above): 7 PB per year in 2019;
- Simulations: 10 PB per year in 2017.

5.8.3 Use case 38: Large Survey Data for Cosmology

Application:

For DES, the data are sent from the mountaintop, via a microwave link, to La Serena, Chile. From there, an optical link forwards them to the NCSA and to NERSC for storage and 'reduction.' Here, galaxies

and stars in both the individual and stacked images are identified and catalogued, and finally their properties are measured and stored in a database.

Current approach:

Subtraction pipelines are run using extant imaging data to find new optical transients through machine learning algorithms. Data technologies and hardware resources are described in [A.7.3](#).

Future:

Techniques are needed for handling Cholesky decomposition for thousands of simulations with matrices of order one million on a side and parallel image storage. LSST will generate 60 PB of imaging data and 15 PB of catalogue data and a correspondingly large (or larger) amount of simulation data. In total, over 20 TB of data will be generated per night.

5.8.4 Use case 39: Particle Physics—Analysis of Large Hadron Collider Data: Discovery of Higgs Particle

Application:

Analysis is conducted on collisions at the European Organization for Nuclear Research (CERN) Large Hadron Collider (LHC) accelerator.

Processed information defines physics properties of events and generates lists of particles with type and momenta. These events are analyzed to find new effects—both new particles (e.g., Higgs), and present evidence that conjectured particles (e.g., Supersymmetry) have not been detected. A few major experiments are being conducted at LHC, including ATLAS and CMS (Compact Muon Solenoid). These experiments have global participants (e.g., CMS has 3,600 participants from 183 institutions in 38 countries), and so the data at all levels are transported and accessed across continents.

Current approach:

The LHC experiments are pioneers of a distributed big data science infrastructure. Several aspects of the LHC experiments' workflow highlight issues that other disciplines will need to solve. These issues include automation of data distribution, high-performance data transfer, and large-scale high-throughput computing. A data grid analysis for Higgs Particle discovery utilised 350,000 cores running near-continuously—over two million jobs per day arranged in three major tiers: CERN, Continents/Countries, and Universities. The analysis uses distributed, high-throughput computing (i.e., pleasing parallel) architecture with facilities integrated across the world by the Worldwide LHC Computing Grid (WLCG) and Open Science Grid in the U.S. Accelerator data and analysis generates 15 PB of data each year for a total of 200 PB. Specifically, in 2012, ATLAS had 8 PB on Tier1 tape and over 10 PB on Tier 1 disk at BNL and 12 PB on disk cache at U.S. Tier 2 centers. CMS has similar data sizes. Over half the resources are used for Monte Carlo simulations as opposed to data analysis.

Future:

In the past, the particle physics community has been able to rely on industry to deliver exponential increases in performance per unit cost over time, as described by Moore's Law. However, the available performance will be much more difficult to exploit in the future since technology limitations, in particular regarding power consumption, have led to profound changes in the architecture of modern central processing unit (CPU) chips.

In the past, software could run unchanged on successive processor generations and achieve performance gains that follow Moore's Law, thanks to the regular increase in clock rate that continued until 2006. The era of scaling sequential applications on an HEP (heterogeneous element processor) is now over. Changes in CPU architectures imply significantly more software parallelism, as well as exploitation of specialized floating point capabilities.

The structure and performance of HEP data processing software need to be changed such that they can continue to be adapted and developed to run efficiently on new hardware. This represents a major paradigm shift in HEP software design and implies large-scale re-engineering of data structures and

algorithms. Parallelism needs to be added simultaneously at all levels: the event level, the algorithm level, and the sub-algorithm level. Components at all levels in the software stack need to interoperate, and therefore the goal is to standardize as much as possible on basic design patterns and on the choice of a concurrency model. This will also help to ensure efficient and balanced use of resources.

5.8.5 Use case 40: Belle II High Energy Physics Experiment

Application:

The Belle experiment is a particle physics experiment with more than 400 physicists and engineers investigating charge parity (CP) violation effects with B meson production at the High Energy Accelerator KEKB e⁺ e⁻ accelerator in Tsukuba, Japan. In particular, numerous decay modes at the Upsilon (4S) resonance are sought to identify new phenomena beyond the standard model of particle physics. This accelerator has the largest intensity of any in the world, but the events are simpler than those from LHC, and so analysis is less complicated, but similar in style to the CERN accelerator analysis.

Current approach:

At this time, this project is in the preliminary planning phases and, therefore, the current approach is not fully developed.

Future:

An upgraded experiment Belle II and accelerator SuperKEKB was starting operation in 2015. Data will increase by a factor of 50, with total integrated raw data of ~120 PB and physics data of ~15 PB and ~100 PB of Monte Carlo samples. The next stage will necessitate a move to a distributed computing model requiring continuous raw data transfer of ~20 GB per second at designed luminosity between Japan and the United States. Required softwares are described in [A.7.5](#).

5.9 Earth, environmental, and polar science

5.9.1 Use case 41: European Incoherent Scatter Scientific Association 3D Incoherent Scatter Radar System

Application:

EISCAT conducts research on the lower, middle, and upper atmosphere and ionosphere using the incoherent scatter radar technique. This technique is the most powerful ground-based tool for these research applications. EISCAT studies instabilities in the ionosphere and investigates the structure and dynamics of the middle atmosphere. EISCAT operates a diagnostic instrument in ionospheric modification experiments with addition of a separate heating facility. Currently, EISCAT operates three of the ten major incoherent radar scattering instruments worldwide; their three systems are located in the Scandinavian sector, north of the Arctic Circle.

Current approach:

The currently running EISCAT radar generates data at rates of terabytes per year. The system does not present special challenges.

Future:

The design of the next-generation radar, EISCAT_3D, will consist of a core site with transmitting and receiving radar arrays and four sites with receiving antenna arrays at some 100 km from the core. The fully operational five-site system will generate several thousand times the number of data of the current EISCAT system, with 40 PB per year in 2022, and is expected to operate for 30 years. EISCAT_3D data e-Infrastructure plans to use high-performance computers for central site data processing and high-throughput computers for mirror site data processing. Downloading the full data is not time-critical, but operations require real-time information about certain pre-defined events, which would be sent from the sites to the operations center, and a real-time link from the operations center to the sites to set the mode of radar operation in real time.

5.9.2 Use case 42: Common Operations of Environmental Research Infrastructure

Application:

ENVRI (Common Operations of Environmental Research Infrastructures) addresses European distributed, long-term, remote-controlled observational networks focused on understanding processes, trends, thresholds, interactions, and feedbacks, as well as increasing the predictive power to address future environmental challenges. The following efforts are part of ENVRI.

- ICOS (Integrated Carbon Observation System) is a European distributed infrastructure dedicated to the monitoring of greenhouse gases (GHGs) through its atmospheric, ecosystem, and ocean networks.
- EURO-Argo is the European contribution to Argo, which is a global ocean observing system.
- EISCAT_3D (described separately) is a European new-generation incoherent scatter research radar system for upper atmospheric science.
- LifeWatch (described separately) is an e-science infrastructure for biodiversity and ecosystem research.
- EPOS (European Plate Observing System) is a European research infrastructure for earthquakes, volcanoes, surface dynamics, and tectonics.
- EMSO (European Multidisciplinary Seafloor and Water Column Observatory) is a European network of seafloor observatories for the long-term monitoring of environmental processes related to ecosystems, climate change, and geo-hazards.
- IAGOS (In-service Aircraft for a Global Observing System) is setting up a network of aircraft for global atmospheric observation.
- SIOS (Svalbard Integrated Arctic Earth Observing System) is establishing an observation system in and around Svalbard that integrates the studies of geophysical, chemical, and biological processes from all research and monitoring platforms.

Current approach:

ENVRI develops a reference model (ENVRI RM) as a common ontological framework and standard for the description and characterization of computational and storage infrastructures. The goal is to achieve seamless interoperability between the heterogeneous resources of different infrastructures. The ENVRI RM serves as a common language for community communication, providing a uniform framework into which the infrastructure's components can be classified and compared. The ENVRI RM also serves to identify common solutions to common problems. Data sizes in a given infrastructure vary from GBs to petabytes per year.

Future:

ENVRI's common environment will empower the users of the collaborating environmental research infrastructures and enable multidisciplinary scientists to access, study, and correlate data from multiple domains for system-level research. Collaboration affects big data requirements coming from interdisciplinary research.

ENVRI analyzed the computational characteristics of the six European Strategy Forum on Research Infrastructures (ESFRI) environmental research infrastructures, and identified five common subsystems. They are defined in the ENVRI RM (<http://www.envri.eu/rm>) and below.

- Data acquisition: Collects raw data from sensor arrays, various instruments, or human observers, and brings the measurements (data streams) into the system.
- Data curation: Facilitates quality control and preservation of scientific data and is typically operated at a data center.

- Data access: Enables discovery and retrieval of data housed in data resources managed by a data curation subsystem.
- Data processing: Aggregates data from various resources and provides computational capabilities and capacities for conducting data analysis and scientific experiments.
- Community support: Manages, controls, and tracks users' activities and supports users in conduct of their community roles.

5.9.3 Use case 43: Radar Data Analysis for the Center for Remote Sensing of Ice Sheets

Application:

The Center for Remote Sensing of Ice Sheets (CRE SIS) effort uses custom radar systems to measure ice sheet bed depths and (annual) snow layers at the North and South Poles and mountainous regions.

Resulting data feed into the Intergovernmental Panel on Climate Change (IPCC). The radar systems are typically flown in by aircraft in multiple paths.

Current approach:

The initial analysis uses Matlab signal processing that produces a set of radar images. These cannot be transported from the field over the Internet and are typically copied onsite to a few removable disks that hold a terabyte of data, then flown to a laboratory for detailed analysis. Image features (i.e., layers) are found using image understanding tools with some human oversight. A typical echogram with detected boundaries differentiates the boundary between air and ice layers and boundary between ice and terrain. This information is stored in a database front-ended by a geographical information system. The ice sheet bed depths are used in simulations of glacier flow. Each trip into the field, usually lasting a few weeks, results in 50 to 100 TB of data.

Future:

With improved instrumentation, an order of magnitude more data (a petabyte per mission) is projected. As the increasing field data must be processed in an environment with constrained power access, low-power or low-performance architectures, such as GPU systems, are indicated.

5.9.4 Use case 44: Unmanned Air Vehicle Synthetic Aperture Radar (UAVSAR) Data Processing, Data Product Delivery, and Data Services

Application:

Synthetic aperture radar (SAR) can identify landscape changes caused by seismic activity, landslides, deforestation, vegetation changes, and flooding. This function can be used to support earthquake science as well as disaster management. This use case supports the storage, image processing application, and visualization of geo-located data with angular specification.

Current approach:

Data from planes and satellites are processed on NASA computers before being stored after substantial data communication. The data are made public upon processing. They require significant curation owing to instrumental glitches. The current data size is approximately 150 TB.

Future:

The data size would increase dramatically if Earth Radar Mission launched. Clouds are suitable hosts but are not used today in production.

5.9.5 Use case 45: NASA Langley Research Center/ Goddard Space Flight Center iRODS Federation Test Bed

Application:

NASA Center for Climate Simulation and NASA Atmospheric Science Data Center have complementary data sets, each containing vast amounts of data that are not easily shared and queried. Climate researchers, weather forecasters, instrument teams, and other scientists need to access data from across multiple datasets in order to compare sensor measurements from various instruments, compare sensor measurements to model outputs, calibrate instruments, look for correlations across multiple parameters, and more.

Current approach:

Data are generated from two products: the Modern Era Retrospective Analysis for Research and Applications (MERRA, described separately) and NASA Clouds and Earth's Radiant Energy System (CERES) EBAF-TOA (Energy Balanced and Filled-Top of Atmosphere) product, which accounts for about 420 MB, and the EBAF-Surface product, which accounts for about 690 MB. Data numbers grow with each version update (about every six months). To analyze, visualize, and otherwise process data from heterogeneous datasets is currently a time-consuming effort. Scientists must separately access, search for, and download data from multiple servers, and often the data are duplicated without an understanding of the authoritative source. Often accessing data takes longer than scientific analysis. Current datasets are hosted on modest-sized (144 to 576 cores) InfiniBand clusters.

Future:

Improved access will be enabled through the use of iRODS. These systems support parallel downloads of datasets from selected replica servers, providing users with worldwide access to the geographically dispersed servers. iRODS operation will be enhanced with semantically organized metadata and managed via a highly precise NASA Earth Science ontology. Cloud solutions will also be explored.

5.9.6 Use case 46: MERRA Analytic Services (MERRA/AS)

Application:

This application produces global temporally and spatially consistent syntheses of 26 key climate variables by combining numerical simulations with observational data. Three-dimensional results are produced every six hours extending from 1979 to the present. The data support important applications such as IPCC research and the NASA/Department of Interior RECOVER wildfire decision support system; these applications typically involve integration of MERRA with other datasets.

Current approach:

Map/Reduce is used to process a current total of 480 TB. The current system is hosted on a 36-node InfiniBand cluster.

Future:

Clouds are being investigated. The data is growing by one TB a month.

5.9.7 Use case 47: Atmospheric Turbulence – Event Discovery and Predictive Analytics

Application:

Data mining is built on top of reanalysis products, including MERRA (described separately) and the North American Regional Reanalysis (NARR), a long-term, high-resolution climate data set for the North American domain. The analytics correlate aircraft reports of turbulence (either from pilot reports or from automated aircraft measurements of eddy dissipation rates) with recently completed atmospheric reanalyses. The information is of value to aviation industry and to weather forecasters. There are no standards for reanalysis products, complicating systems for which Map/Reduce is being

investigated. The reanalysis data are hundreds of terabytes, slowly updated, whereas the turbulence dataset is smaller in size and implemented as a streaming service.

Current approach:

The current 200 TB dataset can be analyzed with Map/Reduce or the like using SciDB or another scientific database.

Future:

The dataset will reach 500 TB in five years. The initial turbulence case can be extended to other ocean/atmosphere phenomena, but the analytics would be different in each case.

5.9.8 Use case 48: Climate Studies Using the Community Earth System Model at the U.S. Department of Energy (DOE) NERSC Center

Application:

Simulations with the Community Earth System Model (CESM) can be used to understand and quantify contributions of natural and anthropogenic-induced patterns of climate variability and change in the 20th and 21st centuries. The results of supercomputer simulations across the world should be stored and compared.

Current approach:

The Earth System Grid (ESG) enables global access to climate science data on a massive scale — petascale, or even exascale — with multiple petabytes of data at dozens of federated sites worldwide. The ESG is recognized as the leading infrastructure for the management and access of large distributed data volumes for climate change research. It supports the Coupled Model Intercomparison Project (CMIP), whose protocols enable the periodic assessments carried out by the IPCC.

Future:

Rapid growth of data is expected, with 30 PB produced at NERSC (assuming 15 end-to-end climate change experiments) in 2017 and many times more than this worldwide.

5.9.9 Use case 49: DOE Biological and Environmental Research (BER) Subsurface Biogeochemistry Scientific Focus Area

Application:

A genome-enabled watershed simulation capability (GEWaSC) is needed to provide a predictive framework for understanding the following:

- how genomic information stored in a subsurface microbiome affects biogeochemical watershed functioning;
- how watershed-scale processes affect microbial functioning;
- how these interactions co-evolve.

Current approach:

Current modeling capabilities can represent processes occurring over an impressive range of scales — from a single bacterial cell to that of a contaminant plume. Data cross all scales from genomics of the microbes in the soil to watershed hydro-biogeochemistry. Data are generated by the different research areas and include simulation data, field data (e.g., hydrological, geochemical, geophysical), 'omics' data, and observations from laboratory experiments.

Future:

Little effort to date has been devoted to develop a framework for systematically connecting scales, as is needed to identify key controls and to simulate important feedbacks. GEWaSC will develop a simulation framework that formally scales from genomes to watersheds and will synthesize diverse and disparate field, laboratory, and simulation datasets across different semantic, spatial, and temporal scales.

5.9.10 Use case 50: DOE BER AmeriFlux and FLUXNET Networks

Application:

AmeriFlux and Flux Tower Network (FLUXNET) are U.S. and world collections, respectively, of sensors that observe trace gas fluxes (e.g., CO₂, water vapor) across a broad spectrum of times (e.g., hours, days, seasons, years, and decades) and space. Moreover, such datasets provide the crucial linkages among organisms, ecosystems, and process-scale studies — at climate-relevant scales of landscapes, regions, and continents—for incorporation into biogeochemical and climate models.

Current approach:

Software information is described in [A.8.10](#). There are approximately 150 towers in AmeriFlux and over 500 towers distributed globally collecting flux measurements.

Future:

Field experiment data-taking would be improved by access to existing data and automated entry of new data via mobile devices. Interdisciplinary studies integrating diverse data sources will be expanded.

5.10 Energy

5.10.1 Use case 51: Consumption Forecasting in Smart Grids

Application:

Smart meters support prediction of energy consumption for customers, transformers, substations and the electrical grid service area. Advanced meters provide measurements every 15 min at the granularity of individual consumers within the service area of smart power utilities. Data to be combined include the head end of smart meters (distributed), utility databases (customer information, network topology; centralized), U.S. Census data (distributed), NOAA weather data (distributed), micro-grid building information systems (centralized), and micro-grid sensor networks (distributed). The central theme is real-time, data-driven analytics for time series from cyber physical systems.

Current approach:

Forecasting uses GIS-based visualization. Data amount to around 4 TB per year for a city such as Los Angeles with 1,4 million sensors. There are significant privacy issues requiring anonymization by aggregation. Real-time and historic data are combined with machine learning to predict consumption. Software information is described in [A.9.1](#).

Future:

Advanced grid technologies will be widely deployed. Smart grids will have new analytics integrating diverse data and supporting curtailment requests. New technologies will support mobile applications for client interactions.

5.10.2 Use case 52: Home Energy Management System

Application:

HEMS (Home Energy Management System) is useful system for energy conservation in private homes. In the HEMS, many kinds of sensors and devices are introduced into private homes, such as, smart meter, electric vehicle, solar power panel, light, air conditioner, fuel cell, water heater, storage battery. Energy manager gathers those data generated at private homes and stores them into cloud database named the

large HEMS information platform. Information manager operates the large HEMS information platform and manages data. Privacy and security of users are responsible to Information manager. Servicer analyzes data and provides valuable information to users as a service.

Current approach:

Services provided by servicer is not restricted by monitoring service of power usage. Other examples of useful services are elderly person life watching service, appropriate energy contract plan suggestion, prediction of PV power generation, coupon incentive-based demand response.

Future:

Standardization of Application Programming Interface will be necessary to increase usefulness of HEMS data

6 Use cases derived technical considerations

Technical considerations are the challenges limiting further use of big data. After collection, processing, and review of the use cases, technical considerations within seven characteristic categories were extracted from the individual use cases. These use case specific technical considerations were then aggregated to produce high-level, general technical considerations, within the seven characteristic categories, that are vendor neutral and technology agnostic. It is emphasized that neither the use case nor the requirements lists are exhaustive.

6.1 Use case specific technical considerations

Each use case was evaluated for technical considerations within the following seven categories:

- **Data source** (e.g., data size, file formats, rate of growth, at rest or in motion);
- **Data transformation** (e.g., data fusion, analytics);
- **Capabilities** (e.g., software tools, platform tools, hardware resources such as storage and networking);
- **Data consumer** (e.g., processed results in text, table, visual, and other formats);
- **Security and privacy;**
- **Life cycle management** (e.g., curation, conversion, quality check, pre-analytic processing);
- **Other technical considerations.**

Some use cases contained technical considerations in all seven categories while others only included technical considerations for a few categories. The complete list of specific technical considerations extracted from the use cases is presented in [Annex D](#). These categories informed the eventual selection of the roles specified in ISO/IEC 20547-3.

6.2 Summary of requirements analysis

There were 35 generic technical considerations [1] summarizing 439 specific technical considerations from the 52 use cases. Column 2 of [Table 1](#) gives the number of specific technical considerations driving this generic technical considerations.

Table 1 — Generic technical considerations with count of number of motivating specific technical considerations

#	Count	Generic technical considerations
Data source considerations		
1	28	Needs to support reliable real time, asynchronous, streaming, and batch processing to collect data from centralized, distributed, and cloud data sources, sensors, or instruments.
2	22	Needs to support slow, bursty, and high-throughput data transmission between data sources and computing clusters.
3	28	Needs to support diversified data content ranging from structured and unstructured text, document, graph, web, geospatial, compressed, timed, spatial, multimedia, simulation, and instrumental data.
Transformation considerations		
1	38	Needs to support diversified compute-intensive, analytic processing, and machine learning techniques.
2	7	Needs to support batch and real-time analytic processing.
3	15	Needs to support processing large diversified data content and modeling.
4	6	Needs to support processing data in motion (streaming, fetching new content, tracking, etc.).
Capability considerations		
1	29	Needs to support legacy and advanced software packages (software).
2	17	Needs to support legacy and advanced computing platforms (platform).
3	23	Needs to support legacy and advanced distributed computing clusters, co-processors, input output (I/O) processing (infrastructure).
4	14	Needs to support elastic data transmission (networking).
5	35	Needs to support legacy, large, and advanced distributed data storage (storage).
6	13	Needs to support legacy and advanced executable programming: applications, tools, utilities, and libraries (software).
Data consumer considerations		
1	4	Needs to support fast searches (~0,1 s) from processed data with high relevancy, accuracy, and high recall.
2	16	Needs to support diversified output file formats for visualization, rendering, and reporting.
3	2	Needs to support visual layout for results presentation.
4	11	Needs to support rich user interface for access using browser, visualization tools.
5	20	Needs to support high-resolution multi-dimension layer of data visualization.
6	1	Needs to support streaming results to clients.
Security and privacy considerations		
1	32	Needs to protect and preserve security and privacy on sensitive data.
2	12	Needs to support multi-level policy-driven, sandbox, access control, authentication on protected data.
Lifecycle management considerations		
1	20	Needs to support data quality curation including pre-processing, data clustering, classification, reduction, format transformation.
2	2	Needs to support dynamic updates on data, user profiles, and links.
3	6	Needs to support data lifecycle and long-term preservation policy, including data provenance.
4	4	Needs to support data validation.
5	4	Needs to support human annotation for data validation.
6	3	Needs to support prevention of data loss or corruption.
7	1	Needs to support multi-site archival.
8	2	Needs to support persistent identifier and data traceability.
9	1	Needs to support standardizing, aggregating, and normalizing data from disparate sources.
Other considerations		

Table 1 (continued)

#	Count	Generic technical considerations
1	6	Needs to support rich user interface from mobile platforms to access processed results
2	2	Needs to support performance monitoring on analytic processing from mobile platforms
3	13	Needs to support rich visual content search and rendering from mobile platforms
4	1	Needs to support mobile device data acquisition.
5	1	Needs to support security across mobile devices

6.3 Features of use cases

Table 2 lists the number of use cases tagged by various properties. This analysis from References [2][3][4] was the basis of the use case tags.

Table 2 — Features of use cases

Abbreviation	#	Description
PP	26	Pleasingly Parallel or Map Only
MR	18	Classic MapReduce MR (add MRStat below for full count)
MRStat	7	Simple version of MR where key computations are simple reduction as found in statistical averages, such as histograms and averages
MRIter	23	Iterative MapReduce or MPI
Graph	9	Complex graph data structure needed in analysis
Fusion	11	Integrate diverse data to aid discovery/decision making; could involve sophisticated algorithms or just be a portal
Streaming	41	Some data comes in incrementally and is processed this way
Classify	30	Classification: divide data into categories
S/Q	12	Index, Search and Query
CF	4	Collaborative Filtering for recommender engines
LML	36	Local Machine Learning (Independent for each parallel entity)
GML	23	Global Machine Learning: Deep Learning, Clustering, LDA, PLSI, MDS, Large Scale Optimizations as in Variational Bayes, MCMC, Lifted Belief Propagation, Stochastic Gradient Descent, L-BFGS, Levenberg-Marquardt. Can call EGO or Exascale Global Optimization with scalable parallel algorithm
	51	Workflow: Universal, so no label
GIS	16	Geotagged data often displayed in ESRI, Microsoft Virtual Earth, Google Earth, GeoServer, etc.
HPC	5	Classic large-scale simulation of cosmos, materials, etc., generating (visualization) data
Agent	2	Simulations of models of data-defined macroscopic entities represented as agents

Using this and an extended analysis, this table was expanded [3] to give 50 tags arranged in 4 views given in Tables 3 to 6.

Table 3 — Problem Architecture View Facets of Ogres (Meta or Macro Pattern)

Pleasingly Parallel	Seen in BLAST, Protein docking, some (bio-) imagery including Local Analytics or Local Machine Learning with pleasingly parallel filtering
Classic MapReduce	Search, Index and Query and Classification algorithms like collaborative filtering
Map Collective	Seen in machine learning – especially with linear algebra kernels
Map P2P	Point to Point Communication seen in parallel simulation and graph algorithms
Map Streaming	Combination of (parallel) long running maps accepting streamed data

Table 3 (continued)

Shared Memory	As opposed to distributed data (memory). Corresponds to problem where shared memory implementations are important. Tend to be dynamic asynchronous
SPMD	Single Program Multiple Data, well-known parallel computing style
BSP	Bulk Synchronous Processing: well-defined compute-communication phases
Fusion	Knowledge discovery often involves fusion of multiple methods or sources
Dataflow	Composite structure with multiple components linked by exchanged data
Agents	As used in epidemiology, discrete event simulations, etc. Swarm approaches
Workflow	Many applications often involve orchestration (workflow) of multiple components

Table 4 — Execution Features View Facets of Ogres

Performance metrics	As measured in benchmarks
Flops per byte	Important for performance
Execution Environment	Cloud or HPC; are Core libraries needed such as matrix-matrix/vector algebra, conjugate gradient, reduction, broadcast
Volume	Data size
Velocity	Measures Streaming
Variety	Multiple data sources are often mixed. See Fusion facet
Veracity	Accuracy of data affecting pre-processing needed and reliability of answer
Communication Structure	Interconnect structure? Is communication Synchronous or Asynchronous? In latter case shared memory may be attractive;
Static or Dynamic?	Does application (graph) change during execution?
Regularity	Most applications consist of a set of interconnected entities; is this regular as a set of pixels or is it a complicated irregular graph?
Iterative or not?	Important algorithm characteristic
Data Abstraction	Key-value, pixel, graph, vector, HDF5, Bag of words, etc.
Data Space?	Are data points in metric or non-metric spaces?
Complexity	Is algorithm $O(N^2)$ or $O(N)$ (up to logs) for N points per iteration?

Table 5 — Data Source and Style View Facets of Ogres

SQL, NoSQL or NewSQL	NoSQL includes Document, Column, Key-value, Graph, Triple store
Enterprise data systems	10 examples from NIST [5] integrate SQL/NoSQL
Files or Objects	Files as managed in iRODS and extremely common in scientific research. Objects most common in ABDS
HDFS/Lustre/GPFS	Are data and compute collocated?
Archive/Batched/Streaming	Streaming is Incremental update of datasets with new algorithms to achieve real-time response
Storage system styles	Styles include Shared, Dedicated, Permanent, and Transient
Metadata/Provenance	Define overall features of data and processing
Internet of Things	24 [5] to 50 (Cisco [7][8]) billion devices on the Internet by 2020
HPC generated data	Simulations generate visualization output that often needs to be mined
GIS	Geographical Information Systems provide access to geospatial data

Table 6 — Processing or Run-time View Facets of Ogres

Micro Benchmarks	A simple kernel or mini-app used to measure core system performance
LML	Local Analytics or Local Machine Learning
GML	Global Analytics or Machine Learning requiring iterative runtime

Table 6 (continued)

Base Statistics	Simple statistics seen in Table 2 as MRStat
Recommendations	Collaborative Filtering and other recommender analytics
Search/Query/Index	Rich set of technologies used in Search, Query and Indexing data
Classification	Technologies to label data (SVM, Bayes, deep learning, clustering)
Learning	Training algorithms
Optimization Methodology	Machine Learning, Nonlinear Optimization, Least Squares, Linear/Quadratic Programming, Combinatorial Optimization, expectation maximization, Monte Carlo, Variational Bayes, Global Inference
Streaming	Growing class of fast online O(N) algorithms
Alignment	Variant of Search seen in sequence comparison as in BLAST
Linear Algebra	Many machine learning algorithms build on linear algebra kernels
Graph	Problem set up as a graph as opposed to vector, grid, etc.
Visualization	Important component of many analysis pipelines

Copyrighted document, no reproduction or circulation allowed
 Click to view the full PDF of ISO/IEC TR 20547-2:2018
 For review by FG on AI in Healthcare
 Oct 2024

Annex A

Submitted use case studies

A.1 Government operation

A.1.1 Use case 1: Big data Archival: Census 2010 and 2000

Use case title	Big data Archival: Census 2010 and 2000 — Title 13 big data	
Vertical (area)	Digital Archives	
Author/company/email	Vivek Navale and Quyen Nguyen (NARA)	
Actors/stakeholders and their roles and responsibilities	NARA's Archivists Public users (after 75 years)	
Goals	Preserve data for a long term in order to provide access and perform analytics after 75 years. Title 13 of U.S. code authorizes the Census Bureau and guarantees that individual and industry specific data is protected.	
Use case description	Maintain data "as-is". No access and no data analytics for 75 years. Preserve the data at the bit-level. Perform curation, which includes format transformation if necessary. Provide access and analytics after nearly 75 years.	
Current solutions	Compute(System)	Linux servers
	Storage	NetApps, Magnetic tapes.
	Networking	
	Software	
Big data characteristics	Data source (distributed/centralized)	Centralized storage.
	Volume (size)	380 TB
	Velocity (e.g. real time)	Static.
	Variety (multiple datasets, mashup)	Scanned documents
	Variability (rate of change)	None
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Cannot tolerate data loss.
	Visualization	TBD
	Data quality (syntax)	Unknown.
	Data types	Scanned documents
	Data analytics	Only after 75 years.
Big data specific challenges (Gaps)	Preserve data for a long time scale.	
Big data specific challenges in mobility	TBD	

Security and privacy technical considerations	Title 13 data.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	
More information (URLs)	

A.1.2 Use case 2: NARA Accession, Search, Retrieve, Preservation

Use case title	National Archives and Records Administration Accession NARA Accession, Search, Retrieve, Preservation	
Vertical (area)	Digital Archives	
Author/company/email	Quyên Nguyen and Vivek Navale (NARA)	
Actors/stakeholders and their roles and responsibilities	Agencies' Records Managers NARA's Records Accessioners NARA's Archivists Public users×	
Goals	Accession, Search, Retrieval, and Long term Preservation of big data.	
Use case description	<ol style="list-style-type: none"> 1) Get physical and legal custody of the data. In the future, if data reside in the cloud, physical custody should avoid transferring big data from Cloud to Cloud or from Cloud to Data Center 2) Pre-process data for virus scan, identifying file format identification, removing empty files 3) Index 4) Categorize records (sensitive, unsensitive, privacy data, etc.) 5) Transform old file formats to modern formats (e.g. WordPerfect to PDF) 6) E-discovery 7) Search and retrieve to respond to special request 8) Search and retrieve of public records by public users 	
Current solutions	Compute(System)	Linux servers
	Storage	NetApps, Hitachi, Magnetic tapes.
	Networking	
	Software	Custom software, commercial search products, commercial databases.
Big data characteristics	Data source (distributed/centralized)	Distributed data sources from federal agencies. Current solution requires transfer of those data to a centralized storage. In the future, those data sources may reside in different Cloud environments.
	Volume (size)	Hundreds of Terabytes, and growing.
	Velocity (e.g. real time)	Input rate is relatively low compared to other use cases, but the trend is bursty. That is the data can arrive in batches of size ranging from GB to hundreds of TB.

	<p>Variety (multiple data-sets, mashup)</p> <p>Variety data types, unstructured and structured data: textual documents, emails, photos, scanned documents, multimedia, social networks, web sites, databases, etc.</p> <p>Variety of application domains, since records come from different agencies.</p> <p>Data come from variety of repositories, some of which can be cloud-based in the future.</p>
	<p>Variability (rate of change)</p> <p>Rate can change especially if input sources are variable, some having audio, video more, some more text, and other images, etc.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p> <p>Search results should have high relevancy and high recall. Categorization of records should be highly accurate.</p>
	<p>Visualization</p> <p>TBD</p>
	<p>Data quality (syntax)</p> <p>Unknown.</p>
	<p>Data types</p> <p>Variety data types: textual documents, emails, photos, scanned documents, multimedia, databases, etc.</p>
	<p>Data analytics</p> <p>Crawl/index; search; ranking; predictive search.</p> <p>Data categorization (sensitive, confidential, etc.)</p> <p>Personally Identifiable Information (PII) data detection and flagging.</p>
<p>Big data specific challenges (Gaps)</p>	<p>Perform pre-processing and manage for long-term of large and varied data.</p> <p>Search huge amount of data.</p> <p>Ensure high relevancy and recall.</p> <p>Data sources may be distributed in different clouds in future.</p>
<p>Big data specific challenges in mobility</p>	<p>Mobile search must have similar interfaces/results</p>
<p>Security and privacy technical considerations</p>	<p>Need to be sensitive to data access restrictions.</p>
<p>Highlight issues for generalizing this Use case (e.g. for ref. architecture)</p>	
<p>More information (URLs)</p>	

A.1.3 Use case 3: Statistical Survey Response Improvement

Use case title	Statistical Survey Response Improvement (Adaptive Design)	
Vertical (area)	Government Statistical Logistics	
Author/company/email	Cavan Capps: U.S. Census Bureau/cavan.paul.capps@census.gov	
Actors/stakeholders and their roles and responsibilities	U.S. statistical agencies are charged to be the leading authoritative sources about the nation's people and economy, while honoring privacy and rigorously protecting confidentiality. This is done by working with states, local governments and other government agencies.	
Goals	To use advanced methods, that are open and scientifically objective, the statistical agencies endeavor to improve the quality, the specificity and the timeliness of statistics provided while reducing operational costs and maintaining the confidentiality of those measured.	
Use case description	Survey costs are increasing as survey response declines. The goal of this work is to use advanced "recommendation system techniques" using data mashed up from several sources and historical survey para-data to drive operational processes in an effort to increase quality and reduce the cost of field surveys.	
Current solutions	Compute(System)	Linux systems
	Storage	SAN and Direct Storage
	Networking	Fiber, 10 gigabit Ethernet, Infiniband 40 gigabit.
	Software	Hadoop, Spark, Hive, R, SAS, Mahout, Allegrograph, MySQL, Oracle, Storm, BigMemory, Cassandra, Pig
Big data characteristics	Data source (distributed/centralized)	Survey data, other government administrative data, geographical positioning data from various sources.
	Volume (size)	For this particular class of operational problem approximately one petabyte.
	Velocity (e.g. real time)	Varies, paradata from field data streamed continuously, during the decennial census approximately 150 million records transmitted.
	Variety (multiple datasets, mashup)	Data is typically defined strings and numerical fields. Data can be from multiple datasets mashed together for analytical use.
	Variability (rate of change)	Varies depending on surveys in the field at a given time. High rate of velocity during a decennial census.

Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Data must have high veracity and systems must be very robust. The semantic integrity of conceptual metadata concerning what exactly is measured and the resulting limits of inference remain a challenge
	Visualization	Data visualization is useful for data review, operational activity and general analysis. It continues to evolve.
	Data quality (syntax)	Data quality should be high and statistically checked for accuracy and reliability throughout the collection process.
	Data types	Pre-defined ASCII strings and numerical data
	Data analytics	Analytics are required for recommendation systems, continued monitoring and general survey improvement.
Big data specific challenges (Gaps)	Improving recommendation systems that reduce costs and improve quality while providing confidentiality safeguards that are reliable and publically auditable.	
Big data specific challenges in mobility	Mobile access is important.	
Security and privacy technical considerations	All data must be both confidential and secure. All processes must be auditable for security and confidentiality as required by various legal statutes	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Recommender systems have features in common to e-commerce like Amazon, Netflix, UPS etc.	
More information (URLs)		

A.1.4 Use case 4: Non Traditional Data in Statistical Survey

Use case title	Non Traditional Data in Statistical Survey Response Improvement (Adaptive Design)
Vertical (area)	Government Statistical Logistics
Author/company/email	Cavan Capps: U.S. Census Bureau/ cavan.paul.capps@census.gov
Actors/stakeholders and their roles and responsibilities	U.S. statistical agencies are charged to be the leading authoritative sources about the nation's people and economy, while honoring privacy and rigorously protecting confidentiality. This is done by working with states, local governments and other government agencies.
Goals	To use advanced methods, that are open and scientifically objective, the statistical agencies endeavor to improve the quality, the specificity and the timeliness of statistics provided while reducing operational costs and maintaining the confidentiality of those measured.

Use case description	Survey costs are increasing as survey response declines. The potential of using non-traditional commercial and public data sources from the web, wireless communication, electronic transactions mashed up analytically with traditional surveys to improve statistics for small area geographies, new measures and to improve the timeliness of released statistics.	
Current solutions	Compute(System)	Linux systems
	Storage	SAN and Direct Storage
	Networking	Fiber, 10 gigabit Ethernet, Infiniband 40 gigabit.
	Software	Hadoop, Spark, Hive, R, SAS, Mahout, Allegro-graph, MySQL, Oracle, Storm, BigMemory, Cassandra, Pig
Big data characteristics	Data source (distributed/centralized)	Survey data, other government administrative data, web scrapped data, wireless data, e-transaction data, potentially social media data and positioning data from various sources.
	Volume (size)	TBD
	Velocity (e.g. real time)	TBD
	Variety (multiple datasets, mashup)	Textual data as well as the traditionally defined strings and numerical fields. Data can be from multiple datasets mashed together for analytical use.
	Variability (rate of change)	TBD.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Data must have high veracity and systems must be very robust. The semantic integrity of conceptual metadata concerning what exactly is measured and the resulting limits of inference remain a challenge
	Visualization	Data visualization is useful for data review, operational activity and general analysis. It continues to evolve.
	Data quality (syntax)	Data quality should be high and statistically checked for accuracy and reliability throughout the collection process.

	Data types	Textual data, pre-defined ASCII strings and numerical data
	Data analytics	Analytics are required to create reliable estimates using data from traditional survey sources, government administrative data sources and non-traditional sources from the digital economy.
Big data specific challenges (Gaps)	Improving analytic and modeling systems that provide reliable and robust statistical estimated using data from multiple sources that are scientifically transparent and while providing confidentiality safeguards that are reliable and publically auditable.	
Big data specific challenges in mobility	Mobile access is important.	
Security and privacy technical considerations	All data must be both confidential and secure. All processes must be auditable for security and confidentiality as required by various legal statutes.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Statistical estimation that provide more detail, on a more near real time basis for less cost. The reliability of estimated statistics from such “mashed up” sources still must be evaluated.	
More information (URLs)		

A.2 Commercial

A.2.1 Use case 5: Cloud Computing in Financial Industries

Use case title	This use case represents one approach to implementing a BD (Big data) strategy, within a Cloud Eco System, for FI (Financial Industries) transacting business within the United States.
Vertical (area)	<p>The following lines of business (LOB) include:</p> <p>Banking, including: Commercial, Retail, Credit Cards, Consumer Finance, Corporate Banking, Transaction Banking, Trade Finance, and Global Payments.</p> <p>Securities and Investments, such as; Retail Brokerage, Private Banking/Wealth Management, Institutional Brokerages, Investment Banking, Trust Banking, Asset Management, Custody and Clearing Services</p> <p>Insurance, including; Personal and Group Life, Personal and Group Property/Casualty, Fixed and Variable Annuities, and Other Investments</p> <p>Please Note: Any Public/Private entity, providing financial services within the regulatory and jurisdictional risk and compliance purview of the United States, are required to satisfy a complex multilayer number of regulatory governance, risk management, and compliance (GRC)/ confidentiality, integrity, and availability (CIA) requirements, as overseen by various jurisdictions and agencies, including; Fed., State, Local and cross-border.</p>
Author/company/email	Pw Carey, Compliance Partners, LLC, pwc.pwcarey@email.com

<p>Actors/stakeholders and their roles and responsibilities</p>	<p>Regulatory and advisory organizations and agencies including the; SEC (Securities and Exchange Commission), FDIC (Federal Deposit Insurance Corporation), CFTC (Commodity Futures Trading Commission), US Treasury, PCAOB (Public Company Accounting and Oversight Board), COSO, CobiT, reporting supply chains and stakeholders, investment community, shareholders, pension funds, executive management, data custodians, and employees.</p> <p>At each level of a financial services organization, an inter-related and inter-dependent mix of duties, obligations and responsibilities are in-place, which are directly responsible for the performance, preparation and transmittal of financial data, thereby satisfying both the regulatory GRC and CIA of their organizations financial data. This same information is directly tied to the continuing reputation, trust and survivability of an organization's business.</p>
<p>Goals</p>	<p>The following represents one approach to developing a workable BD/ FI strategy within the financial services industry. Prior to initiation and switch-over, an organization must perform the following baseline methodology for utilizing BD/FI within a Cloud Eco-system for both public and private financial entities offering financial services within the regulatory confines of the United States; Federal, State, Local and/or cross-border such as the UK, EU and China.</p> <p>Each financial services organization must approach the following disciplines supporting their BD/FI initiative, with an understanding and appreciation for the impact each of the following four overlaying and inter-dependent forces will play in a workable implementation.</p> <p>These four areas are:</p> <ol style="list-style-type: none"> 1) People (resources), 2) Processes (time/cost/ROI), 3) Technology (various operating systems, platforms and footprints) and 4) Regulatory Governance (subject to various and multiple regulatory agencies). <p>In addition, these four areas must work through the process of being; identified, analyzed, evaluated, addressed, tested, and reviewed in preparation for attending to the following implementation phases:</p> <ol style="list-style-type: none"> 1) Project Initiation and Management Buy-in 2) Risk Evaluations and Controls 3) Business Impact Analysis 4) Design, Development and Testing of the Business Continuity Strategies 5) Emergency Response and Operations (aka; Disaster Recovery) 6) Developing and Implementing Business Continuity Plans 7) Awareness and Training Programs 8) Maintaining and Exercising Business Continuity, (aka: Maintaining Regulatory Currency) <p>Please Note: Whenever appropriate, these eight areas should be tailored and modified to fit the requirements of each organizations unique and specific corporate culture and line of financial services.</p>

<p>Use case description</p>	<p>Big data as developed by Google was intended to serve as an Internet Web site indexing tool to help them sort, shuffle, categorize and label the Internet. At the outset, it was not viewed as a replacement for legacy IT data infrastructures. With the spin-off development within OpenGroup and Hadoop, big data has evolved into a robust data analysis and storage tool that is still undergoing development. However, in the end, big data is still being developed as an adjunct to the current IT client/server/big iron data warehouse architectures which is better at some things, than these same data warehouse environments, but not others.</p> <p>Currently within FI, BD/Hadoop is used for fraud detection, risk analysis and assessments as well as improving the organizations knowledge and understanding of the customers via a strategy known as.... 'know your customer', pretty clever, eh?</p> <p>However, this strategy still must following a well thought out taxonomy that satisfies the entities unique, and individual requirements. One such strategy is the following formal methodology which address two fundamental yet paramount questions; "What are we doing"? and "Why are we doing it"?</p> <ol style="list-style-type: none"> 1) Policy Statement/Project Charter (Goal of the Plan, Reasons and Resources....define each), 2) Business Impact Analysis (how does effort improve our business services), 3) Identify System-wide Policies, Procedures and Requirements, 4) Identify Best Practices for Implementation (including Change Management/Configuration Management) and/or Future Enhancements, 5) Plan B-Recovery Strategies (how and what will need to be recovered, if necessary), 6) Plan Development (Write the Plan and Implement the Plan Elements), 7) Plan buy-in and Testing (important everyone Knows the Plan, and Knows What to Do), and 8) Implement the Plan (then identify and fix gaps during first 3 months, 6 months, and annually after initial implementation) 9) Maintenance (Continuous monitoring and updates to reflect the current enterprise environment) 10) Lastly, System Retirement 	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>Currently, big data/Hadoop within a Cloud Eco-system within the FI is operating as part of a hybrid system, with BD being utilized as a useful tool for conducting risk and fraud analysis, in addition to assisting in organizations in the process of ('know your customer'). These are three areas where BD has proven to be good at;</p> <ol style="list-style-type: none"> 1) detecting fraud, 2) associated risks and a 3) 'know your customer' strategy. <p>At the same time, the traditional client/server/data warehouse/RDBMS are used for the handling, processing, storage and archival of the entities financial data. Recently the SEC has approved the initiative for requiring the FI to submit financial statements via the XBRL (extensible Business Related Markup Language), as of May 13th, 2013.</p>

	<p>Storage</p>	<p>The same Federal, State, Local and cross-border legislative and regulatory requirements can impact any and all geographical locations, including; VMware, NetApps, Oracle, IBM, Brocade, et cetera.</p> <p>Please Note: Based upon legislative and regulatory concerns, these storage solutions for FI data must ensure this same data conforms to US regulatory compliance for GRC/CIA, at this point in time.</p> <p>For confirmation, please visit the following agencies web sites: SEC (U.S. Security and Exchange Commission), CFTC (U.S. Commodity Futures Trading Commission), FDIC (U.S. Federal Deposit Insurance Corporation), DOJ (U.S. Department of Justice), and my favorite the PCAOB (Public Company Accounting and Oversight Board).</p>
	<p>Networking</p>	<p>Please Note: The same Federal, State, Local and cross-border legislative and regulatory requirements can impact any and all geographical locations of HW/SW, including but not limited to; WANs, LANs, MANs WiFi, fiber optics, Internet Access, via Public, Private, Community and Hybrid Cloud environments, with or without VPNs.</p> <p>Based upon legislative and regulatory concerns, these networking solutions for FI data must ensure this same data conforms to US regulatory compliance for GRC/CIA, such as the US Treasury Dept., at this point in time.</p> <p>For confirmation, please visit the following agencies web sites: SEC, CFTC, FDIC, US Treasury Dept., DOJ, and my favorite the PCAOB (Public Company Accounting and Oversight Board).</p>
	<p>Software</p>	<p>Please Note: The same legislative and regulatory obligations impacting the geographical location of HW/SW, also restricts the location for; Hadoop, Map/Reduce, Open-source, and/or Vendor Proprietary such as AWS (Amazon Web Services), Google Cloud Services, and Microsoft</p> <p>Based upon legislative and regulatory concerns, these software solutions incorporating both SOAP (Simple Object Access Protocol), for Web development and OLAP (online analytical processing) software language for databases, specifically in this case for FI data, both must ensure this same data conforms to US regulatory compliance for GRC/CIA, at this point in time.</p> <p>For confirmation, please visit the following agencies web sites: SEC, CFTC, U.S. Treasury, FDIC, DOJ, and my favorite the PCAOB (Public Company Accounting and Oversight Board).</p>

Big data characteristics	Data source (distributed/centralized)	<p>Please Note: The same legislative and regulatory obligations impacting the geographical location of HW/SW, also impacts the location for; both distributed/centralized data sources flowing into HA/DR Environment and HVSS (Hosted Virtual Servers), such as the following constructs: DC1---> VMWare/KVM (Clusters, w/Virtual Firewalls), Data link-Vmware Link-Vmotion Link-Network Link, Multiple PB of NaaS (Network as a Service), DC2--->, VMWare/KVM (Clusters w/Virtual Firewalls), DataLink (Vmware Link, Vmotion Link, Network Link), Multiple PB of NaaS, (Requires Fail-Over Virtualization), among other considerations.</p> <p>Based upon legislative and regulatory concerns, these data source solutions, either distributed and/or centralized for FI data, must ensure this same data conforms to US regulatory compliance for GRC/CIA, at this point in time.</p> <p>For confirmation, please visit the following agencies web sites: SEC, CFTC, US Treasury, FDIC, DOJ, and my favorite the PCAOB (Public Company Accounting and Oversight Board).</p>
	Volume (size)	<p>Tera-bytes up to Peta-bytes.</p> <p>Please Note: This is a 'Floppy Free Zone'.</p>
	Velocity (e.g. real time)	<p>Velocity is more important for fraud detection, risk assessments and the 'know your customer' initiative within the BD FI.</p> <p>Please Note: However, based upon legislative and regulatory concerns, velocity is not at issue regarding BD solutions for FI data, except for fraud detection, risk analysis and customer analysis.</p> <p>Based upon legislative and regulatory restrictions, velocity is not at issue, rather the primary concern for FI data, is that it must satisfy all US regulatory compliance obligations for GRC/CIA, at this point in time.</p>
	Variety (multiple datasets, mashup)	<p>Multiple virtual environments either operating within a batch processing architecture or a hot-swappable parallel architecture supporting fraud detection, risk assessments and customer service solutions.</p> <p>Please Note: Based upon legislative and regulatory concerns, variety is not at issue regarding BD solutions for FI data within a Cloud Eco-system, except for fraud detection, risk analysis and customer analysis.</p> <p>Based upon legislative and regulatory restrictions, variety is not at issue, rather the primary concern for FI data, is that it must satisfy all US regulatory compliance obligations for GRC/CIA, at this point in time.</p>

	<p>Variability (rate of change)</p>	<p>Please Note: Based upon legislative and regulatory concerns, variability is not at issue regarding BD solutions for FI data within a Cloud Eco-system, except for fraud detection, risk analysis and customer analysis.</p> <p>Based upon legislative and regulatory restrictions, variability is not at issue, rather the primary concern for FI data, is that it must satisfy all US regulatory compliance obligations for GRC/CIA, at this point in time.</p> <p>Variability with BD FI within a Cloud Eco-System will depending upon the strength and completeness of the SLA agreements, the costs associated with (CapEx), and depending upon the requirements of the business.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Please Note: Based upon legislative and regulatory concerns, veracity is not at issue regarding BD solutions for FI data within a Cloud Eco-system, except for fraud detection, risk analysis and customer analysis.</p> <p>Based upon legislative and regulatory restrictions, veracity is not at issue, rather the primary concern for FI data, is that it must satisfy all US regulatory compliance obligations for GRC/CIA, at this point in time.</p> <p>Within a big data Cloud Eco-System, data integrity is important over the entire life cycle of the organization due to regulatory and compliance issues related to individual data privacy and security, in the areas of CIA and GRC requirements.</p>
	<p>Visualization</p>	<p>Please Note: Based upon legislative and regulatory concerns, visualization is not at issue regarding BD solutions for FI data, except for fraud detection, risk analysis and customer analysis, FI data is handled by traditional client/server/data warehouse big iron servers.</p> <p>Based upon legislative and regulatory restrictions, visualization is not at issue, rather the primary concern for FI data, is that it must satisfy all US regulatory compliance obligations for GRC/CIA, at this point in time.</p> <p>Data integrity within BD is critical and essential over the entire life-cycle of the organization due to regulatory and compliance issues related to CIA and GRC requirements.</p>

	<p>Data quality (syntax) Please Note: Based upon legislative and regulatory concerns, data quality will always be an issue, regardless of the industry or platform.</p> <p>Based upon legislative and regulatory restrictions, data quality is at the core of data integrity, and is the primary concern for FI data, in that it must satisfy all US regulatory compliance obligations for GRC/CIA, at this point in time.</p> <p>For BD/FI data, data integrity is critical and essential over the entire life-cycle of the organization due to regulatory and compliance issues related to CIA and GRC requirements.</p>
	<p>Data types Please Note: Based upon legislative and regulatory concerns, data types is important in that it must have a degree of consistency and especially survivability during audits and digital forensic investigations where the data format deterioration can negatively impact both an audit and a forensic investigation when passed through multiple cycles.</p> <p>For BD/FI data, multiple data types and formats, include but is not limited to; flat files, .txt, .pdf, android application files, .wav, .jpg and VOIP (Voice over IP)</p> <p>Data analytics Please Note: Based upon legislative and regulatory concerns, data analytics is an issue regarding BD solutions for FI data, especially in regards to fraud detection, risk analysis and customer analysis.</p> <p>However, data analytics for FI data is currently handled by traditional client/server/data warehouse big iron servers which must ensure they comply with and satisfy all United States GRC/CIA requirements, at this point in time.</p> <p>For BD/FI data analytics must be maintained in a format that is non-destructive during search and analysis processing and procedures.</p>
<p>Big data specific challenges (Gaps)</p>	<p>Currently, the areas of concern associated with BD/FI with a Cloud Eco-system, include the aggregating and storing of data (sensitive, toxic and otherwise) from multiple sources which can and does create administrative and management problems related to the following:</p> <ul style="list-style-type: none"> — Access control — Management/Administration — Data entitlement and — Data ownership <p>However, based upon current analysis, these concerns and issues are widely known and are being addressed at this point in time, via the Research and Development SDLC/HDLC (Software Development Life Cycle/Hardware Development Life Cycle) sausage makers of technology. Please stay tuned for future developments in this regard</p>

<p>Big data specific challenges in mobility</p>	<p>Mobility is a continuously growing layer of technical complexity; however, not all big data mobility solutions are technical in nature. There are two interrelated and co-dependent parties who required to work together to find a workable and maintainable solution, the FI business side and IT. When both are in agreement sharing a, common lexicon, taxonomy and appreciation and understand for the requirements each is obligated to satisfy, these technical issues can be addressed.</p> <p>Both sides in this collaborative effort will encounter the following current and on-going FI data considerations:</p> <ul style="list-style-type: none"> — Inconsistent category assignments — Changes to classification systems over time — Use of multiple overlapping or — Different categorization schemes <p>In addition, each of these changing and evolving inconsistencies, are required to satisfy the following data characteristics associated with ACID:</p> <ul style="list-style-type: none"> — Atomic - All of the work in a transaction completes (commit) or none of it completes — Consistent - A transmittal transforms the database from one consistent state to another consistent state. Consistency is defined in terms of constraints. — Isolated - The results of any changes made during a transaction are not visible until the transaction has committed. — Durable - The results of a committed transaction survive failures. <p>When each of these data categories is satisfied, well, it's a glorious thing. Unfortunately, sometimes glory is not in the room, however, that does not mean we give up the effort to resolve these issues.</p>
<p>Security and privacy technical considerations</p>	<p>No amount of security and privacy due diligence will make up for the innate deficiencies associated with human nature that creep into any program and/or strategy. Currently, the BD/FI must contend with a growing number of risk buckets, such as:</p> <ul style="list-style-type: none"> — AML-Anti-Money Laundering — CDD- Client Due Diligence — Watch-lists — FCPA – Foreign Corrupt Practices Act <p>...to name a few.</p> <p>For a reality check, please consider Mr. Harry M. Markopolos' nine-year effort to get the SEC among other agencies to do their job and shut down Mr. Bernard Madoff's billion dollar Ponzi scheme.</p>

	<p>However, that aside, identifying and addressing the privacy/security requirements of the FI, providing services within a BD/Cloud Eco-system, via continuous improvements in:</p> <ol style="list-style-type: none">1) technology,2) processes,3) procedures,4) people and5) regulatory jurisdictions <p>...is a far better choice for both the individual and the organization, especially when considering the alternative.</p> <p>Utilizing a layered approach, this strategy can be broken down into the following sub categories:</p> <ol style="list-style-type: none">1) Maintaining operational resilience2) Protecting valuable assets3) Controlling system accounts4) Managing security services effectively, and5) Maintaining operational resilience <p>For additional background security and privacy solutions addressing both security and privacy, we'll refer you to the two following organization's:</p> <ul style="list-style-type: none">— ISACA (International Society of Auditors and Computer Analysts)— isc2 (International Security Computer and Systems Auditors)
--	--

<p>Highlight issues for generalizing this Use case (e.g. for ref. architecture)</p>	<p>Areas of concern include the aggregating and storing data from multiple sources can create problems related to the following:</p> <ul style="list-style-type: none"> — Access control — Management/Administration — Data entitlement and — Data ownership <p>Each of these areas is being improved upon, yet they still must be considered and addressed, via access control solutions, and SIEM (Security Incident/Event Management) tools.</p> <p>I don't believe we're there yet, based upon current security concerns mentioned whenever big data/Hadoop within a Cloud Eco-system is brought up in polite conversation.</p> <p>Current and on-going challenges to implementing BD Finance within a Cloud Eco, as well as traditional client/server data warehouse architectures, include the following areas of Financial Accounting under both US GAAP (U.S. Generally Accepted Accounting Practices) or IFRS (International Financial Reporting Standards):</p> <p>XBRL (extensible Business Related Markup Language)</p> <p>Consistency (terminology, formatting, technologies, regulatory gaps)</p> <p>SEC mandated use of XBRL (extensible Business Related Markup Language) for regulatory financial reporting.</p> <p>SEC, GAAP/IFRS and the yet to be fully resolved new financial legislation impacting reporting requirements are changing and point to trying to improve the implementation, testing, training, reporting and communication best practices required of an independent auditor, regarding:</p> <p>Auditing, Auditor's reports, Control self-assessments, Financial audits, GAAS / ISAs, Internal audits, and the Sarbanes–Oxley Act of 2002 (SOX).</p>
--	--

<p>More information (URLs)</p>	<ol style="list-style-type: none"> 1) Cloud Security Alliance big data Working Group, "Top 10 Challenges in big data Security and Privacy", 2012. 2) The IFRS, Securities and Markets Working Group, http://www.xbrl-eu.org 3) IEEE Big data conference http://www.ischool.drexel.edu/bigdata/bigdata2013/topics.htm 4) Map/Reduce http://www.mapreduce.org. 5) PCAOB http://www.pcaob.org 6) http://www.ey.com/GL/en/Industries/Financial-Services/Insurance 7) http://www.treasury.gov/resource-center/fin-mkts/Pages/default.aspx 8) CFTC http://www.cftc.org 9) SEC http://www.sec.gov 10) FDIC http://www.fdic.gov 11) COSO http://www.coso.org 12) isc2 International Information Systems Security Certification Consortium, Inc.: http://www.isc2.org 13) ISACA Information Systems Audit and Control Association: http://www.isca.org 14) IFARS http://www.ifars.org 15) Apache http://www.opengroup.org 16) http://www.computerworld.com/s/article/print/9221652/IT_must_prepare_for_Hadoop_security_issues?tax ... 17) "No One Would Listen: A True Financial Thriller" (hard-cover book). Hoboken, NJ: John Wiley & Sons. March 2010. Retrieved April 30, 2010. ISBN 978-0-470-55373-2 18) Assessing the Madoff Ponzi Scheme and Regulatory Failures (Archive of: Subcommittee on Capital Markets, Insurance, and Government Sponsored Enterprises Hearing): (http://financialserv.edgeboss.net/wmedia//hearing020409.wvx) (Windows Media). U.S. House Financial Services Committee. February 4, 2009. Retrieved June 29, 2009. 19) COSO, The Committee of Sponsoring Organizations of the Treadway Commission (COSO), Copyright© 2013, http://www.coso.org. 20) (ITIL) Information Technology Infrastructure Library, Copyright© 2007-13 APM Group Ltd. All rights reserved, Registered in England No. 2861902, http://www.itil-officialsite.com. 21) CobiT, Ver. 5.0, 2013, ISACA, Information Systems Audit and Control Association, (a framework for IT Governance and Controls), http://www.isaca.org. 22) TOGAF, Ver. 9.1, The Open Group Architecture Framework (a framework for IT architecture), http://www.opengroup.org. 23) ISO/IEC 27000:2016Info. Security Mgt., International Organization for Standardization and the International Electrotechnical Commission, http://www.standards.iso.org/
---------------------------------------	---

NOTE Please feel free to improve our **INITIAL DRAFT, Ver. 0.1, August 25th, 2013**....as we do not consider our efforts to be pearls, at this point in time.....Respectfully yours, Pw Carey, Compliance Partners, LLC_pwc.pwcarey@gmail.com

A.2.2 Use case 6: Mendeley—An International Network of Research

Use case title	Mendeley – An International Network of Research	
Vertical (area)	Commercial Cloud Consumer Services	
Author/company/email	William Gunn / Mendeley / william.gunn@mendeley.com	
Actors/stakeholders and their roles and responsibilities	Researchers, librarians, publishers, and funding organizations.	
Goals	To promote more rapid advancement in scientific research by enabling researchers to efficiently collaborate, librarians to understand researcher needs, publishers to distribute research findings more quickly and broadly, and funding organizations to better understand the impact of the projects they fund.	
Use case description	Mendeley has built a database of research documents and facilitates the creation of shared bibliographies. Mendeley uses the information collected about research reading patterns and other activities conducted via the software to build more efficient literature discovery and analysis tools. Text mining and classification systems enables automatic recommendation of relevant research, improving the cost and performance of research teams, particularly those engaged in curation of literature on a particular subject, such as the Mouse Genome Informatics group at Jackson Labs, which has a large team of manual curators who scan the literature. Other use cases include enabling publishers to more rapidly disseminate publications, facilitating research institutions and librarians with data management plan compliance, and enabling funders to better understand the impact of the work they fund via real-time data on the access and use of funded research.	
Current solutions	Compute(System)	Amazon EC2
	Storage	HDFS Amazon S3
	Networking	Client-server connections between Mendeley and end user machines, connections between Mendeley offices and Amazon services.
	Software	Hadoop, Scribe, Hive, Mahout, Python
Big data characteristics	Data source (distributed/centralized)	Distributed and centralized
	Volume (size)	15 TB presently, growing about 1 TB/month
	Velocity (e.g. real time)	Currently Hadoop batch jobs are scheduled daily, but work has begun on real-time recommendation

	Variety (multiple datasets, mashup)	PDF documents and log files of social network and client activities
	Variability (rate of change)	Currently a high rate of growth as more researchers sign up for the service, highly fluctuating activity over the course of the year
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Metadata extraction from PDFs is variable, it's challenging to identify duplicates, there's no universal identifier system for documents or authors (though ORCID proposes to be this)
	Visualization	Network visualization via Gephi, scatterplots of readership vs. citation rate, etc.
	Data quality (syntax)	90% correct metadata extraction according to comparison with Crossref, Pubmed, and Arxiv
	Data types	Mostly PDFs, some image, spreadsheet, and presentation files
	Data analytics	Standard libraries for machine learning and analytics, LDA, custom built reporting tools for aggregating readership and social activities per document
Big data specific challenges (Gaps)	The database contains ~400 M documents, roughly 80 M unique documents, and receives 5-700 k new uploads on a weekday. Thus a major challenge is clustering matching documents together in a computationally efficient way (scalable and parallelized) when they're uploaded from different sources and have been slightly modified via third-part annotation tools or publisher watermarks and cover pages	
Big data specific challenges in mobility	Delivering content and services to various computing platforms from Windows desktops to Android and iOS mobile devices	

Security and privacy technical considerations	Researchers often want to keep what they're reading private, especially industry researchers, so the data about who's reading what has access controls.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	This use case could be generalized to providing content-based recommendations to various scenarios of information consumption
More information (URLs)	— Mendeley. http://mendeley.com . Accessed March 3, 2015. — Mendeley. http://dev.mendeley.com . Accessed March 3, 2015.

A.2.3 Use case 7: Netflix Movie Service

Use case title	Netflix Movie Service	
Vertical (area)	Commercial Cloud Consumer Services	
Author/company/email	Geoffrey Fox, Indiana University gcf@indiana.edu	
Actors/stakeholders and their roles and responsibilities	Netflix Company (Grow sustainable Business), Cloud Provider (Support streaming and data analysis), Client user (Identify and watch good movies on demand)	
Goals	Allow streaming of user selected movies to satisfy multiple objectives (for different stakeholders) — especially retaining subscribers. Find best possible ordering of a set of videos for a user (household) within a given context in real time; maximize movie consumption.	
Use case description	Digital movies stored in cloud with metadata; user profiles and rankings for small fraction of movies for each user. Use multiple criteria – content based recommender system, user-based recommender system; diversity. Refine algorithms continuously with A/B testing.	
Current solutions	Compute (System)	Amazon Web Services AWS
	Storage	Uses Cassandra NoSQL technology with Hive, Teradata
	Networking	Need Content Delivery System to support effective streaming video
	Software	Hadoop and Pig; Cassandra; Teradata
Big data characteristics	Data source (distributed/centralized)	Add movies institutionally. Collect user rankings and profiles in a distributed fashion
	Volume (size)	Summer 2012. 25 million subscribers; 4 million ratings per day; 3 million searches per day; 1 billion hours streamed in June 2012. Cloud storage 2 petabytes (June 2013)
	Velocity (e.g. real time)	Media (video and properties) and Rankings continually updated
	Variety (multiple datasets, mashup)	Data varies from digital media to user rankings, user profiles and media properties for content-based recommendations

	Variability (rate of change)	Very competitive business. Need to aware of other companies and trends in both content (which Movies are hot) and technology. Need to investigate new business initiatives such as Netflix sponsored content
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Success of business requires excellent quality of service
	Visualization	Streaming media and quality user-experience to allow choice of content
	Data quality (syntax)	Rankings are intrinsically “rough” data and need robust learning algorithms
	Data types	Media content, user profiles, “bag” of user rankings
	Data analytics	Recommender systems and streaming video delivery. Recommender systems are always personalized and use logistic/linear regression, elastic nets, matrix factorization, clustering, latent Dirichlet allocation, association rules, gradient boosted decision trees and others. Winner of Netflix competition (to improve ratings by 10 %) combined over 100 different algorithms.
Big data specific challenges (Gaps)	Analytics needs continued monitoring and improvement.	
Big data specific challenges in mobility	Mobile access important	
Security and privacy technical considerations	Need to preserve privacy for users and digital rights for media.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Recommender systems have features in common to e-commerce like Amazon. Streaming video has features in common with other content providing services like iTunes, Google Play, Pandora and Last.fm	
More information (URLs)	<p>— Building Large-scale Real-world Recommender Systems - Recsys2012 tutorial. http://www.slideshare.net/xamat/building-largescale-realworld-recommender-systems-recsys2012-tutorial. Accessed March 3, 2015.</p> <p>— RAD – Outlier Detection on big data. http://techblog.netflix.com/. Accessed March 3, 2015.</p>	

A.2.4 Use case 8: Web Search

Use case title	Web Search (Bing, Google, Yahoo...)
Vertical (area)	Commercial Cloud Consumer Services
Author/company/email	Geoffrey Fox, Indiana University gcf@indiana.edu

Actors/stakeholders and their roles and responsibilities	Owners of web information being searched; search engine companies; advertisers; users	
Goals	Return in ~0,1 s, the results of a search based on average of 3 words; important to maximize “precision@10”; number of great responses in top 10 ranked results	
Use case description	1) Crawl the web; 2) Pre-process data to get searchable things (words, positions); 3) Form Inverted Index mapping words to documents; 4) Rank relevance of documents: PageRank; 5) Lots of technology for advertising, “reverse engineering ranking” “preventing reverse engineering”; 6) Clustering of documents into topics (as in Google News) 7) Update results efficiently	
Current solutions	Compute(System)	Large Clouds
	Storage	Inverted Index not huge; crawled documents are petabytes of text – rich media much more
	Networking	Need excellent external network links; most operations pleasingly parallel and I/O sensitive. High performance internal network not needed
	Software	Map/Reduce + Bigtable; Dryad + Cosmos. PageRank. Final step essentially a recommender engine
Big data characteristics	Data source (distributed/centralized)	Distributed web sites
	Volume (size)	45 B web pages total, 500 M photos uploaded each day, 100 h of video uploaded to YouTube each minute
	Velocity (e.g. real time)	Data continually updated
	Variety (multiple datasets, mashup)	Rich set of functions. After processing, data similar for each page (except for media types)
	Variability (rate of change)	Average page has life of a few months
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Exact results not essential but important to get main hubs and authorities for search query
	Visualization	Not important although page layout critical
	Data quality (syntax)	A lot of duplication and spam
	Data types	Mainly text but more interest in rapidly growing image and video

	Data analytics	Crawling; searching including topic based search; ranking; recommending
Big data specific challenges (Gaps)	Search of “deep web” (information behind query front ends) Ranking of responses sensitive to intrinsic value (as in Pagerank) as well as advertising value Link to user profiles and social network data	
Big data specific challenges in mobility	Mobile search must have similar interfaces/results	
Security and privacy technical considerations	Need to be sensitive to crawling restrictions. Avoid Spam results	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Relation to Information retrieval such as search of scholarly works.	
More information (URLs)	<ul style="list-style-type: none"> — Internet Trends D11 Conference. http://www.slideshare.net/kleinerperkins/kpcb-internet-trends-2013. Accessed March 3, 2015. — Introduction to Search Engine Technology. http://webcourse.cs.technion.ac.il/236621/Winter2011-2012/en/ho_Lectures.html. Accessed March 3, 2015. — Lecture “Information Retrieval and Web Search Engines” (SS 2011). http://www.ifis.cs.tu-bs.de/teaching/ss-11/irws. Accessed March 3, 2015. — Recommender Systems Tutorial (Part 1) -Introduction. http://www.slideshare.net/bee chung/recommender-systems-tutorialpart1intro. Accessed March 3, 2015. — The size of the World Wide Web (The Internet). http://www.worldwidewebsite.com/. Accessed March 3, 2015. 	

A.2.5 Use case 9: Cloud-based Continuity and Disaster Recovery

Use case title	IaaS (Infrastructure as a Service) big data BC/DR Within A Cloud Eco-System provided by Cloud Service Providers (CSPs) and Cloud Brokerage Service Providers (CBSPs)
Vertical (area)	Large Scale Reliable Data Storage
Author/company/email	Pw Carey, Compliance Partners, LLC, pwc.pwcarey@email.com
Actors/stakeholders and their roles and responsibilities	Executive Management, Data Custodians, and Employees responsible for the integrity, protection, privacy, confidentiality, availability, safety, security and survivability of a business by ensuring the 3-As of data accessibility to an organizations services are satisfied; anytime, anyplace and on any device.

<p>Goals</p>	<p>The following represents one approach to developing a workable BC/DR strategy. Prior to outsourcing an organizations BC/DR onto the backs/shoulders of a CSP or CBSP, the organization must perform the following use case, which will provide each organization with a baseline methodology for BC/DR best practices, within a Cloud Eco-system for both Public and Private organizations.</p> <p>Each organization must approach the ten disciplines supporting BC/DR, with an understanding and appreciation for the impact each of the following four overlaying and inter-dependent forces will play in ensuring a workable solution to an entity's business continuity plan and requisite disaster recovery strategy. The four areas are; people (resources), processes (time/cost/ROI), technology (various operating systems, platforms and footprints) and governance (subject to various and multiple regulatory agencies).</p>
	<p>These four concerns must be; identified, analyzed, evaluated, addressed, tested, reviewed, addressed during the following ten phases:</p> <ol style="list-style-type: none"> 1) Project Initiation and Management Buy-in 2) Risk Evaluations and Controls 3) Business Impact Analysis 4) Design, Development and Testing of the Business Continuity Strategies 5) Emergency Response and Operations (aka; Disaster Recovery 6) Developing and Implementing Business Continuity Plans 7) Awareness and Training Programs 8) Maintaining and Exercising Business Continuity Plans, (aka: Maintaining Currency) 9) Public Relations (PR) and Crises Management Plans 10) Coordination with Public Agencies <p>Please Note: When appropriate, these ten areas can be tailored to fit the requirements of the organization.</p>
<p>Use case description</p>	<p>Big data as developed by Google was intended to serve as an Internet Web site indexing tool to help them sort, shuffle, categorize and label the Internet. At the outset, it was not viewed as a replacement for legacy IT data infrastructures. With the spin-off development within OpenGroup and Hadoop, big data has evolved into a robust data analysis and storage tool that is still undergoing development. However, in the end, big data is still being developed as an adjunct to the current IT client/server/big iron data warehouse architectures which is better at some things, than these same data warehouse environments, but not others.</p>

	<p>As a result, it is necessary, within this business continuity/disaster recovery use case, we ask good questions, such as; why are we doing this and what are we trying to accomplish? What are our dependencies upon manual practices and when can we leverage them? What systems have been and remain outsourced to other organizations, such as our Telephony and what are their DR/BC business functions, if any? Lastly, we must recognize the functions that can be simplified and what are the preventative steps we can take that do not have a high cost associated with them such as simplifying business practices.</p> <p>We must identify what are the critical business functions that need to be recovered, 1st, 2nd, 3rd in priority, or at a later time/date, and what is the Model of A Disaster we're trying to resolve, what are the types of disasters more likely to occur realizing that we don't need to resolve all types of disasters. When backing up data within a Cloud Eco-system is a good solution, this will shorten the fail-over time and satisfy the requirements of RTO/RPO. In addition, there must be 'Buy-in', as this is not just an IT problem; it is a business services problem as well, requiring the testing of the Disaster Plan via formal walk-throughs, et cetera. There should be a formal methodology for developing a BC/DR Plan, including: 1). Policy Statement (Goal of the Plan, Reasons and Resources....define each), 2). Business Impact Analysis (how does a shutdown impact the business financially and otherwise), 3). Identify Preventive Steps (can a disaster be avoided by taking prudent steps), 4). Recovery Strategies (how and what you will need to recover), 5). Plan Development (Write the Plan and Implement the Plan Elements), 6). Plan buy-in and Testing (very important so that everyone knows the Plan and knows what to do during its execution), and 7). Maintenance (Continuous changes to reflect the current enterprise environment)</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>Cloud Eco-systems, incorporating IaaS (Infrastructure as a Service), supported by Tier 3 Data Centers....Secure Fault Tolerant (Power).... for Security, Power, Air Conditioning et cetera...geographically off-site data recovery centers... providing data replication services, Note: Replication is different from Backup. Replication only moves the changes since the last time a replication, including block level changes. The replication can be done quickly, with a five second window, while the data is replicated every four hours. This data snap shot is retained for seven business days, or longer if necessary. Replicated data can be moved to a Fail-over Center to satisfy the organizations RPO (Recovery Point Objectives) and RTO</p>
	<p>Storage</p>	<p>VMware, NetApps, Oracle, IBM, Brocade,</p>
	<p>Networking</p>	<p>WANs, LANs, WiFi, Internet Access, via Public, Private, Community and Hybrid Cloud environments, with or without VPNs.</p>
	<p>Software</p>	<p>Hadoop, Map/Reduce, Open-source, and/or Vendor Proprietary such as AWS (Amazon Web Services), Google Cloud Services, and Microsoft</p>
<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>Both distributed/centralized data sources flowing into HA/DR Environment and HVSSs, such as the following: DC1---> VMWare/KVM (Clusters, w/Virtual Firewalls), Data link-VMware Link-Vmotion Link-Network Link, Multiple PB of NaaS, DC2--->, VMWare/KVM (Clusters w/ Virtual Firewalls), DataLink (VMware Link, Motion Link, Network Link), Multiple PB of NaaS, (Requires Fail-Over Virtualization)</p>
	<p>Volume (size)</p>	<p>Terabytes up to Petabytes</p>

	<p>Velocity (e.g. real time)</p>	<p>Tier 3 Data Centers with Secure Fault Tolerant (Power) for Security, Power, and Air Conditioning. IaaS (Infrastructure as a Service) in this example, based upon NetApps. Replication is different from Backup; replication requires only moving the CHANGES since the last time a REPLICATION was performed, including the block level changes. The Replication can be done quickly as the data is Replicated every four hours. These replications can be performed within a 5 second window, and this Snap Shot will be kept for 7 business days, or longer if necessary to a Fail-Over Center.....at the RPO and RTO....</p>
	<p>Variety (multiple datasets, mashup)</p>	<p>Multiple virtual environments either operating within a batch processing architecture or a hot-swappable parallel architecture.</p>
	<p>Variability (rate of change)</p>	<p>Depending upon the SLA agreement, the costs (CapEx) increases, depending upon the RTO/RPO and the requirements of the business.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Data integrity is critical and essential over the entire life-cycle of the organization due to regulatory and compliance issues related to data CIA and GRC data requirements.</p>
	<p>Visualization</p>	<p>Data integrity is critical and essential over the entire life-cycle of the organization due to regulatory and compliance issues related to data CIA and GRC data requirements.</p>
	<p>Data quality (syntax)</p>	<p>Data integrity is critical and essential over the entire life-cycle of the organization due to regulatory and compliance issues related to data CIA and GRC data requirements.</p>
	<p>Data types</p>	<p>Multiple data types and formats, including but not limited to; flat files, .txt, .pdf, android application files, .wav, .jpg and VOIP (Voice over IP)</p>
	<p>Data analytics</p>	<p>Must be maintained in a format that is non-destructive during search and analysis processing and procedures.</p>
<p>Big data specific challenges (Gaps)</p>	<p>The complexities associated with migrating from a Primary Site to either a Replication Site or a Backup Site is not fully automated at this point in time. The goal is to enable the user to automatically initiate the Fail Over Sequence, moving Data Hosted within Cloud requires a well-defined and continuously monitored server configuration management. In addition, both organizations must know which servers have to be restored and what are the dependencies and inter-dependencies between the Primary Site servers and Replication and/or Backup Site servers. This requires a continuous monitoring of both, since there are two solutions involved with this process, either dealing with servers housing stored images or servers running hot all the time, as in running parallel systems with hot-swappable functionality, all of which requires accurate and up-to-date information from the client.</p>	

<p>Big data specific challenges in mobility</p>	<p>Mobility is a continuously growing layer of technical complexity; however, not all DR/BC solutions are technical in nature, as there are two sides required to work together to find a solution, the business side and the IT side. When they are in agreement, these technical issues must be addressed by the BC/DR strategy implemented and maintained by the entire organization. One area, which is not limited to mobility challenges, concerns a fundamental issue impacting most BC/DR solutions. If your Primary Servers (A, B, C) understand X, Y, Z....but your Secondary Virtual Replication/Backup Servers (a, b, c) over the passage of time, are not properly maintained (configuration management) and become out of sync with your Primary Servers, and only understand X, and Y, when called upon to perform a Replication or Back-up, well "Houston, we have a problem..."</p> <p>Please Note: Over time all systems can and will suffer from sync-creep, some more than others, when relying upon manual processes to ensure system stability.</p>
<p>Security and privacy technical considerations</p>	<p>Dependent upon the nature and requirements of the organization's industry verticals, such as; Finance, Insurance, and Life Sciences including both public and/or private entities, and the restrictions placed upon them by; regulatory, compliance and legal jurisdictions.</p>

Copyrighted document, no reproduction or distribution allowed without permission from ISO/IEC. For review by FC on PR in healthcare. Click to view the full PDF of ISO/IEC TR 20547-2:2018. Oct 2024

Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Challenges to Implement BC/DR, include the following: 1) Recognition, a). Management Vision, b). Assuming the issue is an IT issue, when it is not just an IT issue, 2). People: a). Staffing levels - Many SMBs are understaffed in IT for their current workload, b). Vision - (Driven from the Top Down) Can the business and IT resources see the whole problem and craft a strategy such a 'Call List' in case of a Disaster, c). Skills - Are there resources that can architect, implement and test a BC/DR Solution, d). Time - Do Resources have the time and does the business have the Windows of Time for constructing and testing a DR/BC Solution as DR/BC is an additional Add-On Project the organization needs the time and resources. 3). Money - This can be turned in to an OpEx Solution rather than a CapEx Solution which and can be controlled by varying RPO/RTO, a). Capital is always a constrained resource, b). BC Solutions need to start with "what is the Risk" and "how does cost constrain the solution"? 4). Disruption - Build BC/DR into the standard "Cloud" infrastructure (IaaS) of the SMB, a). Planning for BC/DR is disruptive to business resources, b). Testing BC is also disruptive.....
More information (URLs)	1) http://www.disasterrecovery.org/ , (March, 2013). 2) BC_DR From the Cloud, Avoid IT Disasters EN POINTE Technologies and dinCloud, Webinar Presenter Barry Weber, http://www.dincloud.com . 3) COSO, The Committee of Sponsoring Organizations of the Treadway Commission (COSO), Copyright© 2013, http://www.coso.org . 4) ITIL Information Technology Infrastructure Library, Copyright© 2007-13 APM Group Ltd. All rights reserved, Registered in England No. 2861902, http://www.itil-officialsite.com . 5) CobiT, Ver. 5.0, 2013, ISACA, Information Systems Audit and Control Association, (a framework for IT Governance and Controls), http://www.isaca.org . 6) TOGAF, Ver. 9.1, The Open Group Architecture Framework (a framework for IT architecture), http://www.opengroup.org . 7) ISO/IEC 27000:2016Info. Security Mgt., International Organization for Standardization and the International Electrotechnical Commission, http://www.standards.iso.org/ . 8) PCAOB, Public Company Accounting and Oversight Board, http://www.pcaobus.org .
NOTE Please feel free to improve our INITIAL DRAFT, Ver. 0.1, August 10 th , 2013....as we do not consider our efforts to be pearls, at this point in time.....Respectfully yours, Pw Carey, Compliance Partners, LLC_pwc.pwcarey@gmail.com	

A.2.6 Use case 10: Cargo Shipping

Use case title	Cargo Shipping
Vertical (area)	Industry
Author/company/email	William Miller/MaCT USA/mact-usa@att.net
Actors/stakeholders and their roles and responsibilities	End-users (Sender/Recipients) Transport Handlers (Truck/Ship/Plane) Telecom Providers (Cellular/SATCOM) Shippers (Shipping and Receiving)
Goals	Retention and analysis of items (Things) in transport

<p>Use case description</p>	<p>The following use case defines the overview of a big data application related to the shipping industry (i.e. FedEx, UPS, DHL, etc.). The shipping industry represents possible the largest potential use case of big data that is in common use today. It relates to the identification, transport, and handling of item (Things) in the supply chain. The identification of an item begins with the sender to the recipients and for all those in between with a need to know the location and time of arrive of the items while in transport. A new aspect will be status condition of the items which will include sensor information, GPS coordinates, and a unique identification schema based upon a new ISO 29161 standards under development within ISO JTC1 SC31 WG2. The data is in near real time being updated when a truck arrives at a depot or upon delivery of the item to the recipient. Intermediate conditions are not currently known, the location is not updated in real time, items lost in a warehouse or while in shipment represent a problem potentially for homeland security. The records are retained in an archive and can be accessed for xx days.</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>Unknown</p>
	<p>Storage</p>	<p>Unknown</p>
	<p>Networking</p>	<p>LAN/T1/Internet Web Pages</p>
	<p>Software</p>	<p>Unknown</p>
<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>Centralized today</p>
	<p>Volume (size)</p>	<p>Large</p>
	<p>Velocity (e.g. real time)</p>	<p>The system is not currently real time.</p>
	<p>Variety (multiple datasets, mashup)</p>	<p>Updated when the driver arrives at the depot and download the time and date the items were picked up. This is currently not real time.</p>
	<p>Variability (rate of change)</p>	<p>Today the information is updated only when the items that were checked with a bar code scanner are sent to the central server. The location is not currently displayed in real time.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	
	<p>Visualization</p>	<p>NONE</p>
	<p>Data quality (syntax)</p>	<p>YES</p>
	<p>Data types</p>	<p>Not Available</p>
	<p>Data analytics</p>	<p>YES</p>
<p>Big data specific challenges (Gaps)</p>	<p>Provide more rapid assessment of the identity, location, and conditions of the shipments, provide detailed analytics and location of problems in the system in real time.</p>	

Big data specific challenges in mobility	Currently conditions are not monitored on-board trucks, ships, and aircraft
Security and privacy technical considerations	Security need to be more robust
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	This use case includes local data bases as well as the requirement to synchronize with the central server. This operation would eventually extend to mobile device and on-board systems which can track the location of the items and provide real-time update of the information including the status of the conditions, logging, and alerts to individuals who have a need to know.
More information (URLs)	

A.2.7 Use case 11: Materials Data

Use case title	Materials Data
Vertical (area)	Manufacturing, Materials Research
Author/company/email	John Rumble, R&R Data Services; jumbleusa@earthlink.net
Actors/stakeholders and their roles and responsibilities	Product Designers (Inputters of materials data in CAE) Materials Researchers (Generators of materials data; users in some cases) Materials Testers (Generators of materials data; standards developers) Data distributors (Providers of access to materials, often for profit)
Goals	Broaden accessibility, quality, and usability; Overcome proprietary barriers to sharing materials data; Create sufficiently large repositories of materials data to support discovery
Use case description	Every physical product is made from a material that has been selected for its properties, cost, and availability. This translates into hundreds of billion dollars of material decisions made every year. In addition, as the Materials Genome Initiative has so effectively pointed out, the adoption of new materials normally takes decades (two to three) rather than a small number of years, in part because data on new materials is not easily available.

	<p>All actors within the materials life cycle today have access to very limited quantities of materials data, thereby resulting in materials-related decision that are non-optimal, inefficient, and costly. While the Materials Genome Initiative is addressing one major and important aspect of the issue, namely the fundamental materials data necessary to design and test materials computationally, the issues related to physical measurements on physical materials (from basic structural and thermal properties to complex performance properties to properties of novel (nanoscale materials) are not being addressed systematically, broadly (cross-discipline and internationally), or effectively (virtually no materials data meetings, standards groups, or dedicated funded programs).</p> <p>One of the greatest challenges that big data approaches can address is predicting the performance of real materials (gram to ton quantities) starting at the atomistic, nanometer, and/or micrometer level of description.</p> <p>As a result of the above considerations, decisions about materials usage are unnecessarily conservative, often based on older rather than newer materials research and development data, and not taking advantage of advances in modeling and simulations. Materials informatics is an area in which the new tools of data science can have major impact.</p>	
Current solutions	Compute(System)	None
	Storage	Widely dispersed with many barriers to access
	Networking	Virtually none
	Software	Narrow approaches based on national programs (Japan, Korea, and China), applications (EU Nuclear program), proprietary solutions (Granta, etc.)
Big data characteristics	Data source (distributed/centralized)	Extremely distributed with data repositories existing only for a very few fundamental properties
	Volume (size)	It is has been estimated (in the 1980s) that there were over 500,000 commercial materials made in the last fifty years. The last three decades has seen large growth in that number.
	Velocity (e.g. real time)	Computer-designed and theoretically design materials (e.g., nanomaterials) are growing over time

	Variety (multiple datasets, mashup)	Many data sets and virtually no standards for mashups
	Variability (rate of change)	Materials are changing all the time, and new materials data are constantly being generated to describe the new materials
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	More complex material properties can require many (100s?) of independent variables to describe accurately. Virtually no activity no exists that is trying to identify and systematize the collection of these variables to create robust data sets.
	Visualization	Important for materials discovery. Potentially important to understand the dependency of properties on the many independent variables. Virtually unaddressed.
	Data quality (syntax)	Except for fundamental data on the structural and thermal properties, data quality is poor or unknown. See Munro's NIST Standard Practice Guide.
	Data types	Numbers, graphical, images
	Data analytics	Empirical and narrow in scope
Big data specific challenges (Gaps)	<ol style="list-style-type: none"> 1) Establishing materials data repositories beyond the existing ones that focus on fundamental data 2) Developing internationally-accepted data recording standards that can be used by a very diverse materials community, including developers materials test standards (such as ASTM and ISO), testing companies, materials producers, and research and development labs 3) Tools and procedures to help organizations wishing to deposit proprietary materials in data repositories to mask proprietary information, yet to maintain the usability of data 4) Multi-variable materials data visualization tools, in which the number of variables can be quite high 	
Big data specific challenges in mobility	Not important at this time	

Security and privacy technical considerations	Proprietary nature of many data very sensitive.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Development of standards; development of large scale repositories; involving industrial users; integration with CAE (don't underestimate the difficulty of this – materials people are generally not as computer savvy as chemists, bioinformatics people, and engineers)
More information (URLs)	

Copyrighted document, no reproduction or circulation
 STANDARDSISO.COM: Click to view the full PDF of ISO/IEC TR 20547 WG - 2:2018
 For review by FG on AI in healthcare
 Oct 2024

A.2.8 Use case 12: Simulation Driven Materials Genomics

Use case title	Simulation driven Materials Genomics	
Vertical (area)	Scientific Research: Materials Science	
Author/company/email	David Skinner/LBNL/deskinner@lbl.gov	
Actors/stakeholders and their roles and responsibilities	<p>Capability providers: National labs and energy hubs provide advanced materials genomics capabilities using computing and data as instruments of discovery.</p> <p>User Community: DOE, industry and academic researchers as a user community seeking capabilities for rapid innovation in materials.</p>	
Goals	Speed the discovery of advanced materials through informatically driven simulation surveys.	
Use case description	Innovation of battery technologies through massive simulations spanning wide spaces of possible design. Systematic computational studies of innovation possibilities in photovoltaics. Rational design of materials based on search and simulation.	
Current solutions	Compute(System)	Hopper.nersc.gov (150K cores), omics-like data analytics hardware resources.
	Storage	GPFS, MongoDB
	Networking	10 GB
	Software	PyMatGen, FireWorks, VASP, ABINIT, NWChem, BerkeleyGW, varied community codes
Big data characteristics	Data source (distributed/centralized)	Gateway-like. Data streams from simulation surveys driven on centralized peta/exascale systems. Widely distributed web of dataflows from central gateway to users.
	Volume (size)	100 TB (current), 500 TB within 5 years. Scalable key-value and object store databases needed.
	Velocity (e.g. real time)	High throughput computing (HTC), fine-grained tasking and queuing. Rapid start/stop for ensembles of tasks. Real-time data analysis for web-like responsiveness.
	Variety (multiple datasets, mashup)	Mashup of simulation outputs across codes and levels of theory. Formatting, registration and integration of datasets. Mashups of data across simulation scales.

	Variability (rate of change)	The targets for materials design will become more search and crowd-driven. The computational backend must flexibly adapt to new targets.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Validation and UQ of simulation with experimental data of varied quality. Error checking and bounds estimation from simulation inter-comparison.
	Visualization	Materials browsers as data from search grows. Visual design of materials.
	Data quality (syntax)	UQ in results based on multiple datasets. Propagation of error in knowledge systems.
	Data types	Key value pairs, JSON, materials file formats
	Data analytics	Map/Reduce and search that join simulation and experimental data.
	Big data specific challenges (Gaps)	HTC at scale for simulation science. Flexible data methods at scale for messy data. Machine learning and knowledge systems that integrate data from publications, experiments, and simulations to advance goal-driven thinking in materials design.
Big data specific challenges in mobility	Potential exists for widespread delivery of actionable knowledge in materials science. Many materials genomics “apps” are amenable to a mobile platform.	
Security and privacy technical considerations	Ability to “sandbox” or create independent working areas between data stakeholders. Policy-driven federation of datasets.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	An OSTP blueprint toward broader materials genomics goals was made available in May 2013.	
More information (URLs)	— The Materials Project. http://www.materialsproject.org . Accessed March 3, 2015.	

A.3 Defense

A.3.1 Use case 13: Large Scale Geospatial Analysis and Visualization

Use case title	Large Scale Geospatial Analysis and Visualization	
Vertical (area)	Defense – but applicable to many others	
Author/company/email	David Boyd/Data Tactics/ dboyd@data-tactics.com	
Actors/stakeholders and their roles and responsibilities	Geospatial Analysts Decision Makers Policy Makers	
Goals	Support large scale geospatial data analysis and visualization.	
Use case description	As the number of geospatially aware sensors increase and the number of geospatially tagged data sources increases the volume geospatial data requiring complex analysis and visualization is growing exponentially. Traditional GIS systems are generally capable of analyzing a millions of objects and easily visualizing thousands. Today's intelligence systems often contain trillions of geospatial objects and need to be able to visualize and interact with millions of objects.	
Current solutions	Compute (System)	Compute and Storage systems - Laptops to Large servers (see notes about clusters) Visualization systems - handhelds to laptops
	Storage	Compute and Storage - local disk or SAN Visualization - local disk, flash ram
	Networking	Compute and Storage - Gigabit or better LAN connection Visualization - Gigabit wired connections, Wireless including WiFi (802.11), Cellular (3g/4g), or Radio Relay
	Software	Compute and Storage – generally Linux or Win Server with Geospatially enabled RDBMS, Geospatial server/analysis software – ESRI ArcServer, Geoserver Visualization – Windows, Android, IOS – browser based visualization. Some laptops may have local ArcMap.
Big data characteristics	Data source (distributed/centralized)	Very distributed.
	Volume (size)	Imagery – 100s of Terabytes Vector Data – 10s of GBs but billions of points

	Velocity (e.g. real time)	Some sensors delivery vector data in NRT. Visualization of changes should be NRT.
	Variety (multiple datasets, mashup)	Imagery (various formats NITF, GeoTiff, CADRG) Vector (various formats shape files, kml, text streams: Object types include points, lines, areas, polylines, circles, ellipses.
	Variability (rate of change)	Moderate to high
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Data accuracy is critical and is controlled generally by three factors: 1) Sensor accuracy is a big issue. 2) datum/spheroid. 3) Image registration accuracy
	Visualization	Displaying in a meaningful way large data sets (millions of points) on small devices (handhelds) at the end of low bandwidth networks.
	Data quality (syntax)	The typical problem is visualization implying quality/accuracy not available in the original data. All data should include metadata for accuracy or circular error probability.
	Data types	Imagery (various formats NITF, GeoTiff, CADRG) Vector (various formats shape files, kml, text streams: Object types include points, lines, areas, polylines, circles, ellipses.
	Data analytics	Closest point of approach, deviation from route, point density over time, PCA and ICA
	Big data specific challenges (Gaps)	Indexing, retrieval and distributed analysis Visualization generation and transmission
Big data specific challenges in mobility	Visualization of data at the end of low bandwidth wireless connections.	
Security and privacy technical considerations	Data is sensitive and must be completely secure in transit and at rest (particularly on handhelds)	

Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Geospatial data requires unique approaches to indexing and distributed analysis.
More information (URLs)	Applicable Standards: http://www.opengeospatial.org/standards http://geojson.org/ http://earth-info.nga.mil/publications/specs/printed/CADRG/cadrg.html Geospatial Indexing: Quad Trees, Space Filling Curves (Hilbert Curves) – You can google these for lots of references.
NOTE There has been some work with in DoD related to this problem set. Specifically, the DCGS-A standard cloud stores, indexes, and analyzes some big data sources. However, many issues still remain with visualization.	

A.3.2 Use case 14: Object Identification and Tracking – Persistent Surveillance

Use case title	Object identification and tracking from Wide Area Large Format Imagery (WALF) Imagery or Full Motion Video (FMV) – Persistent Surveillance	
Vertical (area)	Defense (Intelligence)	
Author/company/email	David Boyd/Data Tactics/dboyd@data-tactics.com	
Actors/stakeholders and their roles and responsibilities	<ol style="list-style-type: none"> 1) Civilian Military decision makers 2) Intelligence Analysts 3) Warfighters 	
Goals	To be able to process and extract/track entities (vehicles, people, packages) over time from the raw image data. Specifically, the idea is to reduce the petabytes of data generated by persistent surveillance down to a manageable size (e.g. vector tracks)	
Use case description	Persistent surveillance sensors can easily collect petabytes of imagery data in the space of a few hours. It is unfeasible for this data to be processed by humans for either alerting or tracking purposes. The data needs to be processed close to the sensor which is likely forward deployed since it is too large to be easily transmitted. The data should be reduced to a set of geospatial object (points, tracks, etc.) which can easily be integrated with other data to form a common operational picture.	
Current solutions	Compute(System)	Various – they range from simple storage capabilities mounted on the sensor, to simple display and storage, to limited object extraction. Typical object extraction systems are currently small (1-20 node) GPU enhanced clusters.
	Storage	Currently flat files persisted on disk in most cases. Sometimes RDBMS indexes pointing to files or portions of files based on metadata/telemetry data.
	Networking	Sensor comms tend to be Line of Sight or Satellite based.
	Software	A wide range custom software and tools including traditional RDBMS and display tools.
Big data characteristics	Data source (distributed/centralized)	Sensors include airframe mounted and fixed position optical, IR, and SAR images.

	Volume (size)	FMV – 30 to 60 frames per/sec at full color 1080P resolution. WALF – 1 to 10 frames per/sec at 10K×10K full color resolution.
	Velocity (e.g. real time)	Real Time
	Variety (multiple datasets, mashup)	Data Typically exists in one or more standard imagery or video formats.
	Variability (rate of change)	Little
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	The veracity of extracted objects is critical. If the system fails or generates false positives people are put at risk.
	Visualization	Visualization of extracted outputs will typically be as overlays on a geospatial display. Overlay objects should be links back to the originating image/video segment.
	Data quality (syntax)	Data quality is generally driven by a combination of sensor characteristics and weather (both obscuring factors - dust/moisture and stability factors - wind).
	Data types	Standard imagery and video formats are input. Output should be in the form of OGC compliant web features or standard geospatial files (shape files, KML).
	Data analytics	<ol style="list-style-type: none"> 1) Object identification (type, size, color) and tracking. 2) Pattern analysis of object (did the truck observed every Weds. afternoon take a different route today or is there a standard route this person takes every day). 3) Crowd behavior/dynamics (is there a small group attempting to incite a riot. Is this person out of place in the crowd or behaving differently?) 4) Economic activity <ol style="list-style-type: none"> a) is the line at the bread store, the butcher, or the ice cream store, b) are more trucks traveling north with goods than trucks going south c) Has activity at or the size of stores in this market place increased or decreased over the past year. 5) Fusion of data with other data to improve quality and confidence.
Big data specific challenges (Gaps)	Processing the volume of data in NRT to support alerting and situational awareness.	
Big data specific challenges in mobility	Getting data from mobile sensor to processing	

Security and privacy technical considerations	Significant – sources and methods cannot be compromised the enemy should not be able to know what we see.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Typically this type of processing fits well into massively parallel computing such as provided by GPUs. Typical problem is integration of this processing into a larger cluster capable of processing data from several sensors in parallel and in NRT. Transmission of data from sensor to system is also a large challenge.
More information (URLs)	<p>Motion Imagery Standards - http://www.gwg.nga.mil/misb/</p> <p>Some of many papers on object identity/tracking: http://www.dabi.temple.edu/~hbling/publication/SPIE12_Dismount_Formatted_v2_BW.pdf</p> <p>http://csce.uark.edu/~jgauch/library/Tracking/Orten.2005.pdf</p> <p>http://www.sciencedirect.com/science/article/pii/S0031320305004863</p> <p>General Articles on the need: http://www.militaryaerospace.com/topics/m/video/79088650/persistent-surveillance-relies-on-extracting-relevant-data-points-and-connecting-the-dots.htm</p> <p>http://www.defencetalk.com/wide-area-persistent-surveillance-revolutionizes-tactical-isr-45745/</p> <p>http://www.defencetalk.com/wide-area-persistent-surveillance-revolutionizes-tactical-isr-45745/</p>

A.3.3 Use case 15: Intelligence Data Processing and Analysis

Use case title	Intelligence Data Processing and Analysis
Vertical (area)	Defense (Intelligence)
Author/company/email	David Boyd/Data Tactics/dboyd@data-tactics.com
Actors/stakeholders and their roles and responsibilities	<p>Senior Civilian/Military Leadership</p> <p>Field Commanders</p> <p>Intelligence Analysts</p> <p>Warfighters</p>
Goals	<ol style="list-style-type: none"> 1) Provide automated alerts to Analysts, Warfighters, Commanders, and Leadership based on incoming intelligence data. 2) Allow Intelligence Analysts to identify in Intelligence data <ol style="list-style-type: none"> a) Relationships between entities (people, organizations, places, equipment) b) Trends in sentiment or intent for either general population or leadership group (state, non-state actors). c) Location of and possibly timing of hostile actions (including implantation of IEDs). d) Track the location and actions of (potentially) hostile actors 3) Ability to reason against and derive knowledge from diverse, disconnected, and frequently unstructured (e.g. text) data sources. 4) Ability to process data close to the point of collection and allow data to be shared easily to/from individual soldiers, forward deployed units, and senior leadership in garrison.

<p>Use case description</p>	<p>1) Ingest/accept data from a wide range of sensors and sources across intelligence disciplines (IMINT, MASINT, GEOINT, HUMINT, SIGINT, OSINT, etc.)</p> <p>2) Process, transform, or align data from disparate sources in disparate formats into a unified data space to permit:</p> <ul style="list-style-type: none"> a) Search b) Reasoning c) Comparison <p>3) Provide alerts to users of significant changes in the state of monitored entities or significant activity within an area.</p> <p>4) Provide connectivity to the edge for the Warfighter (in this case the edge would go as far as a single soldier on dismounted patrol)</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>Fixed and deployed computing clusters ranging from 1000s of nodes to 10s of nodes.</p>
	<p>Storage</p>	<p>10s of Terabytes to 100s of Petabytes for edge and fixed site clusters. Dismounted soldiers would have at most 1-100s of GBs (mostly single digit handheld data storage sizes).</p>
	<p>Networking</p>	<p>Networking with-in and between in garrison fixed sites is robust. Connectivity to forward edge is limited and often characterized by high latency and packet loss. Remote comms might be Satellite based (high latency) or even limited to RF Line of sight radio.</p>
	<p>Software</p>	<p>Currently baseline leverages:</p> <ul style="list-style-type: none"> 1) Hadoop 2) Accumulo (Big Table) 3) Solr 4) NLP (several variants) 5) Puppet (for deployment and security) 6) Storm 7) Custom applications and visualization tools
<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>Very distributed</p>
	<p>Volume (size)</p>	<p>Some IMINT sensors can produce over a petabyte of data in the space of hours. Other data is as small as infrequent sensor activations or text messages.</p>
	<p>Velocity (e.g. real time)</p>	<p>Much sensor data is real time (Full motion video, SIGINT) other is less real time. The critical aspect is to be able ingest, process, and disseminate alerts in NRT.</p>
	<p>Variety (multiple datasets, mashup)</p>	<p>Everything from text files, raw media, imagery, video, audio, electronic data, human generated data.</p>
	<p>Variability (rate of change)</p>	<p>While sensor interface formats tend to be stable, most other data is uncontrolled and may be in any format. Much of the data is unstructured.</p>

<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Data provenance (e.g. tracking of all transfers and transformations) must be tracked over the life of the data.</p> <p>Determining the veracity of “soft” data sources (generally human generated) is a critical requirement.</p>
	<p>Visualization</p>	<p>Primary visualizations will be Geospatial overlays and network diagrams. Volume amounts might be millions of points on the map and thousands of nodes in the network diagram.</p>
	<p>Data quality (syntax)</p>	<p>Data Quality for sensor generated data is generally known (image quality, sig/noise) and good.</p> <p>Unstructured or “captured” data quality varies significantly and frequently cannot be controlled.</p>
	<p>Data types</p>	<p>Imagery, Video, Text, Digital documents of all types, Audio, Digital signal data.</p>
	<p>Data analytics</p>	<p>1) NRT Alerts based on patterns and baseline changes. 2) Link Analysis 3) Geospatial Analysis 4) Text Analytics (sentiment, entity extraction, etc.)</p>
<p>Big data specific challenges (Gaps)</p>	<p>1) Big (or even moderate size data) over tactical networks 2) Data currently exists in disparate silos which must be accessible through a semantically integrated data space. 3) Most critical data is either unstructured or imagery/video which requires significant processing to extract entities and information.</p>	
<p>Big data specific challenges in mobility</p>	<p>The outputs of this analysis and information must be transmitted to or accessed by the dismounted forward soldier.</p>	
<p>Security and privacy technical considerations</p>	<p>Foremost. Data must be protected against:</p> <p>1) Unauthorized access or disclosure 2) Tampering</p>	
<p>Highlight issues for generalizing this Use case (e.g. for ref. architecture)</p>	<p>Wide variety of data types, sources, structures, and quality which will span domains and requires integrated search and reasoning.</p>	
<p>More information (URLs)</p>	<p>http://www.afcea-aberdeen.org/files/presentations/AFCEAAberdeen_DCGSA_COLWells_PS.pdf http://stids.c4i.gmu.edu/papers/STIDSPapers/STIDS2012_T14_SmithEtAl_HorizontalIntegrationOfWarfighterIntel.pdf http://stids.c4i.gmu.edu/STIDS2011/papers/STIDS2011_CR_T1_SalmenEtAl.pdf http://www.youtube.com/watch?v=l4Qii7T8zeg http://dcgsa.apg.army.mil/</p>	

A.4 Healthcare and Life Sciences

A.4.1 Use case 16: Electronic Medical Record Data

Use case title	Electronic Medical Record (EMR) Data	
Vertical (area)	Healthcare	
Author/company/email	Shaun Grannis/Indiana University/sgrannis@regenstrief.org	
Actors/stakeholders and their roles and responsibilities	<p><u>Biomedical informatics research scientists</u> (implement and evaluate enhanced methods for seamlessly integrating, standardizing, analyzing, and operationalizing highly heterogeneous, high-volume clinical data streams);</p> <p><u>Health services researchers</u> (leverage integrated and standardized EMR data to derive knowledge that supports implementation and evaluation of translational, comparative effectiveness, patient-centered outcomes research);</p> <p><u>Healthcare providers – physicians, nurses, public health officials</u> (leverage information and knowledge derived from integrated and standardized EMR data to support direct patient care and population health)</p>	
Goals	Use advanced methods for normalizing patient, provider, facility and clinical concept identification within and among separate health care organizations to enhance models for defining and extracting clinical phenotypes from non-standard discrete and free-text clinical data using feature selection, information retrieval and machine learning decision-models. Leverage clinical phenotype data to support cohort selection, clinical outcomes research, and clinical decision support.	
Use case description	<p>As health care systems increasingly gather and consume EMR data, large national initiatives aiming to leverage such data are emerging, and include developing a digital learning health care system to support increasingly evidence-based clinical decisions with timely accurate and up-to-date patient-centered clinical information; using electronic observational clinical data to efficiently and rapidly translate scientific discoveries into effective clinical treatments; and electronically sharing integrated health data to improve healthcare process efficiency and outcomes. These key initiatives all rely on high-quality, large-scale, standardized and aggregate health data. Despite the promise that increasingly prevalent and ubiquitous EMR data hold, enhanced methods for integrating and rationalizing these data are needed for a variety of reasons. Data from clinical systems evolve over time. This is because the concept space in healthcare is constantly evolving: new scientific discoveries lead to new disease entities, new diagnostic modalities, and new disease management approaches. These in turn lead to new clinical concepts, which drive the evolution of health concept ontologies. Using heterogeneous data from the Indiana Network for Patient Care (INPC), the nation's largest and longest-running health information exchange, which includes more than 4 billion discrete coded clinical observations from more than 100 hospitals for more than 12 million patients, we will use information retrieval techniques to identify highly relevant clinical features from electronic observational data. We will deploy information retrieval and natural language processing techniques to extract clinical features. Validated features will be used to parameterize clinical phenotype decision models based on maximum likelihood estimators and Bayesian networks. Using these decision models we will identify a variety of clinical phenotypes such as diabetes, congestive heart failure, and pancreatic cancer.</p>	
Current solutions	Compute(System)	Big Red II, a new Cray supercomputer at I.U.
	Storage	Teradata, PostgreSQL, MongoDB

	Networking	Various. Significant I/O intensive processing needed.
	Software	Hadoop, Hive, R. Unix-based.
Big data characteristics	Data source (distributed/centralized)	Clinical data from more than 1,100 discrete logical, operational healthcare sources in the Indiana Network for Patient Care (INPC) the nation's largest and longest-running health information exchange.
	Volume (size)	More than 12 million patients, more than 4 billion discrete clinical observations. > 20 TB raw data.
	Velocity (e.g. real time)	Between 500,000 and 1,5 million new real-time clinical transactions added per day.
	Variety (multiple datasets, mashup)	We integrate a broad variety of clinical datasets from multiple sources: free text provider notes; inpatient, outpatient, laboratory, and emergency department encounters; chromosome and molecular pathology; chemistry studies; cardiology studies; hematology studies; microbiology studies; neurology studies; provider notes; referral labs; serology studies; surgical pathology and cytology, blood bank, and toxicology studies.
	Variability (rate of change)	Data from clinical systems evolve over time because the clinical and biological concept space is constantly evolving: new scientific discoveries lead to new disease entities, new diagnostic modalities, and new disease management approaches. These in turn lead to new clinical concepts, which drive the evolution of health concept ontologies, encoded in highly variable fashion.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Data from each clinical source are commonly gathered using different methods and representations, yielding substantial heterogeneity. This leads to systematic errors and bias requiring robust methods for creating semantic interoperability.

	<p>Visualization</p>	<p>Inbound data volume, accuracy, and completeness must be monitored on a routine basis using focus visualization methods. Intrinsic informational characteristics of data sources must be visualized to identify unexpected trends.</p>
	<p>Data quality (syntax)</p>	<p>A central barrier to leveraging EMR data is the highly variable and unique local names and codes for the same clinical test or measurement performed at different institutions. When integrating many data sources, mapping local terms to a common standardized concept using a combination of probabilistic and heuristic classification methods is necessary.</p>
	<p>Data types</p>	<p>Wide variety of clinical data types including numeric, structured numeric, free-text, structured text, discrete nominal, discrete ordinal, discrete structured, binary large blobs (images and video).</p>
	<p>Data analytics</p>	<p>Information retrieval methods to identify relevant clinical features (tf-idf, latent semantic analysis, mutual information). Natural Language Processing techniques to extract relevant clinical features. Validated features will be used to parameterize clinical phenotype decision models based on maximum likelihood estimators and Bayesian networks. Decision models will be used to identify a variety of clinical phenotypes such as diabetes, congestive heart failure, and pancreatic cancer.</p>
<p>Big data specific challenges (Gaps)</p>	<p>Overcoming the systematic errors and bias in large-scale, heterogeneous clinical data to support decision-making in research, patient care, and administrative use-cases requires complex multistage processing and analytics that demands substantial computing power. Further, the optimal techniques for accurately and effectively deriving knowledge from observational clinical data are nascent.</p>	
<p>Big data specific challenges in mobility</p>	<p>Biological and clinical data are needed in a variety of contexts throughout the healthcare ecosystem. Effectively delivering clinical data and knowledge across the healthcare ecosystem will be facilitated by mobile platform such as mHealth.</p>	

Security and privacy technical considerations	Privacy and confidentiality of individuals must be preserved in compliance with federal and state requirements including HIPAA. Developing analytic models using comprehensive, integrated clinical data requires aggregation and subsequent de-identification prior to applying complex analytics.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Patients increasingly receive health care in a variety of clinical settings. The subsequent EMR data is fragmented and heterogeneous. In order to realize the promise of a Learning Health Care system as advocated by the National Academy of Science and the Institute of Medicine, EMR data must be rationalized and integrated. The methods we propose in this use-case support integrating and rationalizing clinical data to support decision-making at multiple levels.
More information (URLs)	Regenstrief Institute (http://www.regenstrief.org); Logical observation identifiers names and codes (http://www.loinc.org); Indiana Health Information Exchange (http://www.ihie.org); Institute of Medicine Learning Healthcare System (http://www.iom.edu/Activities/Quality/LearningHealthcare.aspx)

A.4.2 Use case 17: Pathology Imaging/Digital Pathology

Use case title	Pathology Imaging/digital pathology	
Vertical (area)	Healthcare	
Author/company/email	Fusheng Wang/Emory University/fusheng.wang@emory.edu	
Actors/stakeholders and their roles and responsibilities	Biomedical researchers on translational research; hospital clinicians on imaging guided diagnosis	
Goals	Develop high performance image analysis algorithms to extract spatial information from images; provide efficient spatial queries and analytics, and feature clustering and classification	
Use case description	Digital pathology imaging is an emerging field where examination of high resolution images of tissue specimens enables novel and more effective ways for disease diagnosis. Pathology image analysis segments massive (millions per image) spatial objects such as nuclei and blood vessels, represented with their boundaries, along with many extracted image features from these objects. The derived information is used for many complex queries and analytics to support biomedical research and clinical diagnosis. Recently, 3D pathology imaging is made possible through 3D laser technologies or serially sectioning hundreds of tissue sections onto slides and scanning them into digital images. Segmenting 3D microanatomic objects from registered serial images could produce tens of millions of 3D objects from a single image. This provides a deep “map” of human tissues for next generation diagnosis.	
Current solutions	Compute(System)	Supercomputers; Cloud
	Storage	SAN or HDFS
	Networking	Need excellent external network link
	Software	MPI for image analysis; Map/Reduce + Hive with spatial extension
Big data characteristics	Data source (distributed/centralized)	Digitized pathology images from human tissues

	Volume (size)	1 GB raw image data + 1,5 GB analytical results per 2D image; 1 TB raw image data + 1 TB analytical results per 3D image. 1 PB data per moderated hospital per year
	Velocity (e.g. real time)	Once generated, data will not be changed
	Variety (multiple datasets, mashup)	Image characteristics and analytics depend on disease types
	Variability (rate of change)	No change
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	High quality results validated with human annotations are essential
	Visualization	Needed for validation and training
	Data quality (syntax)	Depend on pre-processing of tissue slides such as chemical staining and quality of image analysis algorithms
	Data types	Raw images are whole slide images (mostly based on BIG-TIFF), and analytical results are structured data (spatial boundaries and features)
	Data analytics	Image analysis, spatial queries and analytics, feature clustering and classification
Big data specific challenges (Gaps)	Extreme large size; multi-dimensional; disease specific analytics; correlation with other data types (clinical data, -omic data)	
Big data specific challenges in mobility	3D visualization of 3D pathology images is not likely in mobile platforms	
Security and privacy technical considerations	Protected health information has to be protected; public data have to be de-identified	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Imaging data; multi-dimensional spatial data analytics	
More information (URLs)	https://web.cci.emory.edu/confluence/display/PAIS https://web.cci.emory.edu/confluence/display/HadoopGIS	

A.4.3 Use case 18: Computational Bioimaging

Use case title	Computational Bioimaging	
Vertical (area)	Scientific Research: Biological Science	
Author/company/email	David Skinner ¹ , deskinner@lbl.gov Joaquin Correa ¹ , JoaquinCorrea@lbl.gov Daniela Ushizima ² , dushizima@lbl.gov Joerg Meyer ² , joergmeyer@lbl.gov ¹ National Energy Scientific Computing Center (NERSC), Lawrence Berkeley National Laboratory, USA ² Computational Research Division, Lawrence Berkeley National Laboratory, USA	
Actors/stakeholders and their roles and responsibilities	<u>Capability providers:</u> Bioimaging instrument operators, microscope developers, imaging facilities, applied mathematicians, and data stewards. <u>User Community:</u> DOE, industry and academic researchers seeking to collaboratively build models from imaging data.	
Goals	Data delivered from bioimaging is increasingly automated, higher resolution, and multi-modal. This has created a data analysis bottleneck that, if resolved, can advance the biosciences discovery through big data techniques. Our goal is to solve that bottleneck with extreme scale computing. Meeting that goal will require more than computing. It will require building communities around data resources and providing advanced algorithms for massive image analysis. High-performance computational solutions can be harnessed by community-focused science gateways to guide the application of massive data analysis toward massive imaging data sets. Workflow components include data acquisition, storage, enhancement, minimizing noise, segmentation of regions of interest, crowd-based selection and extraction of features, and object classification, and organization, and search.	
Use case description	Web-based one-stop-shop for high performance, high throughput image processing for producers and consumers of models built on bio-imaging data.	
Current solutions	Compute(System)	Hopper.nersc.gov (150K cores)
	Storage	Database and image collections
	Networking	10Gb, could use 100Gb and advanced networking (SDN)
	Software	ImageJ, OMER, VolRover, advanced segmentation and feature detection methods from applied math researchers

Big data characteristics	Data source (distributed/centralized)	Distributed experimental sources of bioimages (instruments). Scheduled high volume flows from automated high-resolution optical and electron microscopes.
	Volume (size)	Growing very fast. Scalable key-value and object store databases needed. In database processing and analytics. 50 TB here now, but currently over a petabyte overall. A single scan on emerging machines is 32 TB
	Velocity (e.g. real time)	High throughput computing (HTC), responsive analysis
	Variety (multiple datasets, mashup)	Multi-modal imaging essentially must mash-up disparate channels of data with attention to registration and dataset formats.
	Variability (rate of change)	Biological samples are highly variable and their analysis workflows must cope with wide variation.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Data is messy overall as is training classifiers.
	Visualization	Heavy use of 3D structural models.
	Data quality (syntax)	
	Data types	Imaging file formats
	Data analytics	Machine learning (SVM and RF) for classification and recommendation services.
Big data specific challenges (Gaps)	HTC at scale for simulation science. Flexible data methods at scale for messy data. Machine learning and knowledge systems that drive pixel based data toward biological objects and models.	
Big data specific challenges in mobility		
Security and privacy technical considerations		
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	There is potential in generalizing concepts of search in the context of bioimaging.	
More information (URLs)		

A.4.4 Use case 19: Genomic Measurements

Use case title	Genomic Measurements
Vertical (area)	Healthcare
Author/company/email	Justin Zook/NIST/jzook@nist.gov
Actors/stakeholders and their roles and responsibilities	NIST/Genome in a Bottle Consortium – public/private/academic partnership

Goals	Develop well-characterized Reference Materials, Reference Data, and Reference Methods needed to assess performance of genome sequencing
Use case description	Integrate data from multiple sequencing technologies and methods to develop highly confident characterization of whole human genomes as Reference Materials, and develop methods to use these Reference Materials to assess performance of any genome sequencing run
Current solutions	Compute(System) 72-core cluster for our NIST group, collaboration with >1000 core clusters at FDA, some groups are using cloud
	Storage ~40 TB NFS at NIST, PBs of genomics data at NIH/NCBI
	Networking Varies. Significant I/O intensive processing needed
	Software Open-source sequencing bioinformatics software from academic groups (UNIX-based)
Big data characteristics	Data source (distributed/centralized) Sequencers are distributed across many laboratories, though some core facilities exist.
	Volume (size) 40 TB NFS is full, will need >100 TB in 1 to 2 years at NIST; Healthcare community will need many PBs of storage
	Velocity (e.g. real time) DNA sequencers can generate ~300 GB compressed data/day. Velocity has increased much faster than Moore's Law
	Variety (multiple datasets, mashup) File formats not well-standardized, though some standards exist. Generally structured data.
	Variability (rate of change) Sequencing technologies have evolved very rapidly, and new technologies are on the horizon.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics) All sequencing technologies have significant systematic errors and biases, which require complex analysis methods and combining multiple technologies to understand, often with machine learning

	Visualization	“Genome browsers” have been developed to visualize processed data
	Data quality (syntax)	Sequencing technologies and bioinformatics methods have significant systematic errors and biases
	Data types	Mainly structured text
	Data analytics	Processing of raw data to produce variant calls. Also, clinical interpretation of variants, which is now very challenging.
Big data specific challenges (Gaps)	Processing data requires significant computing power, which poses challenges especially to clinical laboratories as they are starting to perform large-scale sequencing. Long-term storage of clinical sequencing data could be expensive. Analysis methods are quickly evolving. Many parts of the genome are challenging to analyze, and systematic errors are difficult to characterize.	
Big data specific challenges in mobility	Physicians may need access to genomic data on mobile platforms	
Security and privacy technical considerations	Sequencing data in health records or clinical research databases must be kept secure/private, though our Consortium data is public.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	I have some generalizations to medical genome sequencing above, but focus on NIST/Genome in a Bottle Consortium work. Currently, labs doing sequencing range from small to very large. Future data could include other ‘omics’ measurements, which could be even larger than DNA sequencing	
More information (URLs)	Genome in a Bottle Consortium: http://www.genomeinabottle.org	

A.4.5 Use case 20: Comparative Analysis for (meta) Genomes

Use case title	Comparative analysis for metagenomes and genomes
Vertical (area)	Scientific Research: Genomics
Author/company/email	Ernest Szeto / LBNL / eszeto@lbl.gov
Actors/stakeholders and their roles and responsibilities	Joint Genome Institute (JGI) Integrated Microbial Genomes (IMG) project. Heads: Victor M. Markowitz, and Nikos C. Kypides. User community: JGI, bioinformaticians and biologists worldwide.
Goals	Provide an integrated comparative analysis system for metagenomes and genomes. This includes interactive Web UI with core data, backend precomputations, batch job computation submission from the UI.

<p>Use case description</p>	<p>Given a metagenomic sample, (1) determine the community composition in terms of other reference isolate genomes, (2) characterize the function of its genes, (3) begin to infer possible functional pathways, (4) characterize similarity or dissimilarity with other metagenomic samples, (5) begin to characterize changes in community composition and function due to changes in environmental pressures, (6) isolate sub-sections of data based on quality measures and community composition.</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>Linux cluster, Oracle RDBMS server, large memory machines, standard Linux interactive hosts</p>
	<p>Storage</p>	<p>Oracle RDBMS, SQLite files, flat text files, Lucy (a version of Lucene) for keyword searches, BLAST databases, USEARCH databases</p>
	<p>Networking</p>	<p>Provided by NERSC</p>
	<p>Software</p>	<p>Standard bioinformatics tools (BLAST, HMMER, multiple alignment and phylogenetic tools, gene callers, sequence feature predictors...), Perl/Python wrapper scripts, Linux Cluster scheduling</p>
<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>Centralized.</p>
	<p>Volume (size)</p>	<p>50 TB</p>
	<p>Velocity (e.g. real time)</p>	<p>Front end web UI must be real time interactive. Back end data loading processing must keep up with exponential growth of sequence data due to the rapid drop in cost of sequencing technology.</p>

	<p>Variety (multiple datasets, mashup)</p>	<p>Biological data is inherently heterogeneous, complex, structural, and hierarchical. One begins with sequences, followed by features on sequences, such as genes, motifs, regulatory regions, followed by organization of genes in neighborhoods (operons), to proteins and their structural features, to coordination and expression of genes in pathways. Besides core genomic data, new types of “Omics” data such as transcriptomics, methylomics, and proteomics describing gene expression under a variety of conditions must be incorporated into the comparative analysis system.</p>
	<p>Variability (rate of change)</p>	<p>The sizes of metagenomic samples can vary by several orders of magnitude, such as several hundred thousand genes to a billion genes (e.g., latter in a complex soil sample).</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Metagenomic sampling science is currently preliminary and exploratory. Procedures for evaluating assembly of highly fragmented data in raw reads are better defined, but still an open research area.</p>

	<p>Visualization</p>	<p>Interactive speed of web UI on very large data sets is an ongoing challenge. Web UI's still seem to be the preferred interface for most biologists. It is use for basic querying and browsing of data. More specialized tools may be launched from them, e.g. for viewing multiple alignments. Ability to download large amounts of data for offline analysis is another requirement of the system.</p>
	<p>Data quality (syntax)</p>	<p>Improving quality of metagenomic assembly is still a fundamental challenge. Improving the quality of reference isolate genomes, both in terms of the coverage in the phylogenetic tree, improved gene calling and functional annotation is a more mature process, but an ongoing project.</p>
	<p>Data types</p>	<p>Cf. above on "Variety"</p>
	<p>Data analytics</p>	<p>Descriptive statistics, statistical significance in hypothesis testing, discovering new relationships, data clustering and classification is a standard part of the analytics. The less quantitative part includes the ability to visualize structural details at different levels of resolution. Data reduction, removing redundancies through clustering, more abstract representations such as representing a group of highly similar genomes in a pangenome are all strategies for both data management as well as analytics.</p>

Big data specific challenges (Gaps)	The biggest friend for dealing with the heterogeneity of biological data is still the RDBMS. Unfortunately, it does not scale for the current volume of data. NoSQL solutions aim at providing an alternative. Unfortunately, NoSQL solutions do not always lend themselves to real time interactive use, rapid and parallel bulk loading, and sometimes have issues regarding robustness. Our current approach is currently ad hoc, custom, relying mainly on the Linux cluster and the file system to supplement the Oracle RDBMS. The custom solution oftentimes rely in knowledge of the peculiarities of the data allowing us to devise horizontal partitioning schemes as well as inversion of data organization when applicable.
Big data specific challenges in mobility	No special challenges. Just world wide web access.
Security and privacy technical considerations	No special challenges. Data is either public or requires standard login with password.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	A replacement for the RDBMS in big data would be of benefit to everyone. Many NoSQL solutions attempt to fill this role, but have their limitations.
More information (URLs)	http://img.jgi.doe.gov

A.4.6 Use case 21: Individualized Diabetes Management

Use case title	Individualized Diabetes Management
Vertical (area)	Healthcare
Author/company/email	Peter Li, Ying Ding, Philip Yu, Geoffrey Fox, David Wild at Mayo Clinic, Indiana University, UIC; dingying@indiana.edu
Actors/stakeholders and their roles and responsibilities	Mayo Clinic + IU/semantic integration of EHR data UIC/semantic graph mining of EHR data IU cloud and parallel computing
Goals	Develop advanced graph-based data mining techniques applied to EHR to search for these cohorts and extract their EHR data for outcome evaluation. These methods will push the boundaries of scalability and data mining technologies and advance knowledge and practice in these areas as well as clinical management of complex diseases.

<p>Use case description</p>	<p>Diabetes is a growing illness in world population, affecting both developing and developed countries. Current management strategies do not adequately take into account of individual patient profiles, such as co-morbidities and medications, which are common in patients with chronic illnesses. We propose to approach this shortcoming by identifying similar patients from a large Electronic Health Record (EHR) database, i.e. an individualized cohort, and evaluate their respective management outcomes to formulate one best solution suited for a given patient with diabetes.</p> <p>Project under development as below</p> <p>Stage 1: Use the Semantic Linking for Property Values method to convert an existing data warehouse at Mayo Clinic, called the Enterprise Data Trust (EDT), into RDF triples that enables us to find similar patients much more efficiently through linking of both vocabulary-based and continuous values.</p> <p>Stage 2: Needs efficient parallel retrieval algorithms, suitable for cloud or HPC, using open source Hbase with both indexed and custom search to identify patients of possible interest.</p> <p>Stage 3: The EHR, as an RDF graph, provides a very rich environment for graph pattern mining. Needs new distributed graph mining algorithms to perform pattern analysis and graph indexing technique for pattern searching on RDF triple graphs.</p> <p>Stage 4: Given the size and complexity of graphs, mining subgraph patterns could generate numerous false positives and miss numerous false negatives. Needs robust statistical analysis tools to manage false discovery rate and determine true subgraph significance and validate these through several clinical use cases.</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>supercomputers; cloud</p>
	<p>Storage</p>	<p>HDFS</p>
	<p>Networking</p>	<p>Varies. Significant I/O intensive processing needed</p>
	<p>Software</p>	<p>Mayo internal data warehouse called Enterprise Data Trust (EDT)</p>
<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>distributed EHR data</p>
	<p>Volume (size)</p>	<p>The Mayo Clinic EHR dataset is a very large dataset containing over 5 million patients with thousands of properties each and many more that are derived from primary values.</p>

	<p>Velocity (e.g. real time)</p>	<p>not real time but updated periodically</p>
	<p>Variety (multiple datasets, mashup)</p>	<p>Structured data, a patient has controlled vocabulary (CV) property values (demographics, diagnostic codes, medications, procedures, etc.) and continuous property values (lab tests, medication amounts, vitals, etc.). The number of property values could range from less than 100 (new patient) to more than 100,000 (long term patient) with typical patients composed of 100 CV values and 1000 continuous values. Most values are time based, i.e. a timestamp is recorded with the value at the time of observation.</p>
	<p>Variability (rate of change)</p>	<p>Data will be updated or added during each patient visit.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Data are annotated based on domain ontologies or taxonomies. Semantics of data can vary from labs to labs.</p>
	<p>Visualization</p>	<p>no visualization</p>
	<p>Data quality (syntax)</p>	<p>Provenance is important to trace the origins of the data and data quality</p>
	<p>Data types</p>	<p>text, and Continuous Numerical values</p>
	<p>Data analytics</p>	<p>Integrating data into semantic graph, using graph traverse to replace SQL join. Developing semantic graph mining algorithms to identify graph patterns, index graph, and search graph. Indexed Hbase. Custom code to develop new patient properties from stored data.</p>
<p>Big data specific challenges (Gaps)</p>	<p>For individualized cohort, we will effectively be building a datamart for each patient since the critical properties and indices will be specific to each patient. Due to the number of patients, this becomes an impractical approach. Fundamentally, the paradigm changes from relational row-column lookup to semantic graph traversal.</p>	

Big data specific challenges in mobility	Physicians and patient may need access to this data on mobile platforms
Security and privacy technical considerations	Health records or clinical research databases must be kept secure/private.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Data integration: continuous values, ontological annotation, taxonomy Graph Search: indexing and searching graph Validation: Statistical validation
More information (URLs)	

A.4.7 Use case 22: Statistical Relational AI for Health Care

Use case title	Statistical Relational AI for Health Care	
Vertical (area)	Healthcare	
Author/company/email	Sriraam Natarajan/Indiana University /natarasr@indiana.edu	
Actors/stakeholders and their roles and responsibilities	Researchers in Informatics, medicine and practitioners in medicine.	
Goals	The goal of the project is to analyze large, multi-modal, longitudinal data. Analyzing different data types such as imaging, EHR, genetic and natural language data requires a rich representation. This approach employs the relational probabilistic models that have the capability of handling rich relational data and modeling uncertainty using probability theory. The software learns models from multiple data types and can possibly integrate the information and reason about complex queries.	
Use case description	Users can provide a set of descriptions – say for instance, MRI images and demographic data about a particular subject. They can then query for the onset of a particular disease (say Alzheimer's) and the system will then provide a probability distribution over the possible occurrence of this disease.	
Current solutions	Compute(System)	A high performance computer (48 GB RAM) is needed to run the code for a few hundred patients. Clusters for large datasets
	Storage	A 200 GB to 1 TB hard drive typically stores the test data. The relevant data is retrieved to main memory to run the algorithms. Backend data in database or NoSQL stores
	Networking	Intranet.
	Software	Mainly Java based, in house tools are used to process the data.

Big data characteristics	Data source (distributed/centralized)	All the data about the users reside in a single disk file. Sometimes, resources such as published text need to be pulled from internet.
	Volume (size)	Variable due to the different amount of data collected. Typically can be in 100s of GBs for a single cohort of a few hundred people. When dealing with millions of patients, this can be in the order of 1 petabyte.
	Velocity (e.g. real time)	Varied. In some cases, EHRs are constantly being updated. In other controlled studies, the data often comes in batches in regular intervals.
	Variety (multiple datasets, mashup)	This is the key property in medical data sets. That data is typically in multiple tables and need to be merged in order to perform the analysis.
	Variability (rate of change)	The arrival of data is unpredictable in many cases as they arrive in real time.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Challenging due to different modalities of the data, human errors in data collection and validation.
	Visualization	The visualization of the entire input data is nearly impossible. But typically, partially visualizable. The models built can be visualized under some reasonable assumptions.
	Data quality (syntax)	
	Data types	EHRs, imaging, genetic data that are stored in multiple databases.
	Data analytics	

Big data specific challenges (Gaps)	Data is in abundance in many cases of medicine. The key issue is that there can possibly be too much data (as images, genetic sequences etc.) that can make the analysis complicated. The real challenge lies in aligning the data and merging from multiple sources in a form that can be made useful for a combined analysis. The other issue is that sometimes, large amount of data is available about a single subject but the number of subjects themselves is not very high (i.e., data imbalance). This can result in learning algorithms picking up random correlations between the multiple data types as important features in analysis. Hence, robust learning methods that can faithfully model the data are of paramount importance. Another aspect of data imbalance is the occurrence of positive examples (i.e., cases). The incidence of certain diseases may be rare making the ratio of cases to controls extremely skewed making it possible for the learning algorithms to model noise instead of examples.
Big data specific challenges in mobility	
Security and privacy technical considerations	Secure handling and processing of data is of crucial importance in medical domains.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Models learned from one set of populations cannot be easily generalized across other populations with diverse characteristics. This requires that the learned models can be generalized and refined according to the change in the population characteristics.
More information (URLs)	

A.4.8 Use case 23: World Population Scale Epidemiology

Use case title	World Population Scale Epidemiological Study	
Vertical (area)	Epidemiology, Simulation Social Science, Computational Social Science	
Author/company/email	Madhav Marathe Stephen Eubank or Chris Barrett/ Virginia Bioinformatics Institute, Virginia Tech, mmarathe@vbi.vt.edu, seubank@vbi.vt.edu or cbarrett@vbi.vt.edu	
Actors/stakeholders and their roles and responsibilities	Government and non-profit institutions involved in health, public policy, and disaster mitigation. Social Scientist who wants to study the interplay between behavior and contagion.	
Goals	(a) Build a synthetic global population. (b) Run simulations over the global population to reason about outbreaks and various intervention strategies.	
Use case description	Prediction and control of pandemic similar to the 2009 H1N1 influenza.	
Current solutions	Compute(System)	Distributed (MPI) based simulation system written in Charm++. Parallelism is achieved by exploiting the disease residence time period.

	Storage	Network file system. Exploring database driven techniques.
	Networking	Infiniband. High bandwidth 3D Torus.
	Software	Charm++, MPI
Big data characteristics	Data source (distributed/centralized)	Generated from synthetic population generator. Currently centralized. However, could be made distributed as part of post-processing.
	Volume (size)	100 TB
	Velocity (e.g. real time)	Interactions with experts and visualization routines generate large amount of real time data. Data feeding into the simulation is small but data generated by simulation is massive.
	Variety (multiple datasets, mashup)	Variety depends upon the complexity of the model over which the simulation is being performed. Can be very complex if other aspects of the world population such as type of activity, geographical, socio-economic, cultural variations are taken into account.
	Variability (rate of change)	Depends upon the evolution of the model and corresponding changes in the code. This is complex and time intensive. Hence low rate of change.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Robustness of the simulation is dependent upon the quality of the model. However, robustness of the computation itself, although non-trivial, is tractable.
	Visualization	Would require very large amount of movement of data to enable visualization.
	Data quality (syntax)	Consistent due to generation from a model
	Data types	Primarily network data.
	Data analytics	Summary of various runs and replicates of a simulation

Big data specific challenges (Gaps)	Computation of the simulation is both compute intensive and data intensive. Moreover, due to unstructured and irregular nature of graph processing the problem is not easily decomposable. Therefore it is also bandwidth intensive. Hence, a supercomputer is applicable than cloud type clusters.
Big data specific challenges in mobility	None
Security and privacy technical considerations	Several issues at the synthetic population-modeling phase (see social contagion model).
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	In general contagion diffusion of various kinds: information, diseases, social unrest can be modeled and computed. All of them are agent-based model that utilize the underlying interaction network to study the evolution of the desired phenomena.
More information (URLs)	

A.4.9 Use case 24: Social Contagion Modeling

Use case title	Social Contagion Modeling	
Vertical (area)	Social behavior (including national security, public health, viral marketing, city planning, disaster preparedness)	
Author/company/email	Madhav Marathe or Chris Kuhlman/ Virginia Bioinformatics Institute, Virginia Tech Techymarathe@vbi.vt.edu or ckuhlman@vbi.vt.edu	
Actors/stakeholders and their roles and responsibilities		
Goals	<p>Provide a computing infrastructure that models social contagion processes.</p> <p>The infrastructure enables different types of human-to-human interactions (e.g., face-to-face versus online media; mother-daughter relationships versus mother-coworker relationships) to be simulated. It takes not only human-to-human interactions into account, but also interactions among people, services (e.g., transportation), and infrastructure (e.g., internet, electric power).</p>	
Use case description	<p>Social unrest. People take to the streets to voice unhappiness with government leadership. There are citizens that both support and oppose government. Quantify the degrees to which normal business and activities are disrupted owing to fear and anger. Quantify the possibility of peaceful demonstrations, violent protests. Quantify the potential for government responses ranging from appeasement, to allowing protests, to issuing threats against protestors, to actions to thwart protests. To address these issues, must have fine-resolution models and datasets.</p>	
Current solutions	Compute(System)	Distributed processing software running on commodity clusters and newer architectures and systems (e.g., clouds).

	Storage	File servers (including archives), databases.
	Networking	Ethernet, Infiniband, and similar.
	Software	Specialized simulators, open source software, and proprietary modeling environments. Databases.
Big data characteristics	Data source (distributed/centralized)	Many data sources: populations, work locations, travel patterns, utilities (e.g., power grid) and other man-made infrastructures, online (social) media.
	Volume (size)	Easily 10s of TB per year of new data.
	Velocity (e.g. real time)	During social unrest events, human interactions and mobility key to understanding system dynamics. Rapid changes in data; e.g., who follows whom in Twitter.
	Variety (multiple datasets, mashup)	Variety of data seen in wide range of data sources. Temporal data. Data fusion. Data fusion a big issue. How to combine data from different sources and how to deal with missing or incomplete data? Multiple simultaneous contagion processes.
	Variability (rate of change)	Because of stochastic nature of events, multiple instances of models and inputs must be run to ranges in outcomes.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Failover of soft real-time analyses.
	Visualization	Large datasets; time evolution; multiple contagion processes over multiple network representations. Levels of detail (e.g., individual, neighborhood, city, state, country-level).
	Data quality (syntax)	Checks for ensuring data consistency, corruption. Preprocessing of raw data for use in models.

	Data types	Wide-ranging data, from human characteristics to utilities and transportation systems, and interactions among them.
	Data analytics	Models of behavior of humans and hard infrastructures, and their interactions. Visualization of results.
Big data specific challenges (Gaps)	How to take into account heterogeneous features of 100s of millions or billions of individuals, models of cultural variations across countries that are assigned to individual agents? How to validate these large models? Different types of models (e.g. multiple contagions): disease, emotions, behaviors. Modeling of different urban infrastructure systems in which humans act. With multiple replicates required to assess stochasticity, large amounts of output data are produced; storage requirements.	
Big data specific challenges in mobility	How and where to perform these computations? Combinations of cloud computing and clusters. How to realize most efficient computations; move data to compute resources?	
Security and privacy technical considerations	Two dimensions. First, privacy and anonymity issues for individuals used in modeling (e.g., Twitter and Facebook users). Second, securing data and computing platforms for computation.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Fusion of different data types. Different datasets must be combined depending on the particular problem. How to quickly develop, verify, and validate new models for new applications. What is appropriate level of granularity to capture phenomena of interest while generating results sufficiently quickly; i.e., how to achieve a scalable solution. Data visualization and extraction at different levels of granularity.	
More information (URLs)		

A.4.10 Use case 25: LifeWatch Biodiversity

Use case title	LifeWatch – E-Science European Infrastructure for Biodiversity and Ecosystem Research
Vertical (area)	Scientific Research: Life Science
Author/company/email	Wouter Los, Yuri Demchenko (y.demchenko@uva.nl), University of Amsterdam
Actors/stakeholders and their roles and responsibilities	End-users (biologists, ecologists, field researchers) Data analysts, data archive managers, e-Science Infrastructure managers, EU states national representatives
Goals	Research and monitor different ecosystems, biological species, their dynamics and migration.

<p>Use case description</p>	<p>LifeWatch project and initiative intends to provide integrated access to a variety of data, analytical and modeling tools as served by a variety of collaborating initiatives. Another service is offered with data and tools in selected workflows for specific scientific communities. In addition, LifeWatch will provide opportunities to construct personalized 'virtual labs', also allowing to enter new data and analytical tools.</p> <p>New data will be shared with the data facilities cooperating with LifeWatch.</p> <p>Particular case studies: Monitoring alien species, monitoring migrating birds, wetlands</p> <p>LifeWatch operates Global Biodiversity Information facility and Biodiversity Catalogue that is Biodiversity Science Web Services Catalogue</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>Field facilities TBD</p> <p>Data center: General Grid and cloud based resources provided by national e-Science centers</p>
	<p>Storage</p>	<p>Distributed, historical and trends data archiving</p>
	<p>Networking</p>	<p>May require special dedicated or overlay sensor network.</p>
	<p>Software</p>	<p>Web Services based, Grid based services, relational databases</p>
<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>Ecological information from numerous observation and monitoring facilities and sensor network, satellite images/information, climate and weather, all recorded information.</p> <p>Information from field researchers</p>
	<p>Volume (size)</p>	<p>Involves many existing data sets/sources</p> <p>Collected amount of data TBD</p>
	<p>Velocity (e.g. real time)</p>	<p>Data analyzed incrementally, processes dynamics corresponds to dynamics of biological and ecological processes.</p> <p>However may require real-time processing and analysis in case of the natural or industrial disaster.</p> <p>May require data streaming processing.</p>

	<p>Variety (multiple datasets, mashup)</p>	<p>Variety and number of involved databases and observation data is currently limited by available tools; in principle, unlimited with the growing ability to process data for identifying ecological changes, factors/reasons, species evolution and trends.</p> <p>See below in additional information.</p>
	<p>Variability (rate of change)</p>	<p>Structure of the datasets and models may change depending on the data processing stage and tasks</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>In normal monitoring mode are data are statistically processed to achieve robustness.</p> <p>Some biodiversity research is critical to data veracity (reliability / trustworthiness).</p> <p>In case of natural and technogenic disasters data veracity is critical.</p>
	<p>Visualization</p>	<p>Requires advanced and rich visualization, high definition visualization facilities, visualization data</p> <ul style="list-style-type: none"> — 4D visualization — Visualizing effects of parameter change in (computational) models — Comparing model outcomes with actual observations (multi dimensional)
	<p>Data quality (syntax)</p>	<p>Depends on and ensued by initial observation data.</p> <p>Quality of analytical data depends on used mode and algorithms that are constantly improved.</p> <p>Repeating data analytics should be possible to re-evaluate initial observation data.</p> <p>Actionable data are human aided.</p>

	<p>Data types Multi-type Relational data, key-value, complex semantically rich data</p> <p>Data analytics Parallel data streams and streaming analytics</p>
Big data specific challenges (Gaps)	<p>Variety, multi-type data: SQL and no-SQL, distributed multi-source data.</p> <p>Visualization, distributed sensor networks.</p> <p>Data storage and archiving, data exchange and integration; data linkage: from the initial observation data to processed data and reported/visualized data.</p> <ul style="list-style-type: none"> — Historical unique data — Curated (authorized) reference data (i.e. species names lists), algorithms, software code, workflows — Processed (secondary) data serving as input for other researchers — Provenance (and persistent identification [PID]) control of data, algorithms, and workflows
Big data specific challenges in mobility	<p>Require supporting mobile sensors (e.g. birds migration) and mobile researchers (both for information feed and catalogue search)</p> <ul style="list-style-type: none"> — Instrumented field vehicles, Ships, Planes, Submarines, floating buoys, sensor tagging on organisms — Photos, video, sound recording
Security and privacy technical considerations	<p>Data integrity, referral integrity of the datasets.</p> <p>Federated identity management for mobile researchers and mobile sensors</p> <p>Confidentiality, access control and accounting for information on protected species, ecological information, space images, climate information.</p>
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	<ul style="list-style-type: none"> — Support of distributed sensor network — Multi-type data combination and linkage; potentially unlimited data variety — Data life cycle management: data provenance, referral integrity and identification — Access and integration of multiple distributed databases
More information (URLs)	<p>http://www.lifewatch.eu/web/guest/home</p> <p>https://www.biodiversitycatalogue.org/</p>

NOTE

Variety of data used in Biodiversity research

Genetic (genomic) diversity

- DNA sequences and barcodes
- Metabolomics functions

Species information

- species names
- occurrence data (in time and place)
- species traits and life history data
- host-parasite relations
- collection specimen data

Copyrighted document, no reproduction or circulation
STANDARDSISO.COM: Click to view the full PDF of ISO/IEC TR 20547 WG - 2:2018
For review by FG on 21/10/2024
Oct 2024

Ecological information

- biomass, trunk/root diameter and other physical characteristics
- population density etc.
- habitat structures
- C/N/P etc. molecular cycles

Ecosystem data

- species composition and community dynamics
- remote and earth observation data
- CO₂ fluxes
- Soil characteristics
- Algal blooming
- Marine temperature, salinity, pH, currents, etc.

Ecosystem services

- productivity (i.e., biomass production/time)
- fresh water dynamics
- erosion
- climate buffering
- genetic pools

Data concepts

- conceptual framework of each data
- ontologies
- provenance data

Algorithms and workflows

- software code and provenance
- tested workflows

Multiple sources of data and information

- Specimen collection data
- Observations (human interpretations)
- Sensors and sensor networks (terrestrial, marine, soil organisms), bird etc. tagging
- Aerial and satellite observation spectra
- Field * Laboratory experimentation
- Radar and LiDAR
- Fisheries and agricultural data
- Deceases and epidemics

A.5 Deep Learning and Social media

A.5.1 Use case 26: Large-scale Deep Learning

Use case title	Large-scale Deep Learning	
Vertical (area)	Machine Learning/AI	
Author/company/email	Adam Coates / Stanford University / acoates@cs.stanford.edu	
Actors/stakeholders and their roles and responsibilities	Machine learning researchers and practitioners faced with large quantities of data and complex prediction tasks. Supports state-of-the-art development in computer vision as in automatic car driving, speech recognition, and natural language processing in both academic and industry systems.	
Goals	Increase the size of datasets and models that can be tackled with deep learning algorithms. Large models (e.g., neural networks with more neurons and connections) combined with large datasets are increasingly the top performers in benchmark tasks for vision, speech, and NLP.	
Use case description	A research scientist or machine learning practitioner wants to train a deep neural network from a large (>>1TB) corpus of data (typically imagery, video, audio, or text). Such training procedures often require customization of the neural network architecture, learning criteria, and dataset pre-processing. In addition to the computational expense demanded by the learning algorithms, the need for rapid prototyping and ease of development is extremely high.	
Current solutions	Compute(System)	GPU cluster with high-speed interconnects (e.g., Infiniband, 40gE)
	Storage	100 TB Lustre filesystem
	Networking	Infiniband within HPC cluster; 1G ethernet to outside infrastructure (e.g., Web, Lustre).
	Software	In-house GPU kernels and MPI-based communication developed by Stanford CS. C++/Python source.
Big data characteristics	Data source (distributed/centralized)	Centralized filesystem with a single large training dataset. Dataset may be updated with new training examples as they become available.
	Volume (size)	Current datasets typically 1 to 10 TB. With increases in computation that enable much larger models, datasets of 100 TB or more may be necessary in order to exploit the representational power of the larger models. Training a self-driving car could take 100 million images.
	Velocity (e.g. real time)	Much faster than real-time processing is required. Current computer vision applications involve processing hundreds of image frames per second in order to ensure reasonable training times. For demanding applications (e.g., autonomous driving) we envision the need to process many thousand high-resolution (6 megapixels or more) images per second.
	Variety (multiple datasets, mashup)	Individual applications may involve a wide variety of data. Current research involves neural networks that actively learn from heterogeneous tasks (e.g., learning to perform tagging, chunking and parsing for text, or learning to read lips from combinations of video and audio).

	Variability (rate of change)	Low variability. Most data is streamed in at a consistent pace from a shared source. Due to high computational requirements, server loads can introduce burstiness into data transfers.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Datasets for ML applications are often hand-labeled and verified. Extremely large datasets involve crowd-sourced labeling and invite ambiguous situations where a label is not clear. Automated labeling systems still require human sanity-checks. Clever techniques for large dataset construction is an active area of research.
	Visualization	Visualization of learned networks is an open area of research, though partly as a debugging technique. Some visual applications involve visualization predictions on test imagery.
	Data quality (syntax)	Some collected data (e.g., compressed video or audio) may involve unknown formats, codecs, or may be corrupted. Automatic filtering of original source data removes these.
	Data types	Images, video, audio, text. (In practice: almost anything.)
	Data analytics	Small degree of batch statistical pre-processing; all other data analysis is performed by the learning algorithm itself.
Big data specific challenges (Gaps)	Processing requirements for even modest quantities of data are extreme. Though the trained representations can make use of many terabytes of data, the primary challenge is in processing all of the data during training. Current state-of-the-art deep learning systems are capable of using neural networks with more than 10 billion free parameters (akin to synapses in the brain), and necessitate trillions of floating point operations per training example. Distributing these computations over high-performance infrastructure is a major challenge for which we currently use a largely custom software system.	
Big data specific challenges in mobility	After training of large neural networks is completed, the learned network may be copied to other devices with dramatically lower computational capabilities for use in making predictions in real time. (E.g., in autonomous driving, the training procedure is performed using a HPC cluster with 64 GPUs. The result of training, however, is a neural network that encodes the necessary knowledge for making decisions about steering and obstacle avoidance. This network can be copied to embedded hardware in vehicles or sensors.)	

Security and privacy technical considerations	None.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	<p>Deep Learning shares many characteristics with the broader field of machine learning. The paramount requirements are high computational throughput for mostly dense linear algebra operations, and extremely high productivity. Most deep learning systems require a substantial degree of tuning on the target application for best performance and thus necessitate a large number of experiments with designer intervention in between. As a result, minimizing the turn-around time of experiments and accelerating development is crucial.</p> <p>These two requirements (high throughput and high productivity) are dramatically in contention. HPC systems are available to accelerate experiments, but current HPC software infrastructure is difficult to use which lengthens development and debugging time and, in many cases, makes otherwise computationally tractable applications infeasible.</p> <p>The major components needed for these applications (which are currently in-house custom software) involve dense linear algebra on distributed-memory HPC systems. While libraries for single-machine or single-GPU computation are available (e.g., BLAS, CuBLAS, MAGMA, etc.), distributed computation of dense BLAS-like or LAPACK-like operations on GPUs remains poorly developed. Existing solutions (e.g., ScaLapack for CPUs) are not well-integrated with higher level languages and require low-level programming which lengthens experiment and development time.</p>
More information (URLs)	<p>Recent popular press coverage of deep learning technology:</p> <p>http://www.nytimes.com/2012/11/24/science/scientists-see-advances-in-deep-learning-a-part-of-artificial-intelligence.html http://www.nytimes.com/2012/06/26/technology/in-a-big-network-of-computers-evidence-of-machine-learning.html http://www.nytimes.com/2012/06/26/technology/in-a-big-network-of-computers-evidence-of-machine-learning.html http://www.wired.com/wiredenterprise/2013/06/andrew_ng/ http://www.wired.com/wiredenterprise/2013/06/andrew_ng/</p> <p>A recent research paper on HPC for Deep Learning: http://www.stanford.edu/~acoates/papers/CoatesHuvalWangWuNgCatanzaro_icml2013.pdf</p> <p>Widely-used tutorials and references for Deep Learning:</p> <p>http://ufldl.stanford.edu/wiki/index.php/Main_Page http://deeplearning.net/</p>

A.5.2 Use case 27: Large Scale Consumer Photos Organization

Use case title	Organizing large-scale, unstructured collections of consumer photos
Vertical (area)	(Scientific Research: Artificial Intelligence)
Author/company/email	David Crandall, Indiana University, djcran@indiana.edu
Actors/stakeholders and their roles and responsibilities	Computer vision researchers (to push forward state of art), media and social network companies (to help organize large-scale photo collections), consumers (browsing both personal and public photo collections), researchers and others interested in producing cheap 3D models (archaeologists, architects, urban planners, interior designers...)

	<p>Goals Produce 3D reconstructions of scenes using collections of millions to billions of consumer images, where neither the scene structure nor the camera positions are known a priori. Use resulting 3D models to allow efficient and effective browsing of large-scale photo collections by geographic position. Geolocate new images by matching to 3D models. Perform object recognition on each image.</p>								
<p>Use case description</p>	<p>3D reconstruction is typically posed as a robust non-linear least squares optimization problem in which observed (noisy) correspondences between images are constraints and unknowns are 6-d camera pose of each image and 3D position of each point in the scene. Sparsity and large degree of noise in constraints typically makes naïve techniques fall into local minima that are not close to actual scene structure. Typical specific steps are: (1) extracting features from images, (2) matching images to find pairs with common scene structures, (3) estimating an initial solution that is close to scene structure and/or camera parameters, (4) optimizing non-linear objective function directly. Of these, (1) is embarrassingly parallel, (2) is an all-pairs matching problem, usually with heuristics to reject unlikely matches early on. We solve (3) using discrete optimization using probabilistic inference on a graph (Markov Random Field) followed by robust Levenberg-Marquardt in continuous space. Others solve (3) by solving (4) for a small number of images and then incrementally adding new images, using output of last round as initialization for next round. (4) is typically solved with Bundle Adjustment, which is a non-linear least squares solver that is optimized for the particular constraint structure that occurs in 3D reconstruction problems. Image recognition problems are typically embarrassingly parallel, although learning object models involves learning a classifier (e.g. a Support Vector Machine), a process that is often hard to parallelize.</p>								
<p>Current solutions</p>	<table border="1"> <tr> <td data-bbox="715 1406 1054 1480">Compute(System)</td> <td data-bbox="1054 1406 1391 1480">Hadoop cluster (about 60 nodes, 480 core)</td> </tr> <tr> <td data-bbox="715 1480 1054 1525">Storage</td> <td data-bbox="1054 1480 1391 1525">Hadoop DFS and flat files</td> </tr> <tr> <td data-bbox="715 1525 1054 1570">Networking</td> <td data-bbox="1054 1525 1391 1570">Simple Unix</td> </tr> <tr> <td data-bbox="715 1570 1054 1704">Software</td> <td data-bbox="1054 1570 1391 1704">Hadoop Map-reduce, simple hand-written multithreaded tools (ssh and sockets for communication)</td> </tr> </table>	Compute(System)	Hadoop cluster (about 60 nodes, 480 core)	Storage	Hadoop DFS and flat files	Networking	Simple Unix	Software	Hadoop Map-reduce, simple hand-written multithreaded tools (ssh and sockets for communication)
Compute(System)	Hadoop cluster (about 60 nodes, 480 core)								
Storage	Hadoop DFS and flat files								
Networking	Simple Unix								
Software	Hadoop Map-reduce, simple hand-written multithreaded tools (ssh and sockets for communication)								
<p>Big data characteristics</p>	<table border="1"> <tr> <td data-bbox="715 1704 1054 1809">Data source (distributed/centralized)</td> <td data-bbox="1054 1704 1391 1809">Publicly-available photo collections, e.g. on Flickr, Panoramio, etc.</td> </tr> <tr> <td data-bbox="715 1809 1054 1915">Volume (size)</td> <td data-bbox="1054 1809 1391 1915">500+ billion photos on Facebook, 5+ billion photos on Flickr.</td> </tr> <tr> <td data-bbox="715 1915 1054 1984">Velocity (e.g. real time)</td> <td data-bbox="1054 1915 1391 1984">100+ million new photos added to Facebook per day.</td> </tr> </table>	Data source (distributed/centralized)	Publicly-available photo collections, e.g. on Flickr, Panoramio, etc.	Volume (size)	500+ billion photos on Facebook, 5+ billion photos on Flickr.	Velocity (e.g. real time)	100+ million new photos added to Facebook per day.		
Data source (distributed/centralized)	Publicly-available photo collections, e.g. on Flickr, Panoramio, etc.								
Volume (size)	500+ billion photos on Facebook, 5+ billion photos on Flickr.								
Velocity (e.g. real time)	100+ million new photos added to Facebook per day.								

	Variety (multiple datasets, mashup)	Images and metadata including EXIF tags (focal distance, camera type, etc.),
	Variability (rate of change)	Rate of photos varies significantly, e.g. roughly 10x photos to Facebook on New Years versus other days. Geographic distribution of photos follows long-tailed distribution, with 1000 landmarks (totaling only about 100 square km) accounting for over 20 % of photos on Flickr.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Important to make as accurate as possible, subject to limitations of computer vision technology.
	Visualization	Visualize large-scale 3D reconstructions, and navigate large-scale collections of images that have been aligned to maps.
	Data quality (syntax)	Features observed in images are quite noisy due both to imperfect feature extraction and to non-ideal properties of specific images (lens distortions, sensor noise, image effects added by user, etc.)
	Data types	Images, metadata
	Data analytics	
Big data specific challenges (Gaps)	Analytics needs continued monitoring and improvement.	
Big data specific challenges in mobility	Many/most images are captured by mobile devices; eventual goal is to push reconstruction and organization to phone to allow real-time interaction with the user.	
Security and privacy technical considerations	Need to preserve privacy for users and digital rights for media.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Components of this use case including feature extraction, feature matching, and large-scale probabilistic inference appear in many or most computer vision and image processing problems, including recognition, stereo resolution, image denoising, etc.	
More information (URLs)	http://vision.soic.indiana.edu/disco	

A.5.3 Use case 28: Truthy Twitter Data Analysis

Use case title	Truthy: Information diffusion research from Twitter Data	
Vertical (area)	Scientific Research: Complex Networks and Systems research	
Author/company/email	Filippo Menczer, Indiana University, fil@indiana.edu; Alessandro Flammini, Indiana University, aflammin@indiana.edu; Emilio Ferrara, Indiana University, ferrara@indiana.edu;	
Actors/stakeholders and their roles and responsibilities	Research funded by NFS, DARPA, and McDonnell Foundation.	
Goals	Understanding how communication spreads on socio-technical networks. Detecting potentially harmful information spread at the early stage (e.g., deceiving messages, orchestrated campaigns, untrustworthy information, etc.)	
Use case description	(1) Acquisition and storage of a large volume of continuous streaming data from Twitter (~100 million messages per day, ~500 GB data/day increasing over time); (2) near real-time analysis of such data, for anomaly detection, stream clustering, signal classification and online-learning; (3) data retrieval, big data visualization, data-interactive Web interfaces, public API for data querying.	
Current solutions	Compute(System)	Current: in-house cluster hosted by Indiana University. Critical requirement: large cluster for data storage, manipulation, querying and analysis.
	Storage	Current: Raw data stored in large compressed flat files, since August 2010. Need to move towards Hadoop/IndexedHBase and HDFS distributed storage. Redis as an in-memory database as a buffer for real-time analysis.
	Networking	10 GB/Infiniband required.
	Software	Hadoop, Hive, Redis for data management. Python/SciPy/NumPy/MPI for data analysis.
Big data characteristics	Data source (distributed/centralized)	Distributed - with replication/redundancy
	Volume (size)	~30 TB/year compressed data
	Velocity (e.g. real time)	Near real-time data storage, querying and analysis

	<p>Variety (multiple datasets, mashup)</p>	<p>Data schema provided by social media data source. Currently using Twitter only. We plan to expand incorporating Google+, Facebook</p>
	<p>Variability (rate of change)</p>	<p>Continuous real-time data stream incoming from each source.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>99,99 % uptime required for real-time data acquisition. Service outages might corrupt data integrity and significance.</p>
	<p>Visualization</p>	<p>Information diffusion, clustering, and dynamic network visualization capabilities already exist.</p>
	<p>Data quality (syntax)</p>	<p>Data structured in standardized formats, the overall quality is extremely high. We generate aggregated statistics; expand the features set, etc., generating high-quality derived data.</p>
	<p>Data types</p>	<p>Fully-structured data (JSON format) enriched with users meta-data, geo-locations, etc.</p>
	<p>Data analytics</p>	<p>Stream clustering: data are aggregated according to topics, meta-data and additional features, using ad hoc online clustering algorithms. Classification: using multi-dimensional time series to generate, network features, users, geographical, content features, etc., we classify information produced on the platform. Anomaly detection: real-time identification of anomalous events (e.g., induced by exogenous factors). Online learning: applying machine learning/deep learning methods to real-time information diffusion patterns analysis, users profiling, etc.</p>

Big data specific challenges (Gaps)	Dealing with real-time analysis of large volume of data. Providing a scalable infrastructure to allocate resources, storage space, etc. on-demand if required by increasing data volume over time.
Big data specific challenges in mobility	Implementing low-level data storage infrastructure features to guarantee efficient, mobile access to data.
Security and privacy technical considerations	Twitter publicly releases data collected by our platform. Although, data-sources incorporate user meta-data (in general, not sufficient to uniquely identify individuals), therefore some policy for data storage security and privacy protection must be implemented.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Definition of high-level data schema to incorporate multiple data-sources providing similarly structured data.
More information (URLs)	http://truthy.indiana.edu/ http://cnets.indiana.edu/groups/nan/truthy http://cnets.indiana.edu/groups/nan/despic

A.5.4 Use case 29: Crowd Sourcing in the Humanities

Use case title	Crowd Sourcing in the Humanities as Source for Big and Dynamic Data	
Vertical (area)	Humanities, Social Sciences	
Author/company/email	Sebastian Drude <Sebastian.Drude@mpi.nl>, Max Planck Institute for Psycholinguistics	
Actors/stakeholders and their roles and responsibilities	Scientists (Sociologists, Psychologists, Linguists, Politic Scientists, Historians, etc.), data managers and analysts, data archives The general public as data providers and participants	
Goals	Capture information (manually entered, recorded multimedia, reaction times, pictures, sensor information) from many individuals and their devices. Thus capture wide ranging individual, social, cultural and linguistic variation among several dimensions (space, social space, time).	
Use case description	Many different possible use cases: get recordings of language usage (words, sentences, meaning descriptions, etc.), answers to surveys, info on cultural facts, transcriptions of pictures and texts — correlate these with other phenomena, detect new cultural practices, behavior, values and believes, discover individual variation	
Current solutions	Compute(System)	Individual systems for manual data collection (mostly Websites)
	Storage	Traditional servers
	Networking	barely used other than for data entry via web

	Software	XML technology, traditional relational databases for storing pictures, not much multi-media yet.
Big data characteristics	Data source (distributed/centralized)	Distributed, individual contributors via webpages and mobile devices
	Volume (size)	Depends dramatically, from hundreds to millions of data records. Depending on data-type: from GBs (text, surveys, experiment values) to hundreds of terabytes (multi-media)
	Velocity (e.g. real time)	Depends very much on project: dozens to thousands of new data records per day Data has to be analyzed incrementally.
	Variety (multiple datasets, mashup)	so far mostly homogeneous small data sets; expected large distributed heterogeneous datasets which have to be archived as primary data
	Variability (rate of change)	Data structure and content of collections are changing during data life cycle. There is no critical variation of data producing speed, or runtime characteristics variations.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Noisy data is possible, unreliable metadata, identification and pre-selection of appropriate data
	Visualization	important for interpretation, no special visualization techniques
	Data quality (syntax)	validation is necessary; quality of recordings, quality of content, spam

	<p>Data types individual data records (survey answers, reaction times); text (e.g., comments, transcriptions,...); multi-media (pictures, audio, video)</p> <p>Data analytics pattern recognition of all kind (e.g., speech recognition, automatic A&V analysis, cultural patterns), identification of structures (lexical units, linguistic rules, etc.)</p>
Big data specific challenges (Gaps)	Data management (metadata, provenance info, data identification with PIDs) Data curation Digitizing existing audio-video, photo and documents archives
Big data specific challenges in mobility	Include data from sensors of mobile devices (position, etc.); Data collection from expeditions and field research.
Security and privacy technical considerations	Privacy issues may be involved (A/V from individuals), anonymization may be necessary but not always possible (A/V analysis, small speech communities) Archive and metadata integrity, long term preservation
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Many individual data entries from many individuals, constant flux of data entry, metadata assignment, etc. Offline vs. online use, to be synchronized later with central database. Giving significant feedback to contributors.
More information (URLs)	—
<p>NOTE Crowd sourcing has been barely started to be used on a larger scale. With the availability of mobile devices, now there is a huge potential for collecting much data from many individuals, also making use of sensors in mobile devices. This has not been explored on a large scale so far; existing projects of crowd sourcing are usually of a limited scale and web-based.</p>	

A.5.5 Use case 30: CINET Network Science Cyberinfrastructure

Use case title	CINET: Cyberinfrastructure for Network (Graph) Science and Analytics
Vertical (area)	Network Science
Author/company/email	Team lead by Virginia Tech and comprising of researchers from Indiana University, University at Albany, North Carolina AT, Jackson State University, University at Houston Downtown, Argonne National Laboratory Point of Contact: Madhav Marathe or Keith Bisset, Network Dynamics and Simulation Science Laboratory, Virginia Bio-informatics Institute Virginia Tech, mmarathe@vbi.vt.edu / kbisset@vbi.vt.edu
Actors/stakeholders and their roles and responsibilities	Researchers, practitioners, educators and students interested in the study of networks.

	<p>Goals CINET cyberinfrastructure middleware to support network science. This middleware will give researchers, practitioners, teachers and students access to a computational and analytic environment for research, education and training. The user interface provides lists of available networks and network analysis modules (implemented algorithms for network analysis). A user, who can be a researcher in network science area, can select one or more networks and analysis them with the available network analysis tools and modules. A user can also generate random networks following various random graph models. Teachers and students can use CINET for classroom use to demonstrate various graph theoretic properties and behaviors of various algorithms. A user is also able to add a network or network analysis module to the system. This feature of CINET allows it to grow easily and remain up-to-date with the latest algorithms.</p> <p>The goal is to provide a common web-based platform for accessing various (i) network and graph analysis tools such as SNAP, NetworkX, Galib, etc. (ii) real-world and synthetic networks, (iii) computing resources and (iv) data management systems to the end-user in a seamless manner.</p>	
	<p>Use case description Users can run one or more structural or dynamic analysis on a set of selected networks. The domain specific language allows users to develop flexible high level workflows to define more complex network analysis.</p>	
	<p>Current solutions</p> <p>Compute(System)</p>	<p>A high performance computing cluster (DELL C6100), named Shadowfax, of 60 compute nodes and 12 processors (Intel Xeon X5670 2,93GHz) per compute node with a total of 720 processors and 4 GB main memory per processor.</p> <p>Shared memory systems; EC2 based clouds are also used</p> <p>Some of the codes and networks can utilize single node systems and thus are being currently mapped to Open Science Grid</p>
	<p>Storage</p>	<p>628 TB GPFS</p>
	<p>Networking</p>	<p>Internet, infiniband. A loose collection of supercomputing resources.</p>

	Software	Graph libraries: Galib, NetworkX. Distributed Workflow Management: Simfrastructure, databases, semantic web tools
Big data characteristics	Data source (distributed/centralized)	A single network remains in a single disk file accessible by multiple processors. However, during the execution of a parallel algorithm, the network can be partitioned and the partitions are loaded in the main memory of multiple processors.
	Volume (size)	Can be hundreds of GB for a single network.
	Velocity (e.g. real time)	Two types of changes: (i) the networks are very dynamic and (ii) as the repository grows, we expect at least a rapid growth to lead to over 1 000 to 5 000 networks and methods in about a year
	Variety (multiple datasets, mashup)	Data sets are varied: (i) directed as well as undirected networks, (ii) static and dynamic networks, (iii) labeled, (iv) can have dynamics over these networks,
	Variability (rate of change)	The rate of graph-based data is growing at increasing rate. Moreover, increasingly other life sciences domains are using graph-based techniques to address problems. Hence, we expect the data and the computation to grow at a significant pace.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Challenging due to a synchronous distributed computation. Current systems are designed for real-time synchronous response.
	Visualization	As the input graph size grows the visualization system on client side is stressed heavily both in terms of data and compute.
	Data quality (syntax)	

	Data types
	Data analytics
Big data specific challenges (Gaps)	<p>Parallel algorithms are necessary to analyze massive networks. Unlike many structured data, network data is difficult to partition. The main difficulty in partitioning a network is that different algorithms require different partitioning schemes for efficient operation. Moreover, most of the network measures are global in nature and require either i) huge duplicate data in the partitions or ii) very large communication overhead resulted from the required movement of data. These issues become significant challenges for big networks.</p> <p>Computing dynamics over networks is harder since the network structure often interacts with the dynamical process being studied.</p> <p>CINET enables large class of operations across wide variety, both in terms of structure and size, of graphs. Unlike other compute + data intensive systems, such as parallel databases or CFD, performance on graph computation is sensitive to underlying architecture. Hence, a unique challenge in CINET is manage the mapping between workload (graph type + operation) to a machine whose architecture and runtime is conducive to the system.</p> <p>Data manipulation and bookkeeping of the derived for users is another big challenge since unlike enterprise data there is no well-defined and effective models and tools for management of various graph data in a unified fashion.</p>
Big data specific challenges in mobility	
Security and privacy technical considerations	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	HPC as a service. As data volume grows increasingly large number of applications such as biological sciences need to use HPC systems. CINET can be used to deliver the compute resource necessary for such domains.
More information (URLs)	http://cinet.vbi.vt.edu/cinet_new/

A.5.6 Use case 31: NIST Analytic Technology Measurement and Evaluations

Use case title	NIST Information Access Division analytic technology performance measurement, evaluations, and standards
Vertical (area)	Analytic technology performance measurement and standards for government, industry, and academic stakeholders
Author/company/email	John Garofolo (john.garofolo@nist.gov)
Actors/stakeholders and their roles and responsibilities	NIST developers of measurement methods, data contributors, analytic algorithm developers, users of analytic technologies for unstructured, semi-structured data, and heterogeneous data across all sectors.

	<p>Goals Accelerate the development of advanced analytic technologies for unstructured, semi-structured, and heterogeneous data through performance measurement and standards. Focus communities of interest on analytic technology challenges of importance, create consensus-driven measurement metrics and methods for performance evaluation, evaluate the performance of the performance metrics and methods via community-wide evaluations which foster knowledge exchange and accelerate progress, and build consensus towards widely-accepted standards for performance measurement.</p>						
<p>Use case description</p>	<p>Develop performance metrics, measurement methods, and community evaluations to ground and accelerate the development of advanced analytic technologies in the areas of speech and language processing, video and multimedia processing, biometric image processing, and heterogeneous data processing as well as the interaction of analytics with users. Typically employ one of two processing models: 1) Push test data out to test participants and analyze the output of participant systems, 2) Push algorithm test harness interfaces out to participants and bring in their algorithms and test them on internal computing clusters. Developing approaches to support scalable Cloud-based developmental testing. Also perform usability and utility testing on systems with users in the loop.</p>						
<p>Current solutions</p>	<table border="1"> <tr> <td data-bbox="670 1048 1008 1223"> <p>Compute(System)</p> </td> <td data-bbox="1008 1048 1391 1223"> <p>Linux and OS-10 clusters; distributed computing with stakeholder collaborations; specialized image processing architectures.</p> </td> </tr> <tr> <td data-bbox="670 1223 1008 1397"> <p>Storage</p> </td> <td data-bbox="1008 1223 1391 1397"> <p>RAID arrays, and distribute data on 1 to 2 TB drives, and occasionally FTP. Distributed data distribution with stakeholder collaborations.</p> </td> </tr> <tr> <td data-bbox="670 1397 1008 1617"> <p>Networking</p> </td> <td data-bbox="1008 1397 1391 1617"> <p>Fiber channel disk storage, Gigabit Ethernet for system-system communication, general intra- and Internet resources within NIST and shared networking resources with its stakeholders.</p> </td> </tr> </table>	<p>Compute(System)</p>	<p>Linux and OS-10 clusters; distributed computing with stakeholder collaborations; specialized image processing architectures.</p>	<p>Storage</p>	<p>RAID arrays, and distribute data on 1 to 2 TB drives, and occasionally FTP. Distributed data distribution with stakeholder collaborations.</p>	<p>Networking</p>	<p>Fiber channel disk storage, Gigabit Ethernet for system-system communication, general intra- and Internet resources within NIST and shared networking resources with its stakeholders.</p>
<p>Compute(System)</p>	<p>Linux and OS-10 clusters; distributed computing with stakeholder collaborations; specialized image processing architectures.</p>						
<p>Storage</p>	<p>RAID arrays, and distribute data on 1 to 2 TB drives, and occasionally FTP. Distributed data distribution with stakeholder collaborations.</p>						
<p>Networking</p>	<p>Fiber channel disk storage, Gigabit Ethernet for system-system communication, general intra- and Internet resources within NIST and shared networking resources with its stakeholders.</p>						
	<table border="1"> <tr> <td data-bbox="670 1626 1008 1796"> <p>Software</p> </td> <td data-bbox="1008 1626 1391 1796"> <p>PERL, Python, C/C++, Matlab, R development tools. Create ground-up test and measurement applications.</p> </td> </tr> </table>	<p>Software</p>	<p>PERL, Python, C/C++, Matlab, R development tools. Create ground-up test and measurement applications.</p>				
<p>Software</p>	<p>PERL, Python, C/C++, Matlab, R development tools. Create ground-up test and measurement applications.</p>						
<p>Big data characteristics</p>	<table border="1"> <tr> <td data-bbox="670 1800 1008 2076"> <p>Data source (distributed/centralized)</p> </td> <td data-bbox="1008 1800 1391 2076"> <p>Large annotated corpora of unstructured/semi-structured text, audio, video, images, multimedia, and heterogeneous collections of the above including ground truth annotations for training, developmental testing, and summative evaluations.</p> </td> </tr> </table>	<p>Data source (distributed/centralized)</p>	<p>Large annotated corpora of unstructured/semi-structured text, audio, video, images, multimedia, and heterogeneous collections of the above including ground truth annotations for training, developmental testing, and summative evaluations.</p>				
<p>Data source (distributed/centralized)</p>	<p>Large annotated corpora of unstructured/semi-structured text, audio, video, images, multimedia, and heterogeneous collections of the above including ground truth annotations for training, developmental testing, and summative evaluations.</p>						

	<p>Volume (size)</p>	<p>The test corpora exceed 900 M Web pages occupying 30 TB of storage, 100 M tweets, 100 M ground-truthed biometric images, several hundred thousand partially ground-truthed video clips, and terabytes of smaller fully ground-truthed test collections. Even larger data collections are being planned for future evaluations of analytics involving multiple data streams and very heterogeneous data.</p>
	<p>Velocity (e.g. real time)</p>	<p>Most legacy evaluations are focused on retrospective analytics. Newer evaluations are focusing on simulations of real-time analytic challenges from multiple data streams.</p>
	<p>Variety (multiple datasets, mashup)</p>	<p>The test collections span a wide variety of analytic application types including textual search/extraction, machine translation, speech recognition, image and voice biometrics, object and person recognition and tracking, document analysis, human-computer dialogue, and multimedia search/extraction. Future test collections will include mixed type data and applications.</p>
	<p>Variability (rate of change)</p>	<p>Evaluation of tradeoffs between accuracy and data rates as well as variable numbers of data streams and variable stream quality.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>The creation and measurement of the uncertainty associated with the ground-truthing process – especially when humans are involved – is challenging. The manual ground-truthing processes that have been used in the past are not scalable. Performance measurement of complex analytics must include measurement of intrinsic uncertainty as well as ground truthing error to be useful.</p>

	<p>Visualization Visualization of analytic technology performance results and diagnostics including significance and various forms of uncertainty. Evaluation of analytic presentation methods to users for usability, utility, efficiency, and accuracy.</p>
	<p>Data quality (syntax) The performance of analytic technologies is highly impacted by the quality of the data they are employed against with regard to a variety of domain- and application-specific variables. Quantifying these variables is a challenging research task in itself. Mixed sources of data and performance measurement of analytic flows pose even greater challenges with regard to data quality.</p>
	<p>Data types Unstructured and semi-structured text, still images, video, audio, multimedia (audio+video).</p>
	<p>Data analytics Information extraction, filtering, search, and summarization; image and voice biometrics; speech recognition and understanding; machine translation; video person/object detection and tracking; event detection; imagery/document matching; novelty detection; a variety of structural/semantic/temporal analytics and many subtypes of the above.</p>
<p>Big data specific challenges (Gaps)</p>	<p>Scaling ground-truthing to larger data, intrinsic and annotation uncertainty measurement, performance measurement for incompletely annotated data, measuring analytic performance for heterogeneous data and analytic flows involving users.</p>
<p>Big data specific challenges in mobility</p>	<p>Moving training, development, and test data to evaluation participants or moving evaluation participants' analytic algorithms to computational testbeds for performance assessment. Providing developmental tools and data. Supporting agile developmental testing approaches.</p>

<p>Security and privacy technical considerations</p>	<p>Analytic algorithms working with written language, speech, human imagery, etc. must generally be tested against real or realistic data. It's extremely challenging to engineer artificial data that sufficiently captures the variability of real data involving humans. Engineered data may provide artificial challenges that may be directly or indirectly modeled by analytic algorithms and result in overstated performance. The advancement of analytic technologies themselves is increasing privacy sensitivities. Future performance testing methods will need to isolate analytic technology algorithms from the data the algorithms are tested against. Advanced architectures are needed to support security requirements for protecting sensitive data while enabling meaningful developmental performance evaluation. Shared evaluation testbeds must protect the intellectual property of analytic algorithm developers.</p>
<p>Highlight issues for generalizing this Use case (e.g. for ref. architecture)</p>	<p>Scalability of analytic technology performance testing methods, source data creation, and ground truthing; approaches and architectures supporting developmental testing; protecting intellectual property of analytic algorithms and PII and other personal information in test data; measurement of uncertainty using partially-annotated data; composing test data with regard to qualities impacting performance and estimating test set difficulty; evaluating complex analytic flows involving multiple analytics, data types, and user interactions; multiple heterogeneous data streams and massive numbers of streams; mixtures of structured, semi-structured, and unstructured data sources; agile scalable developmental testing approaches and mechanisms.</p>
<p>More information (URLs)</p>	<p>http://www.nist.gov/itl/iad/</p>

A.6 The Ecosystem for Research

A.6.1 Use case 32: DataNet Federation Consortium (DFC)

Use case title	DataNet Federation Consortium (DFC)	
Vertical (area)	Collaboration Environments	
Author/company/email	Reagan Moore / University of North Carolina at Chapel Hill / rwmoore@renci.org	
Actors/stakeholders and their roles and responsibilities	National Science Foundation research projects: Ocean Observatories Initiative (sensor archiving); Temporal Dynamics of Learning Center (Cognitive science data grid); the iPlant Collaborative (plant genomics); Drexel engineering digital library; Odum Institute for social science research (data grid federation with Dataverse).	
Goals	Provide national infrastructure (collaboration environments) that enables researchers to collaborate through shared collections and shared workflows. Provide policy-based data management systems that enable the formation of collections, data grid, digital libraries, archives, and processing pipelines. Provide interoperability mechanisms that federate existing data repositories, information catalogs, and web services with collaboration environments.	
Use case description	Promote collaborative and interdisciplinary research through federation of data management systems across federal repositories, national academic research initiatives, institutional repositories, and international collaborations. The collaboration environment runs at scale: petabytes of data, hundreds of millions of files, hundreds of millions of metadata attributes, tens of thousands of users, and a thousand storage resources.	
Current solutions	Compute(System)	Interoperability with workflow systems (NCSA Cyber-integrator, Kepler, Taverna)
	Storage	Interoperability across file systems, tape archives, cloud storage, object-based storage
	Networking	Interoperability across TCP/IP, parallel TCP/IP, RBUDP, HTTP
	Software	Integrated Rule Oriented Data System (iRODS)
Big data characteristics	Data source (distributed/centralized)	Manage internationally distributed data
	Volume (size)	Petabytes, hundreds of millions of files
	Velocity (e.g. real time)	Support sensor data streams, satellite imagery, simulation output, observational data, experimental data

	Variety (multiple datasets, mashup)	Support logical collections that span administrative domains, data aggregation in containers, metadata, and workflows as objects										
	Variability (rate of change)	Support active collections (mutable data), versioning of data, and persistent identifiers										
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Provide reliable data transfer, audit trails, event tracking, periodic validation of assessment criteria (integrity, authenticity), distributed debugging										
	Visualization	Support execution of external visualization systems through automated workflows (GRASS)										
	Data quality (syntax)	Provide mechanisms to verify quality through automated workflow procedures										
	Data types	Support parsing of selected formats (NetCDF, HDF5, Dicom), and provide mechanisms to invoke other data manipulation methods										
	Data analytics	Provide support for invoking analysis workflows, tracking workflow provenance, sharing of workflows, and re-execution of workflows										
	Big data specific challenges (Gaps)	Provide standard policy sets that enable a new community to build upon data management plans that address federal agency requirements										
Big data specific challenges in mobility	Capture knowledge required for data manipulation, and apply resulting procedures at either the storage location, or a computer server.											
Security and privacy technical considerations	Federate across existing authentication environments through Generic Security Service API and Pluggable Authentication Modules (GSI, Kerberos, InCommon, Shibboleth). Manage access controls on files independently of the storage location.											
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	<p>Currently 25 science and engineering domains have projects that rely on the iRODS policy-based data management system:</p> <table border="0"> <tr> <td>Astrophysics</td> <td>Auger supernova search</td> </tr> <tr> <td>Atmospheric science</td> <td>NASA Langley Atmospheric Sciences Center</td> </tr> <tr> <td>Biology</td> <td>Phylogenetics at CC IN2P3</td> </tr> <tr> <td>Climate</td> <td>NOAA National Climatic Data Center</td> </tr> <tr> <td>Cognitive Science</td> <td>Temporal Dynamics of Learning Center</td> </tr> </table>		Astrophysics	Auger supernova search	Atmospheric science	NASA Langley Atmospheric Sciences Center	Biology	Phylogenetics at CC IN2P3	Climate	NOAA National Climatic Data Center	Cognitive Science	Temporal Dynamics of Learning Center
Astrophysics	Auger supernova search											
Atmospheric science	NASA Langley Atmospheric Sciences Center											
Biology	Phylogenetics at CC IN2P3											
Climate	NOAA National Climatic Data Center											
Cognitive Science	Temporal Dynamics of Learning Center											

	<p>Computer Science GENI experimental network</p> <p>Cosmic Ray AMS experiment on the International Space Station</p> <p>Dark Matter Physics Edelweiss II</p> <p>Earth Science NASA Center for Climate Simulations</p> <p>Ecology CEED Caveat Emptor Ecological Data</p> <p>Engineering CIBER-U</p> <p>High Energy Physics BaBar</p> <p>Hydrology Institute for the Environment, UNC-CH; Hydroshare</p> <p>Genomics Broad Institute, Wellcome Trust Sanger Institute</p> <p>Medicine Sick Kids Hospital</p> <p>Neuroscience International Neuroinformatics Coordinating Facility</p> <p>Neutrino Physics T2K and dChooz neutrino experiments</p> <p>Oceanography Ocean Observatories Initiative</p> <p>Optical Astronomy National Optical Astronomy Observatory</p> <p>Particle Physics Indra</p> <p>Plant genetics the iPlant Collaborative</p> <p>Quantum IN2P3</p> <p>Chromodynamics</p> <p>Radio Astronomy Cyber Square Kilometer Array, TREND, BAOradio</p> <p>Seismology Southern California Earthquake Center</p> <p>Social Science Odum Institute for Social Science Research, TerraPop</p>
More information (URLs)	<p>The DataNet Federation Consortium: http://www.datafed.org</p> <p>iRODS: http://www.irods.org</p>
<p>NOTE A major challenge is the ability to capture knowledge needed to interact with the data products of a research domain. In policy-based data management systems, this is done by encapsulating the knowledge in procedures that are controlled through policies. The procedures can automate retrieval of data from external repositories, or execute processing workflows, or enforce management policies on the resulting data products. A standard application is the enforcement of data management plans and the verification that the plan has been successfully applied.</p>	

A.6.2 Use case 33: The 'Discinnet Process'

Use case title	The 'Discinnet process', metadata <-> big data global experiment	
Vertical (area)	Scientific Research: Interdisciplinary Collaboration	
Author/company/email	P. Journeau / Discinnet Labs / phjourneau@discinnet.org	
Actors/stakeholders and their roles and responsibilities	Actors Richeact, Discinnet Labs and I4OpenResearch fund France/Europe. American equivalent pending. Richeact is fundamental research and development epistemology, Discinnet Labs applied in web 2.0 http://www.discinnet.org , 14 non-profit warrant.	
Goals	Richeact scientific goal is to reach predictive interdisciplinary model of research fields' behavior (with related meta-grammar). Experimentation through global sharing of now multidisciplinary, later interdisciplinary Discinnet process/web mapping and new scientific collaborative communication and publication system. Expected sharp impact to reducing uncertainty and time between theoretical, applied, technology research and development steps.	
Use case description	<p>Currently 35 clusters started, close to 100 awaiting more resources and potentially much more open for creation, administration and animation by research communities. Examples range from optics, cosmology, materials, microalgae, health to applied maths, computation, rubber and other chemical products/issues.</p> <p>How does a typical case currently work:</p> <ul style="list-style-type: none"> — A researcher or group wants to see how a research field is faring and in a minute defines the field on Discinnet as a 'cluster' — Then it takes another 5 to 10 mn to parameter the first/main dimensions, mainly measurement units and categories, but possibly later on some variable limited time for more dimensions — Cluster then may be filled either by doctoral students or reviewing researchers and/or communities/researchers for projects/progress <p>Already significant value but now needs to be disseminated and advertised although maximal value to come from interdisciplinary/projective next version. Value is to detect quickly a paper/project of interest for its results and next step is trajectory of the field under types of interactions from diverse levels of oracles (subjects/objects) + from interdisciplinary context.</p>	
Current solutions	Compute(System)	Currently on OVH (Hosting company http://www.ovh.co.uk/) servers (mix shared + dedicated)
	Storage	OVH
	Networking	To be implemented with desired integration with others
	Software	Current version with Symfony-PHP, Linux, MySQL
Big data characteristics	Data source (distributed/centralized)	Currently centralized, soon distributed per country and even per hosting institution interested by own platform
	Volume (size)	Not significant : this is a metadata base, not big data

	Velocity (e.g. real time)	Real time
	Variety (multiple datasets, mashup)	Link to big data still to be established in a Meta<->Big relationship not yet implemented (with experimental databases and already 1 st level related metadata)
	Variability (rate of change)	Currently real time, for further multiple locations and distributed architectures, periodic (such as nightly)
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Methods to detect overall consistency, holes, errors, misstatements, known but mostly to be implemented
	Visualization	Multidimensional (hypercube)
	Data quality (syntax)	A priori correct (directly human captured) with sets of checking + evaluation processes partly implemented
	Data types	'cluster displays' (image), vectors, categories, PDFs
	Data analytics	
Big data specific challenges (Gaps)	<p>Our goal is to contribute to Big 2 Metadata challenge by systematic reconciling between metadata from many complexity levels with ongoing input from researchers from ongoing research process.</p> <p>Current relationship with Richeact is to reach the interdisciplinary model, using meta-grammar itself to be experimented and its extent fully proven to bridge efficiently the gap between as remote complexity levels as semantic and most elementary (big) signals. Example with cosmological models versus many levels of intermediary models (particles, gases, galactic, nuclear, geometries). Others with computational versus semantic levels.</p>	
Big data specific challenges in mobility	Appropriate graphic interface power	
Security and privacy technical considerations	Several levels already available and others planned, up to physical access keys and isolated servers. Optional anonymity, usual protected exchanges	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Through 2011 to 2013, we have shown on http://www.discinnet.org that all kinds of research fields could easily get into Discinnet type of mapping, yet developing and filling a cluster requires time and/or dedicated workers.	

More information (URLs)	On http://www.discinnet.org the already started or starting clusters can be watched in one click on 'cluster' (field) title and even more detail is available through free registration (more resource available when registering as researcher (publications) or pending (doctoral student))
	Maximum level of detail is free for contributing researchers in order to protect communities but available to external observers for symbolic fee: all suggestions for improvements and better sharing welcome. We are particularly open to provide and support experimental appropriation by doctoral schools to build and study the past and future behavior of clusters in Earth sciences, Cosmology, Water, Health, Computation, Energy/Batteries, Climate models, Space, etc..
NOTE We are open to facilitate wide appropriation of both global, regional and local versions of the platform (for instance by research institutions, publishers, networks with desirable maximal data sharing for the greatest benefit of advancement of science.	

A.6.3 Use case 34: Graph Search on Scientific Data

Use case title	Enabling Face-Book like Semantic Graph-search on Scientific Chemical and Text-based Data
Vertical (area)	Management of Information from Research Articles
Author/company/email	Talapady Bhat, bhat@nist.gov
Actors/stakeholders and their roles and responsibilities	Chemical structures, Protein Data Bank, Material Genome Project, Open-GOV initiative, Semantic Web, Integrated Data-graphs, Scientific social media
Goals	Establish infrastructure, terminology and semantic data-graphs to annotate and present technology information using 'root' and rule-based methods used primarily by some Indo-European languages like Sanskrit and Latin.
Use case description	<ul style="list-style-type: none"> — Social media hype — Internet and social media play a significant role in modern information exchange. Every day most of us use social-media both to distribute and receive information. Two of the special features of many social media like Face-Book are — the community is both data-providers and data-users — they store information in a pre-defined 'data-shelf' of a data-graph — Their core infrastructure for managing information is reasonably language free — What this has to do with managing scientific information? <p>During the last few decades science has truly evolved to become a community activity involving every country and almost every household. We routinely 'tune-in' to internet resources to share and seek scientific information.</p> <ul style="list-style-type: none"> — What are the challenges in creating social media for science

	<ul style="list-style-type: none">— Creating a social media of scientific information needs an infrastructure where many scientists from various parts of the world can participate and deposit results of their experiment. Some of the issues that one has to resolve prior to establishing a scientific social media are:<ul style="list-style-type: none">— How to minimize challenges related to local language and its grammar?— How to determining the 'data-graph' to place an information in an intuitive way without knowing too much about the data management?— How to find relevant scientific data without spending too much time on the internet?
--	--

Copyrighted document, no reproduction or circulation
STANDARDSISO.COM: Click to view the full PDF of ISO/IEC TR 20547 WG-2:2018
For review by FC on PR in Healthcare
Oct 2024

	<p>Approach: Most languages and more so Sanskrit and Latin use a novel 'root'-based method to facilitate the creation of on-demand, discriminating words to define concepts. Some such examples from English are Bio-logy, Bio-chemistry. Youga, Yogi, Yogendra, Yogesh are examples from Sanskrit. Genocide is an example from Latin. These words are created on-demand based on best-practice terms and their capability to serve as node in a discriminating data-graph with self-explained meaning.</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>Cloud for the participation of community</p>
	<p>Storage</p>	<p>Requires expandable on-demand based resource that is suitable for global users location and requirements</p>
	<p>Networking</p>	<p>Needs good network for the community participation</p>
	<p>Software</p>	<p>Good database tools and servers for data-graph manipulation are needed</p>
<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>Distributed resource with a limited centralized capability</p>
	<p>Volume (size)</p>	<p>Undetermined. May be few terabytes at the beginning</p>
	<p>Velocity (e.g. real time)</p>	<p>Evolving with time to accommodate new best-practices</p>
	<p>Variety (multiple datasets, mashup)</p>	<p>Wildly varying depending on the types available technological information</p>
	<p>Variability (rate of change)</p>	<p>Data-graphs are likely to change in time based on customer preferences and best-practices</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Technological information is likely to be stable and robust</p>
	<p>Visualization</p>	<p>Efficient data-graph based visualization is needed</p>
	<p>Data quality (syntax)</p>	<p>Expected to be good</p>
	<p>Data types</p>	<p>All data types, image to text, structures to protein sequence</p>
	<p>Data analytics</p>	<p>Data-graphs is expected to provide robust data-analysis methods</p>
<p>Big data specific challenges (Gaps)</p>	<p>This is a community effort similar to many social media. Providing a robust, scalable, on-demand infrastructures in a manner that is use-case and user-friendly is a real challenge by any existing conventional methods</p>	

Big data specific challenges in mobility	A community access is required for the data and thus it has to be media and location independent and thus requires high mobility too.
Security and privacy technical considerations	None since the effort is initially focused on publicly accessible data provided by open-platform projects like open-gov, MGI and protein data bank.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	This effort includes many local and networked resources. Developing an infrastructure to automatically integrate information from all these resources using data-graphs is a challenge that we are trying to solve.
More information (URLs)	http://www.eurekalert.org/pub_releases/2013-07/aiop-ffm071813.php http://xpdb.nist.gov/chemblast/pdb.pl http://xpdb.nist.gov/chemblast/pdb.pl
<p>NOTE Many reports, including a recent one on Material Genome Project finds that exclusive top-down solutions to facilitate data sharing and integration are not desirable for federated multi-disciplinary efforts. However, a bottom-up approach can be chaotic. For this reason, there is need for a balanced blend of the two approaches to support easy-to-use techniques to metadata creation, integration and sharing. This challenge is very similar to the challenge faced by language developer at the beginning. One of the successful effort used by many prominent languages is that of 'roots' and rules that form the framework for creating on-demand words for communication. In this approach a top-down method is used to establish a limited number of highly re-usable words called 'roots' by surveying the existing best practices in building terminology. These 'roots' are combined using few 'rules' to create terms on-demand by a bottom-up step.</p> <p>Y(uj) (join), O (creator, God, brain), Ga (motion, initiation) -leads to 'Yoga' in Sanskrit, English Geno (genos)-cide-race based killing - Latin, English Bio-technology -English, Latin Red-light, red-laser-light -English.</p> <p>A press release by the American Institute of Physics on this approach is at http://www.eurekalert.org/pub_releases/2013-07/aiop-ffm071813.php</p> <p>Our efforts to develop automated and rule and root-based methods (Chem-BLAST - http://xpdb.nist.gov/chemblast/pdb.pl) to identify and use best-practice, discriminating terms in generating semantic data-graphs for science started almost a decade back with a chemical structure database. This database has millions of structures obtained from the Protein Data Bank and the PubChem used world-wide. Subsequently we extended our efforts to build root-based terms to text-based data of cell-images. In this work we use few simple rules to define and extend terms based on best-practice as decided by weaning through millions of popular use-cases chosen from over hundred biological ontologies.</p> <p>Currently we are working on extending this method to publications of interest to Material Genome, Open-Gov and NIST-wide publication archive - NIKE. - http://xpdb.nist.gov/nike/term.pl. These efforts are a component of Research Data Alliance Working Group on Metadata https://www.rd-alliance.org/filedepot_download/694/160 and https://rd-alliance.org/poster-session-rda-2nd-plenary-meeting.html</p>	

A.6.4 Use case 35: Light Source Beamlines

Use case title	Light source beamlines
Vertical (area)	Research (Biology, Chemistry, Geophysics, Materials Science, others)
Author/company/email	Eli Dart, LBNL (eddart@lbl.gov)
Actors/stakeholders and their roles and responsibilities	Research groups from a variety of scientific disciplines (see above)

Goals	Use of a variety of experimental techniques to determine structure, composition, behavior, or other attributes of a sample relevant to scientific enquiry.	
Use case description	Samples are exposed to X-rays in a variety of configurations depending on the experiment. Detectors (essentially high-speed digital cameras) collect the data. The data are then analyzed to reconstruct a view of the sample or process being studied. The reconstructed images are used by scientists analysis.	
Current solutions	Compute(System)	Computation ranges from single analysis hosts to high-throughput computing systems at computational facilities
	Storage	Local storage on the order of 1-40 TB on Windows or Linux data servers at facility for temporary storage, over 60 TB on disk at NERSC, over 300 TB on tape at NERSC
	Networking	10 Gbps Ethernet at facility, 100 Gbps to NERSC
	Software	A variety of commercial and open source software is used for data analysis – examples include: — Octopus (http://www.inct.be/en/software/octopus) for Tomographic Reconstruction — Avizo (http://vsg3d.com) and FIJI (a distribution of ImageJ; http://fiji.sc) for Visualization and Analysis Data transfer is accomplished using physical transport of portable media (severely limits performance) or using high-performance GridFTP, managed by Globus Online or workflow systems such as SPADE.
Big data characteristics	Data source (distributed/centralized)	Centralized (high resolution camera at facility). Multiple beamlines per facility with high-speed detectors.
	Volume (size)	3 GB to 30 GB per sample – up to 15 samples/day
	Velocity (e.g. real time)	Near real-time analysis needed for verifying experimental parameters (lower resolution OK). Automation of analysis would dramatically improve scientific productivity.
	Variety (multiple datasets, mashup)	Many detectors produce similar types of data (e.g. TIFF files), but experimental context varies widely
	Variability (rate of change)	Detector capabilities are increasing rapidly. Growth is essentially Moore's Law. Detector area is increasing exponentially (1k × 1k, 2k × 2k, 4k × 4k, ...) and read-out is increasing exponentially (1 Hz, 10 Hz, 100 Hz, 1 kHz, ...). Single detector data rates are expected to reach 1 GB per second within 2 years.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Near real-time analysis required to verify experimental parameters. In many cases, early analysis can dramatically improve experiment productivity by providing early feedback. This implies high-throughput computing, high-performance data transfer, and high-speed storage are routinely available.
	Visualization	Visualization is key to a wide variety of experiments at all light source facilities
	Data quality (syntax)	Data quality and precision are critical (especially since beam time is scarce, and re-running an experiment is often impossible).
	Data types	Many beamlines generate image data (e.g. TIFF files)

	Data analytics	Volume reconstruction, feature identification, others
Big data specific challenges (Gaps)	Rapid increase in camera capabilities, need for automation of data transfer and near-real-time analysis.	
Big data specific challenges in mobility	Data transfer to large-scale computing facilities is becoming necessary because of the computational power required to conduct the analysis on time scales useful to the experiment. Large number of beamlines (e.g. 39 at LBNL ALS) means that aggregate data load is likely to increase significantly over the coming years.	
Security and privacy technical considerations	Varies with project.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	There will be significant need for a generalized infrastructure for analyzing GBs per second of data from many beamline detectors at multiple facilities. Prototypes exist now, but routine deployment will require additional resources.	
More information (URLs)	http://www-als.lbl.gov/ http://www.aps.anl.gov/ https://portal.slac.stanford.edu/sites/lcls_public/Pages/Default.aspx	

STANDARDSISO.COM: Click to view the full PDF of ISO/IEC TR 20547-2:2018

Copyrighted document, no reproduction or circulation

Oct 2024

A.7 Astronomy and Physics

A.7.1 Use case 36: Catalina Digital Sky Survey for Transients

Use case title	Catalina Real-Time Transient Survey (CRTS): a digital, panoramic, synoptic sky survey
Vertical (area)	Scientific Research: Astronomy
Author/company/email	S. G. Djorgovski / Caltech / george@astro.caltech.edu
Actors/stakeholders and their roles and responsibilities	<p>The survey team: data processing, quality control, analysis and interpretation, publishing, and archiving.</p> <p>Collaborators: a number of research groups world-wide: further work on data analysis and interpretation, follow-up observations, and publishing.</p> <p>User community: all of the above plus the astronomical community world-wide: further work on data analysis and interpretation, follow-up observations, and publishing.</p>
Goals	<p>The survey explores the variable universe in the visible light regime, on time scales ranging from minutes to years, by searching for variable and transient sources. It discovers a broad variety of astrophysical objects and phenomena, including various types of cosmic explosions (e.g. Supernovae), variable stars, phenomena associated with accretion to massive black holes (active galactic nuclei) and their relativistic jets, high proper motion stars, etc.</p>
Use case description	<p>The data are collected from 3 telescopes (2 in Arizona and 1 in Australia), with additional ones expected in the near future (in Chile). The original motivation is a search for near-Earth (NEO) and potential planetary hazard (PHO) asteroids, funded by NASA, and conducted by a group at the Lunar and Planetary Laboratory (LPL) at the Univ. of Arizona (UA); that is the Catalina Sky Survey proper (CSS). The data stream is shared by the CRTS for the purposes for exploration of the variable universe, beyond the Solar system, led by the Caltech group. Approximately 83 % of the entire sky is being surveyed through multiple passes (crowded regions near the Galactic plane, and small areas near the celestial poles are excluded).</p> <p>The data are preprocessed at the telescope, and transferred to LPL/UA, and hence to Caltech, for further analysis, distribution, and archiving. The data are processed in real time, and detected transient events are published electronically through a variety of dissemination mechanisms, with no proprietary period (CRTS has a completely open data policy).</p>

	<p>Further data analysis includes automated and semi-automated classification of the detected transient events, additional observations using other telescopes, scientific interpretation, and publishing. In this process, it makes a heavy use of the archival data from a wide variety of geographically distributed resources connected through the Virtual Observatory (VO) framework.</p> <p>Light curves (flux histories) are accumulated for ~ 500 million sources detected in the survey, each with a few hundred data points on average, spanning up to 8 years, and growing. These are served to the community from the archives at Caltech, and shortly from IUCAA, India. This is an unprecedented data set for the exploration of time domain in astronomy, in terms of the temporal and area coverage and depth.</p> <p>CRTS is a scientific and methodological testbed and precursor of the grander surveys to come, notably the Large Synoptic Survey Telescope (LSST), expected to operate in 2020's.</p>								
<p>Current solutions</p>	<table border="1"> <tr> <td data-bbox="708 875 1054 1256"> <p>Compute(System)</p> </td> <td data-bbox="1054 875 1398 1256"> <p>Instrument and data processing computers: a number of desktop and small server class machines, although more powerful machinery is needed for some data analysis tasks.</p> <p>This is not so much a computationally-intensive project, but rather a data-handling-intensive one.</p> </td> </tr> <tr> <td data-bbox="708 1256 1054 1330"> <p>Storage</p> </td> <td data-bbox="1054 1256 1398 1330"> <p>Several multi-TB / tens of TB servers.</p> </td> </tr> <tr> <td data-bbox="708 1330 1054 1404"> <p>Networking</p> </td> <td data-bbox="1054 1330 1398 1404"> <p>Standard inter-university internet connections.</p> </td> </tr> <tr> <td data-bbox="708 1404 1054 1608"> <p>Software</p> </td> <td data-bbox="1054 1404 1398 1608"> <p>Custom data processing pipeline and data analysis software, operating under Linux. Some archives on Windows machines, running a MS SQL server databases.</p> </td> </tr> </table>	<p>Compute(System)</p>	<p>Instrument and data processing computers: a number of desktop and small server class machines, although more powerful machinery is needed for some data analysis tasks.</p> <p>This is not so much a computationally-intensive project, but rather a data-handling-intensive one.</p>	<p>Storage</p>	<p>Several multi-TB / tens of TB servers.</p>	<p>Networking</p>	<p>Standard inter-university internet connections.</p>	<p>Software</p>	<p>Custom data processing pipeline and data analysis software, operating under Linux. Some archives on Windows machines, running a MS SQL server databases.</p>
<p>Compute(System)</p>	<p>Instrument and data processing computers: a number of desktop and small server class machines, although more powerful machinery is needed for some data analysis tasks.</p> <p>This is not so much a computationally-intensive project, but rather a data-handling-intensive one.</p>								
<p>Storage</p>	<p>Several multi-TB / tens of TB servers.</p>								
<p>Networking</p>	<p>Standard inter-university internet connections.</p>								
<p>Software</p>	<p>Custom data processing pipeline and data analysis software, operating under Linux. Some archives on Windows machines, running a MS SQL server databases.</p>								
<p>Big data characteristics</p>	<table border="1"> <tr> <td data-bbox="708 1608 1054 1917"> <p>Data source (distributed/centralized)</p> </td> <td data-bbox="1054 1608 1398 1917"> <p>Distributed:</p> <ol style="list-style-type: none"> 1) Survey data from 3 (soon more?) telescopes 2) Archival data from a variety of resources connected through the VO framework 3) Follow-up observations from separate telescopes </td> </tr> </table>	<p>Data source (distributed/centralized)</p>	<p>Distributed:</p> <ol style="list-style-type: none"> 1) Survey data from 3 (soon more?) telescopes 2) Archival data from a variety of resources connected through the VO framework 3) Follow-up observations from separate telescopes 						
<p>Data source (distributed/centralized)</p>	<p>Distributed:</p> <ol style="list-style-type: none"> 1) Survey data from 3 (soon more?) telescopes 2) Archival data from a variety of resources connected through the VO framework 3) Follow-up observations from separate telescopes 								

	Volume (size)	The survey generates up to ~ 0,1 TB per clear night; ~ 100 TB in current data holdings. Follow-up observational data amount to no more than a few % of that. Archival data in external (VO-connected) archives are in PBs, but only a minor fraction is used.
	Velocity (e.g. real time)	Up to ~ 0,1 TB / night of the raw survey data.
	Variety (multiple datasets, mashup)	The primary survey data in the form of images, processed to catalogs of sources (db tables), and time series for individual objects (light curves). Follow-up observations consist of images and spectra. Archival data from the VO data grid include all of the above, from a wide variety of sources and different wavelengths.
	Variability (rate of change)	Daily data traffic fluctuates from ~ 0,01 to ~ 0,1 TB / day, not including major data transfers between the principal archives (Caltech, UA, and IUCAA).
	Veracity (Robustness Issues, semantics)	A variety of automated and human inspection quality control mechanisms is implemented at all stages of the process.
Big data science (collection, curation, analysis, action)	Visualization	Standard image display and data plotting packages are used. We are exploring visualization mechanisms for highly dimensional data parameter spaces.
	Data quality (syntax)	It varies, depending on the observing conditions, and it is evaluated automatically: error bars are estimated for all relevant quantities.
	Data types	Images, spectra, time series, catalogs.

	Data analytics	A wide variety of the existing astronomical data analysis tools, plus a large amount of custom developed tools and software, some of it a research project in itself.
Big data specific challenges (Gaps)	<p>Development of machine learning tools for data exploration, and in particular for an automated, real-time classification of transient events, given the data sparsity and heterogeneity.</p> <p>Effective visualization of hyper-dimensional parameter spaces is a major challenge for all of us.</p>	
Big data specific challenges in mobility	Not a significant limitation at this time.	
Security and privacy technical considerations	None.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	<ul style="list-style-type: none"> — Real-time processing and analysis of massive data streams from a distributed sensor network (in this case telescopes), with a need to identify, characterize, and respond to the transient events of interest in (near) real time. — Use of highly distributed archival data resources (in this case VO-connected archives) for data analysis and interpretation. — Automated classification given the very sparse and heterogeneous data, dynamically evolving in time as more data come in, and follow-up decision making given limited and sparse resources (in this case follow-up observations with other telescopes). 	
More information (URLs)	<p>CRTS survey: http://crts.caltech.edu</p> <p>CSS survey: http://www.lpl.arizona.edu/css</p> <p>For an overview of the classification challenges, see, e.g., http://arxiv.org/abs/1209.1681</p> <p>For a broader context of sky surveys, past, present, and future, see, e.g., the review http://arxiv.org/abs/1209.1681</p>	
<p>NOTE CRTS can be seen as a good precursor to the astronomy’s flagship project, the Large Synoptic Sky Survey (LSST; http://www.lsst.org), now under development. Their anticipated data rates (~ 20TB to 30TB per clear night, tens of PB over the duration of the survey) are directly on the Moore’s law scaling from the current CRTS data rates and volumes, and many technical and methodological issues are very similar.</p> <p>It is also a good case for real-time data mining and knowledge discovery in massive data streams, with distributed data sources and computational resources.</p>		

A.7.2 Use case 37: Cosmological Sky Survey and Simulations

Use case title	DOE Extreme Data from Cosmological Sky Survey and Simulations	
Vertical (area)	Scientific Research: Astrophysics	
Author/company/email	PIs: Salman Habib, Argonne National Laboratory; Andrew Connolly, University of Washington	
Actors/stakeholders and their roles and responsibilities	Researchers studying dark matter, dark energy, and the structure of the early universe.	
Goals	Clarify the nature of dark matter, dark energy, and inflation, some of the most exciting, perplexing, and challenging questions facing modern physics. Emerging, unanticipated measurements are pointing toward a need for physics beyond the successful Standard Model of particle physics.	
Use case description	<p>This investigation requires an intimate interplay between big data from experiment and simulation as well as massive computation. The melding of all will</p> <ol style="list-style-type: none"> 1) Provide the direct means for cosmological discoveries that require a strong connection between theory and observations ('precision cosmology'); 2) Create an essential 'tool of discovery' in dealing with large datasets generated by complex instruments; and, 3) Generate and share results from high-fidelity simulations that are necessary to understand and control systematics, especially astrophysical systematics. 	
Current solutions	Compute(System)	Hours: 24M (NERSC / Berkeley Lab), 190M (ALCF / Argonne), 10M (OLCF / Oak Ridge)
	Storage	180 TB (NERSC / Berkeley Lab)
	Networking	ESNet connectivity to the national labs is adequate today.
	Software	MPI, OpenMP, C, C++, F90, FFTW, viz packages, python, FFTW, numpy, Boost, OpenMP, ScaLAP-CK, PSQL and MySQL databases, Eigen, cfitsio, astrometry.net, and Minuit2
Big data characteristics	Data source (distributed/centralized)	Observational data will be generated by the Dark Energy Survey (DES) and the Zwicky Transient Factory in 2015 and by the Large Synoptic Sky Survey starting in 2019. Simulated data will generated at DOE super-computing centers.
	Volume (size)	DES: 4 PB, ZTF 1 PB/year, LSST 7 PB/year, Simulations > 10 PB in 2017
	Velocity (e.g. real time)	LSST: 20 TB/day

	Variety (multiple datasets, mashup)	1) Raw Data from sky surveys 2) Processed Image data 3) Simulation data
	Variability (rate of change)	Observations are taken nightly; supporting simulations are run throughout the year, but data can be produced sporadically depending on access to resources
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	
	Visualization	Interpretation of results from detailed simulations requires advanced analysis and visualization techniques and capabilities. Supercomputer I/O subsystem limitations are forcing researchers to explore “in-situ” analysis to replace post-processing methods.
	Data quality (syntax)	
	Data types	Image data from observations must be reduced and compared with physical quantities derived from simulations. Simulated sky maps must be produced to match observational formats.
	Data analytics	
Big data specific challenges (Gaps)	Storage, sharing, and analysis of 10s of PBs of observational and simulated data.	
Big data specific challenges in mobility	LSST will produce 20 TB of data per day. This must be archived and made available to researchers world-wide.	
Security and privacy technical considerations		
Highlight issues for generalizing this Use case (e.g. for ref. architecture)		
More information (URLs)	http://www.lsst.org/lsst/ http://www.nersc.gov/ http://science.energy.gov/hep/research/non-accelerator-physics/ http://www.nersc.gov/assets/Uploads/HabibcosmosimV2.pdf	

A.7.3 Use case 38: Large Survey Data for Cosmology

Use case title	Large Survey Data for Cosmology	
Vertical (area)	Scientific Research: Cosmic Frontier	
Author/company/email	Peter Nugent / LBNL / eszeto@lbl.gov	
Actors/stakeholders and their roles and responsibilities	Dark Energy Survey, Dark Energy Spectroscopic Instrument, Large Synoptic Survey Telescope. ANL, BNL, FNAL, LBL and SLAC: Create the instruments/telescopes, run the survey and perform the cosmological analysis.	
Goals	Provide a way to reduce photometric data in real time for supernova discovery and follow-up and to handle the large volume of observational data (in conjunction with simulation data) to reduce systematic uncertainties in the measurement of the cosmological parameters via baryon acoustic oscillations, galaxy cluster counting and weak lensing measurements.	
Use case description	For DES the data are sent from the mountaintop via a microwave link to La Serena, Chile. From there, an optical link forwards them to the NCSA as well as NERSC for storage and "reduction". Subtraction pipelines are run using extant imaging data to find new optical transients through machine learning algorithms. Then galaxies and stars in both the individual and stacked images are identified, catalogued, and finally their properties measured and stored in a database.	
Current solutions	Compute(System)	Linux cluster, Oracle RDBMS server, large memory machines, standard Linux interactive hosts. For simulations, HPC resources.
	Storage	Oracle RDBMS, PostgreSQL, as well as GPFS and Lustre file systems and tape archives.
	Networking	Provided by NERSC
	Software	Standard astrophysics reduction software as well as Perl/Python wrapper scripts, Linux Cluster scheduling and comparison to large amounts of simulation data via techniques like Cholesky decomposition.
Big data characteristics	Data source (distributed/centralized)	Distributed. Typically between observation and simulation data.
	Volume (size)	LSST will generate 60 PB of imaging data and 15 PB of catalog data and a correspondingly large (or larger) amount of simulation data. Over 20 TB of data per night.

	Velocity (e.g. real time)	20 TB of data will have to be subtracted each night in as near real time as possible in order to maximize the science for supernovae.
	Variety (multiple datasets, mashup)	While the imaging data is similar, the analysis for the 4 different types of cosmological measurements and comparisons to simulation data is quite different.
	Variability (rate of change)	Weather and sky conditions can radically change both the quality and quantity of data.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Astrophysical data is a statistician's nightmare as the both the uncertainties in a given measurement change from night-to-night in addition to the cadence being highly unpredictable. Also, most all of the cosmological measurements are systematically limited, and thus understanding these as best possible is the highest priority for a given survey.
	Visualization	Interactive speed of web UI on very large data sets is an ongoing challenge. Basic querying and browsing of data to find new transients as well as monitoring the quality of the survey is a must. Ability to download large amounts of data for offline analysis is another requirement of the system. Ability to combine both simulation and observational data is also necessary.
	Data quality (syntax)	Understanding the systematic uncertainties in the observational data is a prerequisite to a successful cosmological measurement. Beating down the uncertainties in the simulation data to under this level is a huge challenge for future surveys.
	Data types	Cf. above on "Variety"
	Data analytics	

Big data specific challenges (Gaps)	New statistical techniques for understanding the limitations in simulation data would be beneficial. Often it is the case where there is not enough computing time to generate all the simulations one wants and thus there is a reliance on emulators to bridge the gaps. Techniques for handling Cholesky decomposition for thousands of simulations with matrices of order 1M on a side.
Big data specific challenges in mobility	Performing analysis on both the simulation and observational data simultaneously.
Security and privacy technical considerations	No special challenges. Data is either public or requires standard login with password.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Parallel databases which could handle imaging data would be an interesting avenue for future research.
More information (URLs)	http://www.lsst.org/lsst , http://desi.lbl.gov , and http://www.darkenergysurvey.org

A.7.4 Use case 39: Analysis of LHC (Large Hadron Collider) Data

Use case title	Particle Physics: Analysis of LHC (Large Hadron Collider) Data (Discovery of Higgs particle)	
Vertical (area)	Scientific Research: Physics	
Author/company/email	Michael Ernst mernst@bnl.gov , Lothar Bauerdick bauerdick@fnal.gov based on an initial version written by Geoffrey Fox, Indiana University gcf@indiana.edu , Eli Dart, LBNL eddart@lbl.gov	
Actors/stakeholders and their roles and responsibilities	Physicists (Design and Identify need for Experiment, Analyze Data) Systems Staff (Design, Build and Support distributed Computing Grid), Accelerator Physicists (Design, Build and Run Accelerator), Government (funding based on long term importance of discoveries in field)	
Goals	Understanding properties of fundamental particles	
Use case description	CERN LHC Detectors and Monte Carlo producing events describing particle-apparatus interaction. Processed information defines physics properties of events (lists of particles with type and momenta). These events are analyzed to find new effects; both new particles (Higgs) and present evidence that conjectured particles (Supersymmetry) not seen.	
Current solutions	Compute(System)	WLCG and Open Science Grid in the US integrate computer centers worldwide that provide computing and storage resources into a single infrastructure accessible by all LHC physicists. 350,000 cores running "continuously" arranged in 3 tiers (CERN, "Continents/Countries". "Universities"). Uses "Distributed High Throughput Computing (DHTC)"; 200 PB storage, >2 million jobs/day.

	<p>Storage</p> <p>ATLAS:</p> <ul style="list-style-type: none"> — Brookhaven National Laboratory Tier1 tape: 10PB ATLAS data on tape managed by HPSS (incl. RHIC/NP the total data volume is 35PB) — Brookhaven National Laboratory Tier1 disk: 11 PB; using dCache to virtualize a set of ~60 heterogeneous storage servers with high-density disk backend systems — US Tier2 centers, disk cache: 16PB <p>CMS:</p> <ul style="list-style-type: none"> — Fermilab US Tier1, reconstructed, tape/cache: 20,4 PB — US Tier2 centers, disk cache: 7 PB — US Tier3 sites, disk cache: 1,04 PB
	<p>Networking</p> <ul style="list-style-type: none"> — As experiments have global participants (CMS has 3600 participants from 183 institutions in 38 countries), the data at all levels is transported and accessed across continents. — Large scale automated data transfers occur over science networks across the globe. LHCOPN and LHCONE network overlay provide dedicated network allocations and traffic isolation for LHC data traffic — ATLAS Tier1 data center at BNL has 160Gbps internal paths (often fully loaded). 70Gbps WAN connectivity provided by ESnet. — CMS Tier1 data center at FNAL has 90Gbps WAN connectivity provided by ESnet — Aggregate wide area network traffic for LHC experiments is about 25Gbps steady state worldwide <p>Software</p> <p>The scalable ATLAS workload/workflow management system PanDA manages ~1 million production and user analysis jobs on globally distributed computing resources (~100 sites) per day.</p> <p>The new ATLAS distributed data management system Rucio is the core component keeping track of an inventory of currently ~130PB of data distributed across grid resources and to orchestrate data movement between sites. The data volume is expected to grow to exascale size in the next few years. Based on the xrootd system ATLAS has developed FAX, a federated storage system that allows remote data access.</p> <p>Similarly, CMS is using the OSG glideinWMS infrastructure to manage its workflows for production and data analysis the PhEDEx system to orchestrate data movements, and the AAA/xrootd system to allow remote data access.</p> <p>Experiment-specific physics software including simulation packages, data processing, advanced statistic packages, etc.</p>

Big data characteristics	Data source (distributed/centralized)	High speed detectors produce large data volumes: <ul style="list-style-type: none"> — ATLAS detector at CERN: Originally 1 PB/sec raw data rate, reduced to 300 MB/sec by multi-stage trigger. — CMS detector at CERN: similar Data distributed to Tier1 centers globally, which serve as data sources for Tier2 and Tier3 analysis centers
	Volume (size)	15 Petabytes per year from Detectors and Analysis
	Velocity (e.g. real time)	<ul style="list-style-type: none"> — Real time with some long LHC "shut downs" (to improve accelerator and detectors) with no data except Monte Carlo. — Besides using programmatically and dynamically replicated datasets, real-time remote I/O (using XrootD) is increasingly used by analysis which requires reliable high-performance networking capabilities to reduce file copy and storage system overhead
	Variety (multiple datasets, mashup)	Lots of types of events with from 2- few hundred final particle but all data is collection of particles after initial analysis. Events are grouped into datasets; real detector data is segmented into ~20 datasets (with partial overlap) on the basis of event characteristics determined through real-time trigger system, while different simulated datasets are characterized by the physics process being simulated.
	Variability (rate of change)	Data accumulates and does not change character. What you look for may change based on physics insight. As understanding of detectors increases, large scale data reprocessing tasks are undertaken.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	One can lose modest amount of data without much pain as errors proportional to 1/SquareRoot(Events gathered), but such data loss must be carefully accounted. Importance that accelerator and experimental apparatus work both well and in understood fashion. Otherwise data too "dirty"/"uncorrectable".
	Visualization	Modest use of visualization outside histograms and model fits. Nice event displays but discovery requires lots of events so this type of visualization of secondary importance
	Data quality (syntax)	Huge effort to make certain complex apparatus well understood (proper calibrations) and "corrections" properly applied to data. Often requires data to be re-analyzed
	Data types	Raw experimental data in various binary forms with conceptually a name: value syntax for name spanning "chamber readout" to "particle momentum". Reconstructed data is processed to produce dense data formats optimized for analysis

	<p>Data analytics</p> <p>Initial analysis is processing of experimental data specific to each experiment (ALICE, ATLAS, CMS, LHCb) producing summary information. Second step in analysis uses “exploration” (histograms, scatter-plots) with model fits. Substantial Monte-Carlo computations are necessary to estimate analysis quality.</p> <p>A large fraction (~60 %) of the available CPU resources available to the ATLAS collaboration at the Tier-1 and the Tier-2 centers is used for simulated event production. The ATLAS simulation requirements are completely driven by the physics community in terms of analysis needs and corresponding physics goals. The current physics analyses are looking at real data samples of roughly 2 billion (B) events taken in 2011 and 3B events taken in 2012 (this represents ~5 PB of experimental data), and ATLAS has roughly 3,5B MC events for 2011 data, and 2,5B MC events for 2012 (this represents ~6 PB of simulated data). Given the resource requirements to fully simulate an event using the GEANT 4 package, ATLAS can currently produce about 4 million events per day using the entire processing capacity available to production worldwide.</p> <p>Due to its high CPU cost, the outputs of full Geant4 simulation (HITS) are stored in one custodial tape copy on Tier1 tapes to be re-used in several Monte-Carlo re-processings. The HITS from faster simulation flavors will be only of transient nature in LHC Run 2.</p>
<p>Big data specific challenges (Gaps)</p>	<p>The translation of scientific results into new knowledge, solutions, policies and decisions is foundational to the science mission associated with LHC data analysis and HEP in general. However, while advances in experimental and computational technologies have led to an exponential growth in the volume, velocity, and variety of data available for scientific discovery, advances in technologies to convert this data into actionable knowledge have fallen far short of what the HEP community needs to deliver timely and immediately impacting outcomes. Acceleration of the scientific knowledge discovery process is essential if DOE scientists are to continue making major contributions in HEP.</p> <p>Today's worldwide analysis engine, serving several thousand scientists, will have to be commensurately extended in the cleverness of its algorithms, the automation of the processes, and the reach (discovery) of the computing, to enable scientific understanding of the detailed nature of the Higgs boson. E.g. the approximately forty different analysis methods used to investigate the detailed characteristics of the Higgs boson (many using machine learning techniques) must be combined in a mathematically rigorous fashion to have an agreed upon publishable result.</p> <p><i>Specific challenges: Federated semantic discovery:</i> Interfaces, protocols and environments that support access to, use of, and interoperation across federated sets of resources governed and managed by a mix of different policies and controls that interoperate across streaming and “at rest” data sources. These include: models, algorithms, libraries, and reference implementations for a distributed non-hierarchical discovery service; semantics, methods, interfaces for life-cycle management (subscription, capture, provenance, assessment, validation, rejection) of heterogeneous sets of distributed tools, services and resources; a global environment that is robust in the face of failures and outages; and flexible high-performance data stores (going beyond schema driven) that scale and are friendly to interactive analytics</p>

	<p><i>Resource description and understanding:</i> Distributed methods and implementations that allow resources (people, software, computing incl. data) to publish varying state and function for use by diverse clients. Mechanisms to handle arbitrary entity types in a uniform and common framework – including complex types such as heterogeneous data, incomplete and evolving information, and rapidly changing availability of computing, storage and other computational resources. Abstract data streaming and file-based data movement over the WAN/LAN and on exascale architectures to allow for real-time, collaborative decision making for scientific processes.</p>
<p>Big data specific challenges in mobility</p>	<p>The agility to use any appropriate available resources and to ensure that all data needed is dynamically available at that resource is fundamental to future discoveries in HEP. In this context “resource” has a broad meaning and includes data and people as well as computing and other non-computer based entities: thus, any kind of data—raw data, information, knowledge, etc., and any type of resource—people, computers, storage systems, scientific instruments, software, resource, service, etc. In order to make effective use of such resources, a wide range of management capabilities must be provided in an efficient, secure, and reliable manner, encompassing for example collection, discovery, allocation, movement, access, use, release, and reassignment. These capabilities must span and control large ensembles of data and other resources that are constantly changing and evolving, and will often be in-deterministic and fuzzy in many aspects.</p> <p><i>Specific Challenges: Globally optimized dynamic allocation of resources:</i> These need to take account of the lack of strong consistency in knowledge across the entire system.</p> <p><i>Minimization of time-to-delivery of data and services:</i> Not only to reduce the time to delivery of the data or service but also allow for a predictive capability, so physicists working on data analysis can deal with uncertainties in the real-time decision making processes.</p>
<p>Security and privacy technical considerations</p>	<p>While HEP data itself is not proprietary unintended alteration and/or cyber-security related facility service compromises could potentially be very disruptive to the analysis process. Besides the need of having personal credentials and the related virtual organization credential management systems to maintain access rights to a certain set of resources, a fair amount of attention needs to be devoted to the development and operation of the many software components the community needs to conduct computing in this vastly distributed environment.</p> <p>The majority of software and systems development for LHC data analysis is carried out inside the HEP community or by adopting software components from other parties which involves numerous assumptions and design decisions from the early design stages throughout its life cycle. Software systems make a number of assumptions about their environment - how they are deployed, configured, who runs it, what sort of network is it on, is its input or output sensitive, can it trust its input, does it preserve privacy, etc.? When multiple software components are interconnected, for example in the deep software stacks used in DHTC, without clear understanding of their security assumptions, the security of the resulting system becomes an unknown.</p> <p>A trust framework is a possible way of addressing this problem. A DHTC trust framework, by describing what software, systems and organizations provide and expect of their environment regarding policy enforcement, security and privacy, allows for a system to be analyzed for gaps in trust, fragility and fault tolerance.</p>

<p>Highlight issues for generalizing this Use case (e.g. for ref. architecture)</p>	<p>Large scale example of an event based analysis with core statistics needed. Also highlights importance of virtual organizations as seen in global collaboration. The LHC experiments are pioneers of distributed big data science infrastructure, and several aspects of the LHC experiments' workflow highlight issues that other disciplines will need to solve. These include automation of data distribution, high performance data transfer, and large-scale high-throughput computing.</p>
<p>More information (URLs)</p>	<p>http://grids.ucs.indiana.edu/ptliupages/publications/Where%20does%20all%20the%20data%20come%20from%20v7.pdf http://www.es.net/assets/pubs_presos/High-throughput-lessons-from-the-LHC-experience.Johnston.TNC2013.pdf</p>
<p>Note:</p>	

Copyrighted document, no reproduction or circulation
 STANDARDSISO.COM: Click to view the full PDF of ISO/IEC TR 20547 WG-2:2018
 For review by FC on 21/11/2018
 Oct 2024

Use case Stages	Data Sources	Data Usage	Transformations (Data Analytics)	Infrastructure	Security and Privacy
Particle Physics: Analysis of LHC Large Hadron Collider Data, Discovery of Higgs particle (Scientific Research: Physics)					
Record Raw Data	CERN LHC Accelerator	This data is staged at CERN and then distributed across the globe for next stage in processing	LHC has 10^9 collisions per second; the hardware + software trigger selects "interesting events". Other utilities distribute data across the globe with fast transport	Accelerator and sophisticated data selection (trigger process) that uses ~7000 cores at CERN to record ~100-500 events each second (~1 megabyte each)	N/A
Process Raw Data to Information	Disk Files of Raw Data	Iterative calibration and checking of analysis which has for example "heuristic" track finding algorithms. Produce "large" full physics files and stripped down Analysis Object Data (AOD) files that are ~10 % original size	Full analysis code that builds in complete understanding of complex experimental detector. Also Monte Carlo codes to produce simulated data to evaluate efficiency of experimental detection.	~300,000 cores arranged in 3 tiers. Tier 0: CERN Tier 1: "Major Countries" Tier 2: Universities and laboratories. Note processing is compute and data intensive	N/A

<p>Physics Analysis Information to Knowledge/Discovery</p>	<p>Disk Files of Information including accelerator and Monte Carlo data. Include wisdom from lots of physicists (papers) in analysis choices</p>	<p>Use simple statistical techniques (like histogramming, multi-variate analysis methods and other data analysis techniques and model fits to discover new effects (particles) and put limits on effects not seen</p>	<p>Data reduction and processing steps with advanced physics algorithms to identify event properties, particle hypothesis etc. For interactive data analysis of those reduced and selected data sets the classic program is Root from CERN that reads multiple event (AOD, NTUP) files from selected data sets and use physicist generated C++ code to calculate new quantities such as implied mass of an unstable (new) particle</p>	<p>While the bulk of data processing is done at Tier 1 and Tier 2 resources, the end stage analysis is usually done by users at a local Tier 3 facility. The scale of computing resources at Tier 3 sites range from workstations to small clusters. ROOT is the most common software stack used to analyze compact data formats generated on distributed computing resources. Data transfer is done using ATLAS and CMS DDM tools, which mostly rely on gridFTP middleware. XROOTD based direct data access is also gaining importance wherever high network bandwidth is available.</p>	<p>Physics discoveries and results are confidential until certified by group and presented at meeting/journal. Data preserved so results reproducible</p>
--	---	---	--	---	---

A.7.5 Use case 40: Belle II Experiment

Use case title	Belle II Experiment
Vertical (area)	Scientific Research: High Energy Physics
Author/company/email	David Asner and Malachi Schram, PNNL, david.asner@pnnl.gov and malachi.schram@pnnl.gov
Actors/stakeholders and their roles and responsibilities	David Asner is the Chief Scientist for the US Belle II Project Malachi Schram is Belle II network and data transfer coordinator and the PNNL Belle II computing center manager
Goals	Perform precision measurements to search for new phenomena beyond the Standard Model of Particle Physics
Use case description	Study numerous decay modes at the Upsilon(4S) resonance to search for new phenomena beyond the Standard Model of Particle Physics

Current solutions	Compute(System)	Distributed (Grid computing using DIRAC)
	Storage	Distributed (various technologies)
	Networking	Continuous RAW data transfer of ~20 Gbps at designed luminosity between Japan and US Additional transfer rates are currently being investigated
	Software	Open Science Grid, Geant4, DIRAC, FTS, Belle II framework
Big data characteristics	Data source (distributed/centralized)	Distributed data centers Primary data centers are in Japan (KEK) and US (PNNL)
	Volume (size)	Total integrated RAW data ~120 PB and physics data ~15 PB and ~100PB MC samples
	Velocity (e.g. real time)	Data will be re-calibrated and analyzed incrementally Data rates will increase based on the accelerator luminosity
	Variety (multiple datasets, mashup)	Data will be re-calibrated and distributed incrementally.
	Variability (rate of change)	Collisions will progressively increase until the designed luminosity is reached (3000 BB pairs per sec). Expected event size is ~300 kB per events.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Validation will be performed using known reference physics processes
	Visualization	N/A
	Data quality (syntax)	Output data will be re-calibrated and validated incrementally
	Data types	Tuple based output
	Data analytics	Data clustering and classification is an integral part of the computing model. Individual scientists define event level analytics.
Big data specific challenges (Gaps)	Data movement and bookkeeping (file and event level meta-data).	
Big data specific challenges in mobility	Network infrastructure required for continuous data transfer between Japan (KEK) and US (PNNL).	

Security and privacy technical considerations	No special challenges. Data is accessed using grid authentication.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	
More information (URLs)	http://belle2.kek.jp

A.8 Earth, Environmental and Polar Science

A.8.1 Use case 41: EISCAT 3D Incoherent Scatter Radar System

Use case title	EISCAT 3D incoherent scatter radar system	
Vertical (area)	Environmental Science	
Author/company/email	Yin Chen /Cardiff University/ chenY58@cardiff.ac.uk Ingemar Häggström, Ingrid Mann, Craig Heinselmann/ EISCAT Science Association/{Ingemar.Haggstrom,%20 Ingrid.mann,%20Craig.Heinselmann%7d@eiscat.se	
Actors/stakeholders and their roles and responsibilities	The EISCAT Scientific Association is an international research organization operating incoherent scatter radar systems in Northern Europe. It is funded and operated by research councils of Norway, Sweden, Finland, Japan, China and the United Kingdom (collectively, the EISCAT Associates). In addition to the incoherent scatter radars, EISCAT also operates an Ionospheric Heater facility, as well as two Dynasondes.	
Goals	EISCAT , the <i>European Incoherent Scatter Scientific Association</i> , is established to conduct research on the lower, middle and upper atmosphere and ionosphere using the incoherent scatter radar technique. This technique is the most powerful ground-based tool for these research applications. EISCAT is also being used as a coherent scatter radar for studying instabilities in the ionosphere, as well as for investigating the structure and dynamics of the middle atmosphere and as a diagnostic instrument in ionospheric modification experiments with the Heating facility.	
Use case description	The design of the next generation incoherent scatter radar system, EISCAT_3D, opens up opportunities for physicists to explore many new research fields. On the other hand, it also introduces significant challenges in handling large-scale experimental data which will be massively generated at great speeds and volumes. This challenge is typically referred to as a big data problem and requires solutions from beyond the capabilities of conventional database technologies.	
Current solutions	Compute(System)	EISCAT 3D data e-Infrastructure plans to use the high performance computers for central site data processing and high throughput computers for mirror sites data processing
	Storage	32 TB

	<p>Networking</p> <p>The estimated data rates in local networks at the active site run from 1 GB/s to 10 GB/s. Similar capacity is needed to connect the sites through dedicated high-speed network links. Downloading the full data is not time critical, but operations require real-time information about certain pre-defined events to be sent from the sites to the operation centre and a real-time link from the operation centre to the sites to set the mode of radar operation on with immediate action.</p>
	<p>Software</p> <ul style="list-style-type: none"> — Mainstream operating systems, e.g., Windows, Linux, Solaris, HP/UX, or FreeBSD — Simple, flat file storage with required capabilities e.g., compression, file striping and file journaling — Self-developed software <ul style="list-style-type: none"> — Control and monitoring tools including, system configuration, quick-look, fault reporting, etc. — Data dissemination utilities — User software e.g., for cyclic buffer, data cleaning, RFI detection and excision, auto-correlation, data integration, data analysis, event identification, discovery and retrieval, calculation of value-added data products, ingestion/extraction, plot — User-oriented computing — APIs into standard software environments — Data processing chains and workflow

Big data characteristics	Data source (distributed/centralized)	EISCAT_3D will consist of a core site with a transmitting and receiving radar arrays and four sites with receiving antenna arrays at some 100 km from the core.
	Volume (size)	<ul style="list-style-type: none"> — The fully operational 5-site system will generate 40 PB/year in 2022. — It is expected to operate for 30 years, and data products to be stored at less 10 years
	Velocity (e.g. real time)	<p>At each of 5-receiver-site:</p> <ul style="list-style-type: none"> — each antenna generates 30 Msamples/s (120 MB/s); — each antenna group (consists of 100 antennas) to form beams at speed of 2 Gbit/s/group; — these data are temporary stored in a ringbuffer: 160 groups ->125 TB/h.
	Variety (multiple datasets, mashup)	<ul style="list-style-type: none"> — Measurements: different versions, formats, replicas, external sources ... — System information: configuration, monitoring, logs/provenance ... — Users' metadata/data: experiments, analysis, sharing, communications ...
	Variability (rate of change)	<p>In time, instantly, a few ms.</p> <p>Along the radar beams, 100ns.</p>

<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<ul style="list-style-type: none"> — Running 24/7, EIS-CAT_3D have very high demands on robustness. — Data and performance assurance is vital for the ring-buffer and archive systems. These systems must be able to guarantee to meet minimum data rate acceptance at all times or scientific data will be lost. — Similarly the systems must guarantee that data held is not volatile or corrupt. This latter requirement is particularly vital at the permanent archive where data is most likely to be accessed by scientific users and least easy to check; data corruption here has a significant possibility of being non-recoverable and of poisoning the scientific literature.
	<p>Visualization</p>	<ul style="list-style-type: none"> — Real-time visualization of analyzed data, e.g., with a figure of updating panels showing electron density, temperatures and ion velocity to those data for each beam. — Non-real-time (post-experiment) visualization of the physical parameters of interest, e.g., <ul style="list-style-type: none"> — by standard plots, — using three-dimensional block to show to spatial variation (in the user selected cuts), — using animations to show the temporal variation,

		<ul style="list-style-type: none"> — allow the visualization of 5 or higher dimensional data, e.g., using the 'cut up and stack' technique to reduce the dimensionality, that is take one or more independent coordinates as discrete; or volume rendering technique to display a 2D projection of a 3D discretely sampled data set. — (Interactive) Visualization. E.g., to allow users to combine the information on several spectral features, e.g., by using color coding, and to provide real-time visualization facility to allow the users to link or plug in tailor-made data visualization functions, and more importantly functions to signal for special observational conditions.
	<p>Data quality (syntax)</p>	<ul style="list-style-type: none"> — Monitoring software will be provided which allows The Operator to see incoming data via the Visualization system in real-time and react appropriately to scientifically interesting events. — Control software will be developed to time-integrate the signals and reduce the noise variance and the total data throughput of the system that reached the data archive.
	<p>Data types</p>	<p>HDF-5</p>
	<p>Data analytics</p>	<p>Pattern recognition, demanding correlation routines, high level parameter extraction</p>
<p>Big data specific challenges (Gaps)</p>	<ul style="list-style-type: none"> — High throughput of data for reduction into higher levels. — Discovery of meaningful insights from low-value-density data needs new approaches to the deep, complex analysis e.g., using machine learning, statistical modelling, graph algorithms etc. which go beyond traditional approaches to the space physics. 	
<p>Big data specific challenges in mobility</p>	<p>Is not likely in mobile platforms</p>	

Security and privacy technical considerations	Lower level of data has restrictions for 1 year within the associate countries. All data open after 3 years.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	EISCAT 3D data e-Infrastructure shares similar architectural characteristics with other ISR radars, and many existing big data systems, such as LOFAR, LHC, and SKA
More information (URLs)	https://www.eiscat3d.se/

A.8.2 Use case 42: Common Environmental Research Infrastructure

Use case title	ENVRI (Common Operations of Environmental Research Infrastructure)
Vertical (area)	Environmental Science
Author/company/email	Yin Chen/ Cardiff University / ChenY58@cardiff.ac.uk
Actors/stakeholders and their roles and responsibilities	<p>The ENVRI project is a collaboration conducted within the European Strategy Forum on Research Infrastructures (ESFRI) Environmental Cluster. The ESFRI Environmental research infrastructures involved in ENVRI including:</p> <ul style="list-style-type: none"> — ICOS is a European distributed infrastructure dedicated to the monitoring of greenhouse gases (GHG) through its atmospheric, ecosystem and ocean networks — EURO-Argo is the European contribution to Argo, which is a global ocean observing system. — EISCAT-3D is a European new-generation incoherent-scatter research radar for upper atmospheric science. — LifeWatch is an e-science Infrastructure for biodiversity and ecosystem research. — EPOS is a European Research Infrastructure on earthquakes, volcanoes, surface dynamics and tectonics. — EMSO is a European network of seafloor observatories for the long-term monitoring of environmental processes related to ecosystems, climate change and geo-hazards. <p>ENVRI also maintains close contact with the other not-directly involved ESFRI Environmental research infrastructures by inviting them for joint meetings. These projects are:</p>
Goals	The ENVRI project gathers 6 EU ESFRI environmental science infra-structures (ICOS, EURO-Argo, EISCAT-3D, LifeWatch, EPOS, and EMSO) in order to develop common data and software services. The results will accelerate the construction of these infrastructures and improve interoperability among them.

	<p>The primary goal of ENVRI is to agree on a reference model for joint operations. The ENVRI RM is a common ontological framework and standard for the description and characterisation of computational and storage infrastructures in order to achieve seamless interoperability between the heterogeneous resources of different infrastructures. The ENVRI RM serves as a common language for community communication, providing a uniform framework into which the infrastructure's components can be classified and compared, also serving to identify common solutions to common problems. This may enable reuse, share of resources and experiences, and avoid duplication of efforts.</p>	
<p>Use case description</p>	<p>ENVRI project implements harmonized solutions and draws up guidelines for the common needs of the environmental ESFRI projects, with a special focus on issues as architectures, metadata frameworks, data discovery in scattered repositories, visualization and data curation. This will empower the users of the collaborating environmental research infrastructures and enable multidisciplinary scientists to access, study and correlate data from multiple domains for "system level" research.</p> <p>ENVRI investigates a collection of representative research infrastructures for environmental sciences, and provides a projection of Europe-wide requirements they have; identifying in particular, requirements they have in common. Based on the analysis evidence, the ENVRI Reference Model (http://www.envri.eu/rm) is developed using ISO standard Open Distributed Processing. Fundamentally the model serves to provide a universal reference framework for discussing many common technical challenges facing all of the ESFRI-environmental research infrastructures. By drawing analogies between the reference components of the model and the actual elements of the infrastructures (or their proposed designs) as they exist now, various gaps and points of overlap can be identified</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	
	<p>Storage</p>	<p>File systems and relational databases</p>
	<p>Networking</p>	
	<p>Software</p>	<p>Own</p>

<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>Most of the ENVRI Research Infrastructures (ENV RIs) are <i>distributed, long-term, remote controlled observational networks</i> focused on understanding processes, trends, thresholds, interactions and feedbacks and increasing the predictive power to address future environmental challenges. They are spanning from the Arctic areas to the European Southernmost areas and from Atlantic on west to the Black Sea on east. More precisely:</p> <ul style="list-style-type: none"> — <i>EMSO</i>, network of fixed-point, deep-seafloor and water column observatories, is geographically distributed in key sites of European waters, presently consisting of thirteen sites. — <i>EPOS</i> aims at integrating the existing European facilities in solid Earth science into one coherent multidisciplinary RI, and to increase the accessibility and usability of multidisciplinary data from seismic and geodetic monitoring networks, volcano observatories, laboratory experiments and computational simulations enhancing worldwide interoperability in Earth Science.
--	---	--

Copyrighted document, no reproduction or circulation
 STANDARDSISO.COM: Click to view the full PDF of ISO/IEC TR 20547-2:2018
 Oct 2024

	<p>— <i>ICOS</i> dedicates to the monitoring of greenhouse gases (GHG) through its atmospheric, ecosystem and ocean networks. The <i>ICOS</i> network includes more than 30 atmospheric and more than 30 ecosystem primary long term sites located across Europe, and additional secondary sites. It also includes three Thematic Centres to process the data from all the stations from each network, and provide access to these data.</p> <p>— <i>LifeWatch</i> is a “virtual” infrastructure for biodiversity and ecosystem research with services mainly provided through the Internet. Its Common Facilities is coordinated and managed at a central European level; and the <i>LifeWatch Centres</i> serve as specialized facilities from member countries (regional partner facilities) or research communities.</p> <p>— <i>Euro-Argo</i> provides, deploys and operates an array of around 800 floats contributing to the global array (3,000 floats) and thus provide enhanced coverage in the European regional seas.</p> <p>— <i>EISCAT-3D</i>, makes continuous measurements of the geospace environment and its coupling to the Earth's atmosphere from its location in the auroral zone at the southern edge of the northern polar vortex, and is a distributed infrastructure.</p>
--	---

	Volume (size)	<p>Variable data size. e.g.,</p> <ul style="list-style-type: none"> — The amount of data within the <i>EMSO</i> is depending on the instrumentation and configuration of the observatory between several MBs to several GB per data set. — Within <i>EPOS</i>, the EIDA network is currently providing access to continuous raw data coming from approximately more than 1 000 stations recording about 40 GB per day, so over 15 TB per year. EMSC stores a Database of 1,85 GB of earthquake parameters, which is constantly growing and updated with refined information. — 222 705 - events — 632 327 - origins — 642 555 - magnitudes — Within <i>EISCAT 3D</i> raw voltage data will reach 40PB/year in 2023.
	Velocity (e.g. real time)	Real-time data handling is a common request of the environmental research infrastructures
	Variety (multiple datasets, mashup)	Highly complex and heterogeneous
	Variability (rate of change)	Relative low rate of change
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Normal

	<p>Visualization</p> <p>Most of the projects have not yet developed the visualization technique to be fully operational.</p> <ul style="list-style-type: none"> — <i>EMSO</i> is not yet fully operational, currently only simple graph plotting tools. — Visualization techniques are not yet defined for <i>EPOS</i>. — Within <i>ICOS</i> Level-1.b data products such as near real time GHG measurements are available to users via ATC web portal. Based on Google Chart Tools, an interactive time series line chart with optional annotations allows user to scroll and zoom inside a time series of CO₂ or CH₄ measurement at an ICOS Atmospheric station. The chart is rendered within the browser using Flash. Some Level-2 products are also available to ensure instrument monitoring to PIs. It is mainly instrumental and comparison data plots automatically generated (R language and Python Matplotlib 2D plotting library) and daily pushed on ICOS web server. Level-3 data products such as gridded GHG fluxes derived from ICOS observations increase the scientific impact of ICOS.
--	---

	<p>For this purpose ICOS supports its community of users. The Carbon portal is expected to act as a platform that will offer visualization of the flux products that incorporate ICOS data. Example of candidate Level-3 products from future ICOS GHG concentration data are for instance maps of European high-resolution CO₂ or CH₄ fluxes obtained by atmospheric inversion modellers in Europe. Visual tools for comparisons between products will be developed by the Carbon Portal. Contributions will be open to any product of high scientific quality.</p> <p>— <i>LifeWatch</i> will provide common visualization techniques, such as the plotting of species on maps. New techniques will allow visualizing the effect of changing data and/or parameters in models</p>	
	Data quality (syntax)	Highly important
	Data types	<ul style="list-style-type: none"> — Measurements (often in file formats), — Metadata, — Ontology, — Annotations
	Data analytics	<ul style="list-style-type: none"> — Data assimilation — (Statistical) analysis, — Data mining, — Data extraction, — Scientific modeling and simulation, — Scientific workflow
Big data specific challenges (Gaps)	<ul style="list-style-type: none"> — Real-time handling of extreme high volume of data — Data staging to mirror arc Hives — Integrated Data access and discovery — Data processing and analysis 	

<p>Big data specific challenges in mobility</p>	<p>The need for efficient and high performance mobile detectors and instrumentation is common:</p> <ul style="list-style-type: none"> — In ICOS, various mobile instruments are used to collect data from marine observations, atmospheric observations, and ecosystem monitoring. — In Euro-Argo, thousands of submersible robots to obtain observations of all of the oceans — In Lifewatch, biologists use mobile instruments for observations and measurements.
<p>Security and privacy technical considerations</p>	<p>Most of the projects follow the open data sharing policy. E.g.,</p> <ul style="list-style-type: none"> — The vision of EMSO is to allow scientists all over the world to access observatories data following an open access model. — Within EPOS, EIDA data and Earthquake parameters are generally open and free to use. Few restrictions are applied on few seismic networks and the access is regulated depending on email based authentication/authorization. — The ICOS data will be accessible through a license with full and open access. No particular restriction in the access and eventual use of the data is anticipated, expected the inability to redistribute the data. Acknowledgement of ICOS and traceability of the data will be sought in a specific, way (e.g. DOI of dataset). A large part of relevant data and resources are generated using public funding from national and international sources. — LifeWatch is following the appropriate European policies, such as: the European Research Council (ERC) requirement; the European Commission's open access pilot mandate in 2008. For publications, initiatives such as Dryad instigated by publishers and the Open Access Infrastructure for Research in Europe (OpenAIRE). The private sector may deploy their data in the LifeWatch infrastructure. A special company will be established to manage such commercial contracts. — In EISCAT 3D, lower level of data has restrictions for 1 year within the associate countries. All data open after 3 years.

<p>Highlight issues for generalizing this Use case (e.g. for ref. architecture)</p>	<p>Different research infrastructures are designed for different purposes and evolve over time. The designers describe their approaches from different points of view, in different levels of detail and using different typologies. The documentation provided is often incomplete and inconsistent. What is needed is a uniform platform for interpretation and discussion, which helps to unify understanding</p> <p>In ENVRI, we choose to use a standard model, Open Distributed Processing (ODP), to interpret the design of the research infrastructures, and place their requirements into the ODP framework for further analysis and comparison.</p>
<p>More information (URLs)</p>	<ul style="list-style-type: none"> — ENVRI Project website: http://www.envri.eu — ENVRI Reference Model http://www.envri.eu/rm — ENVRI deliverable D3.2: Analysis of common requirements of Environmental Research Infrastructures — ICOS: http://www.icos-infrastructure.eu/ — Euro-Argo: http://www.euro-argo.eu/ — EISCAT 3D: http://www.eiscat3d.se/ — LifeWatch: http://www.lifewatch.com/ — EPOS: http://www.epos-eu.org/ — EMSO http://www.emso-eu.org/management/

A.8.3 Use case 43: Radar Data Analysis for CReSIS

Use case title	Radar Data Analysis for CReSIS	
Vertical (area)	Scientific Research: Polar Science and Remote Sensing of Ice Sheets	
Author/company/email	Geoffrey Fox, Indiana University gcf@indiana.edu	
Actors/stakeholders and their roles and responsibilities	Research funded by NSF and NASA with relevance to near and long term climate change. Engineers designing novel radar with “field expeditions” for 1 to 2 months to remote sites. Results used by scientists building models and theories involving Ice Sheets	
Goals	Determine the depths of glaciers and snow layers to be fed into higher level scientific analyses	
Use case description	Build radar; build UAV or use piloted aircraft; overfly remote sites (Arctic, Antarctic, Himalayas). Check in field that experiments configured correctly with detailed analysis later. Transport data by air-shipping disk as poor Internet connection. Use image processing to find ice/snow sheet depths. Use depths in scientific discovery of melting ice caps etc	
Current solutions	Compute(System)	Field is a low power cluster of rugged laptops plus classic 2 to 4 CPU servers with ~40 TB removable disk array. Offline is about 2 500 cores
	Storage	Removable disk in field. (Disks suffer in field so 2 copies made) Lustre or equivalent for offline
	Networking	Terrible Internet linking field sites to continental USA.
	Software	Radar signal processing in Matlab. Image analysis is Map/Reduce or MPI plus C/Java. User Interface is a Geographical Information System
Big data characteristics	Data source (distributed/centralized)	Aircraft flying over ice sheets in carefully planned paths with data downloaded to disks.
	Volume (size)	~0,5 Petabytes per year raw data
	Velocity (e.g. real time)	All data gathered in real time but analyzed incrementally and stored with a GIS interface
	Variety (multiple datasets, mashup)	Lots of different datasets – each needing custom signal processing but all similar in structure. This data needs to be used with wide variety of other polar data.

	Variability (rate of change)	Data accumulated in ~100 TB chunks for each expedition
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Essential to monitor field data and correct instrumental problems. Implies must analyze fully portion of data in field
	Visualization	Rich user interface for layers and glacier simulations
	Data quality (syntax)	Main engineering issue is to ensure instrument gives quality data
	Data types	Radar Images
	Data analytics	Sophisticated signal processing; novel new image processing to find layers (can be 100's one per year)
Big data specific challenges (Gaps)	Data volumes increasing. Shipping disks clumsy but no other obvious solution. Image processing algorithms still very active research	
Big data specific challenges in mobility	Smart phone interfaces not essential but LOW power technology essential in field	
Security and privacy technical considerations	Himalaya studies fraught with political issues and require UAV. Data itself open after initial study	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Loosely coupled clusters for signal processing. Must support Matlab.	
More information (URLs)	http://polargrid.org/polargrid https://www.cresis.ku.edu/ See movie at http://polargrid.org/polargrid/gallery	
NOTE		

Use case Stages	Data Sources	Data Usage	Transformations (Data Analytics)	Infrastructure	Security and Privacy
Radar Data Analysis for CReSIS (Scientific Research: Polar Science and Remote Sensing of Ice Sheets)					
Raw Data: Field Trip	Raw Data from Radar instrument on Plane/Vehicle	Capture Data on Disks for L1B. Check Data to monitor instruments.	Robust Data Copying Utilities. Version of Full Analysis to check data.	Rugged Laptops with small server (~2 CPU with ~40TB removable disk system)	N/A
Information: Offline Analysis L1B	Transported Disks copied to (LUSTRE) File System	Produce processed data as radar images	Matlab Analysis code running in parallel and independently on each data sample	~2500 cores running standard cluster tools	N/A except results checked before release on CReSIS web site
Information: L2/L3 Geolocation and Layer Finding	Radar Images from L1B	Input to Science as database with GIS frontend	GIS and Metadata Tools Environment to support automatic and/or manual layer determination	GIS (Geographical Information System). Cluster for Image Processing.	As above
Knowledge, Wisdom, Discovery: Science	GIS interface to L2/L3 data	Polar Science Research integrating multiple data sources e.g. for Climate change. Glacier bed data used in simulations of glacier flow.		Exploration on a cloud style GIS supporting access to data. Simulation is 3D partial differential equation solver on large cluster.	Varies according to science use. Typically results open after research complete.

A.8.4 Use case 44: UAVSAR Data Processing

Use case title	UAVSAR Data Processing, Data Product Delivery, and Data Services	
Vertical (area)	Scientific Research: Earth Science	
Author/company/email	Andrea Donnellan, NASA JPL, andrea.donnellan@jpl.nasa.gov; Jay Parker, NASA JPL, jay.w.parker@jpl.nasa.gov	
Actors/stakeholders and their roles and responsibilities	NASA UAVSAR team, NASA QuakeSim team, ASF (NASA SAR DAAC), USGS, CA Geological Survey	
Goals	Use of Synthetic Aperture Radar (SAR) to identify landscape changes caused by seismic activity, landslides, deforestation, vegetation changes, flooding, etc.; increase its usability and accessibility by scientists.	
Use case description	A scientist who wants to study the after effects of an earthquake examines multiple standard SAR products made available by NASA. The scientist may find it useful to interact with services provided by intermediate projects that add value to the official data product archive.	
Current solutions	Compute(System)	Raw data processing at NASA AMES Pleiades, Endeavour. Commercial clouds for storage and service front ends have been explored.
	Storage	File based.
	Networking	Data require one time transfers between instrument and JPL, JPL and other NASA computing centers (AMES), and JPL and ASF. Individual data files are not too large for individual users to download, but entire data set is unwieldy to transfer. This is a problem to downstream groups like QuakeSim who want to reformat and add value to data sets.
	Software	ROI_PAC, GeoServer, GDAL, GeoTIFF-supporting tools.
Big data characteristics	Data source (distributed/centralized)	Data initially acquired by unmanned aircraft. Initially processed at NASA JPL. Archive is centralized at ASF (NASA DAAC). QuakeSim team maintains separate downstream products (GeoTIFF conversions).

	<p>Volume (size)</p>	<p>Repeat Pass Interferometry (RPI) Data: ~ 3 TB. Increasing about 1 TB to 2 TB/year. Polarimetric Data: ~40 TB (processed) Raw Data: 110 TB Proposed satellite missions (Earth Radar Mission, formerly DESDynI) could dramatically increase data volumes (TBs per day).</p>
	<p>Velocity (e.g. real time)</p>	<p>RPI Data: 1 to 2 TB/year. Polarimetric data is faster.</p>
	<p>Variety (multiple datasets, mashup)</p>	<p>Two main types: Polarimetric and RPI. Each RPI product is a collection of files (annotation file, unwrapped, etc.). Polarimetric products also consist of several files each.</p>
	<p>Variability (rate of change)</p>	<p>Data products change slowly. Data occasionally get re-processed: new processing methods or parameters. There may be additional quality assurance and quality control issues.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Provenance issues need to be considered. This provenance has not been transparent to downstream consumers in the past. Versioning used now; versions described in the UAVSAR web page in notes.</p>
	<p>Visualization</p>	<p>Uses Geospatial Information System tools, services, standards.</p>
	<p>Data quality (syntax)</p>	<p>Many frames and collections are found to be unusable due to unforeseen flight conditions.</p>
	<p>Data types</p>	<p>GeoTIFF and related imagery data</p>
	<p>Data analytics</p>	<p>Done by downstream consumers (such as edge detections): research issues.</p>
<p>Big data specific challenges (Gaps)</p>	<p>Data processing pipeline requires human inspection and intervention. Limited downstream data pipelines for custom users. Cloud architectures for distributing entire data product collections to downstream consumers should be investigated, adopted.</p>	

Big data specific challenges in mobility	Some users examine data in the field on mobile devices, requiring interactive reduction of large data sets to understandable images or statistics.
Security and privacy technical considerations	Data is made immediately public after processing (no embargo period).
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Data is geolocated, and may be angularly specified. Categories: GIS; standard instrument data processing pipeline to produce standard data products.
More information (URLs)	http://uavsar.jpl.nasa.gov/ , http://www.asf.alaska.edu/program/sdc , http://quakesim.org

A.8.5 Use case 45: NASA LARC/GSFC iRODS Federation Testbed

Use case title	NASA LARC/GSFC iRODS Federation Testbed
Vertical (area)	Earth Science Research and Applications
Author/company/email	Michael Little, Roger Dubois, Brandi Quam, Tiffany Mathews, Andrei Vakhnin, Beth Huffer, Christian Johnson / NASA Langley Research Center (LaRC) / M.M.Little@NASA.gov, Roger.A.Dubois@nasa.gov, Brandi.M.Quam@NASA.gov, Tiffany.J.Mathews@NASA.gov, and Andrei.A.Vakhnin@NASA.gov
	John Schnase, Daniel Duffy, Glenn Tamkin, Scott Sinno, John Thompson, and Mark McInerney / NASA Goddard Space Flight Center (GSFC) / John.L.Schnase@NASA.gov, Daniel.Q.Duffy@NASA.gov, Glenn.S.Tamkin@nasa.gov, Scott.S.Sinno@nasa.gov, John.H.Thompson@nasa.gov, and Mark.McInerney@nasa.gov
Actors/stakeholders and their roles and responsibilities	NASA's Atmospheric Science Data Center (ASDC) at Langley Research Center (LaRC) in Hampton, Virginia, and the Center for Climate Simulation (NCCS) at Goddard Space Flight Center (GSFC) both ingest, archive, and distribute data that is essential to stakeholders including the climate research community, science applications community, and a growing community of government and private-sector customers who have a need for atmospheric and climatic data.
Goals	To implement a data federation ability to improve and automate the discovery of heterogeneous data, decrease data transfer latency, and meet customizable criteria based on data content, data quality, metadata, and production. To support/enable applications and customers that require the integration of multiple heterogeneous data collections.

<p>Use case description</p>	<p>ASDC and NCCS have complementary data sets, each containing vast amounts of data that is not easily shared and queried. Climate researchers, weather forecasters, instrument teams, and other scientists need to access data from across multiple datasets in order to compare sensor measurements from various instruments, compare sensor measurements to model outputs, calibrate instruments, look for correlations across multiple parameters, etc. To analyze, visualize and otherwise process data from heterogeneous datasets is currently a time consuming effort that requires scientists to separately access, search for, and download data from multiple servers and often the data is duplicated without an understanding of the authoritative source. Many scientists report spending more time in accessing data than in conducting research. Data consumers need mechanisms for retrieving heterogeneous data from a single point-of-access. This can be enabled through the use of iRODS, a Data grid software system that enables parallel downloads of datasets from selected replica servers that can be geographically dispersed, but still accessible by users worldwide. Using iRODS in conjunction with semantically enhanced metadata, managed via a highly precise Earth Science ontology, the ASDC's Data Products Online (DPO) will be federated with the data at the NASA Center for Climate Simulation (NCCS) at Goddard Space Flight Center (GSFC). The heterogeneous data products at these two NASA facilities are being semantically annotated using common concepts from the NASA Earth Science ontology. The semantic annotations will enable the iRODS system to identify complementary datasets and aggregate data from these disparate sources, facilitating data sharing between climate modelers, forecasters, Earth scientists, and scientists from other disciplines that need Earth science data. The iRODS data federation system will also support cloud-based data processing services in the Amazon Web Services (AWS) cloud.</p>	
<p>Current solutions</p>	<p>Compute(System)</p>	<p>NASA Center for Climate Simulation (NCCS) and NASA Atmospheric Science Data Center (ASDC): Two GPFS systems</p>

	<p>Storage</p> <p>The ASDC's Data Products Online (DPO) GPFS File system consists of 12 x IBM DC4800 and 6 x IBM DCS3700 Storage subsystems, 144 Intel 2,4 GHz cores, 1,400 TB usable storage. NCCS data is stored in the NCCS MERRA cluster, which is a 36 node Dell cluster, 576 Intel 2,6 GHz SandyBridge cores, 1,300 TB raw storage, 1,250 GB RAM, 11,7 TF theoretical peak compute capacity.</p>
	<p>Networking</p> <p>A combination of Fibre Channel SAN and 10GB LAN. The NCCS cluster nodes are connected by an FDR Infiniband network with peak TCP/IP speeds >20 Gbps.</p>
	<p>Software</p> <p>SGE Univa Grid Engine Version 8.1, iRODS version 3.2 and/or 3.3, IBM General Parallel File System (GPFS) version 3.4, Cloudera version 4.5.2-1</p>
<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p> <p>iRODS will be leveraged to share data collected from CERES Level 3B data products including: CERES EBAF-TOA and CERES-Surface products.</p> <p>Surface fluxes in EBAF-Surface are derived from two CERES data products: 1) CERES SYN1deg-Month Ed3 - which provides computed surface fluxes to be adjusted and 2) CERES EBAF-TOA Ed2.7 - which uses observations to provide CERES-derived TOA flux constraints. Access to these products will enable the NCCS at GSFC to run data from the products in a simulation model in order to produce an assimilated flux.</p>

	<p>The NCCS will introduce Modern-Era Retrospective Analysis for Research and Applications (MERRA) data to the iRODS federation. MERRA integrates observational data with numerical models to produce a global temporally and spatially consistent synthesis of 26 key climate variables. MERRA data files are created from the Goddard Earth Observing System version 5 (GEOS-5) model and are stored in HDF-EOS and (Network Common Data Form) NetCDF formats</p> <p>Spatial resolution is 1/2° latitude × 2/3° longitude.</p> <p>Each file contains a single grid with multiple 2D and 3D variables. All data are stored on a longitude-latitude grid with a vertical dimension applicable for all 3D variables. The GEOS-5 MERRA products are divided into 25 collections: 18 standard products, chemistry products. The collections comprise monthly means files and daily files at six-hour intervals running from 1979 to 2012. MERRA data are typically packaged as multi-dimensional binary data within a self-describing NetCDF file format. Hierarchical metadata in the NetCDF header contain the representation information that allows NetCDF-aware software to work with the data. It also contains arbitrary preservation description and policy information that can be used to bring the data into use-specific compliance.</p>
--	---

	<p>Volume (size)</p>	<p>Currently, Data from the EBAF-TOA Product is about 420 MB and Data from the EBAF-Surface Product is about 690 MB. Data grows with each version update (about every six months). The MERRA collection represents about 160 TB of total data (uncompressed); compressed is ~80 TB.</p>
	<p>Velocity (e.g. real time)</p>	<p>Periodic since updates are performed with each new version update.</p>
	<p>Variety (multiple datasets, mashup)</p>	<p>There is a need in many types of applications to combine MERRA reanalysis data with other reanalyses and observational data such as CERES. The NCCS is using the Climate Model Intercomparison Project (CMIP5) Reference standard for ontological alignment across multiple, disparate data sets.</p>
	<p>Variability (rate of change)</p>	<p>The MERRA reanalysis grows by approximately one TB per month.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Validation and testing of semantic metadata, and of federated data products will be provided by data producers at NASA Langley Research Center and at Goddard through regular testing. Regression testing will be implemented to ensure that updates and changes to the iRODS system, newly added data sources, or newly added metadata do not introduce errors to federated data products. MERRA validation is provided by the data producers, NASA Goddard's Global Modeling and Assimilation Office (GMAO).</p>

	<p>Visualization</p>	<p>There is a growing need in the scientific community for data management and visualization services that can aggregate data from multiple sources and display it in a single graphical display. Currently, such capabilities are hindered by the challenge of finding and downloading comparable data from multiple servers, and then transforming each heterogeneous dataset to make it usable by the visualization software. Federation of NASA datasets using iRODS will enable scientists to quickly find and aggregate comparable datasets for use with visualization software.</p>
	<p>Data quality (syntax)</p>	<p>For MERRA, quality controls are applied by the data producers, GMAO.</p>
	<p>Data types</p>	<p>See above</p>
	<p>Data analytics</p>	<p>Pursuant to the first goal of increasing accessibility and discoverability through innovative technologies, the ASDC and NCCS are exploring a capability to improve data access capabilities. Using iRODS, the ASDC's Data Products Online (DPO) can be federated with data at GSFC's NCCS creating a data access system that can serve a much broader customer base than is currently being served. Federating and sharing information will enable the ASDC and NCCS to fully utilize multi-year and multi-instrument data and will improve and automate the discovery of heterogeneous data, increase data transfer latency, and meet customizable criteria based on data content, data quality, metadata, and production.</p>
<p>Big data specific challenges (Gaps)</p>		

Big data specific challenges in mobility	A major challenge includes defining an enterprise architecture that can deliver real-time analytics via communication with multiple APIs and cloud computing systems. By keeping the computation resources on cloud systems, the challenge with mobility resides in not overpowering mobile devices with displaying CPU intensive visualizations that may hinder the performance or usability of the data being presented to the user.
Security and privacy technical considerations	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	<p>This federation builds on several years of iRODS research and development performed at the NCCS. During this time, the NCCS vetted the iRODS features while extending its core functions with domain-specific extensions. For example, the NCCS created and installed Python-based scientific kits within iRODS that automatically harvest metadata when the associated data collection is registered. One of these scientific kits was developed for the MERRA collection. This kit in conjunction with iRODS bolsters the strength of the LaRC/GSFC federation by providing advanced search capabilities. LaRC is working through the establishment of an advanced architecture that leverages multiple technology pilots and tools (access, discovery, and analysis) designed to integrate capabilities across the earth science community – the research and development completed by both data centers is complementary and only further enhances this use case.</p> <p>Other scientific kits that have been developed include: NetCDF, Intergovernmental Panel on Climate Change (IPCC), and Ocean Modeling and Data Assimilation (ODAS). The combination of iRODS and these scientific kits has culminated in a configurable technology stack called the virtual Climate Data Server (vCDS), meaning that this runtime environment can be deployed to multiple destinations (e.g., bare metal, virtual servers, cloud) to support various scientific needs. The vCDS, which can be viewed as a reference architecture for easing the federation of disparate data repositories, is leveraged by but not limited to LaRC and GSFC.</p>
More information (URLs)	Please contact the authors for additional information

A.8.6 Use case 46: MERRA Analytic Services

Use case title	MERRA Analytic Services (MERRA/AS)	
Vertical (area)	Scientific Research: Earth Science	
Author/company/email	John L. Schnase and Daniel Q. Duffy / NASA Goddard Space Flight Center John.L.Schnase@NASA.gov, Daniel.Q.Duffy@NASA.gov	
Actors/stakeholders and their roles and responsibilities	NASA's Modern-Era Retrospective Analysis for Research and Applications (MERRA) integrates observational data with numerical models to produce a global temporally and spatially consistent synthesis of 26 key climate variables. Actors and stakeholders who have an interest in MERRA include the climate research community, science applications community, and a growing number of government and private-sector customers who have a need for the MERRA data in their decision support systems.	
Goals	Increase the usability and use of large-scale scientific data collections, such as MERRA.	
Use case description	MERRA Analytic Services enables Map/Reduce analytics over the MERRA collection. MERRA/AS is an example of cloud-enabled climate analytics as a service (CAaaS), which is an approach to meeting the big data challenges of climate science through the combined use of 1) high performance, data proximal analytics, (2) scalable data management, (3) software appliance virtualization, (4) adaptive analytics, and (5) a domain-harmonized API. The effectiveness of MERRA/AS is being demonstrated in several applications, including data publication to the Earth System Grid Federation (ESGF) in support of Intergovernmental Panel on Climate Change (IPCC) research, the NASA/Department of Interior RECOVER wild land fire decision support system, and data interoperability testbed evaluations between NASA Goddard Space Flight Center and the NASA Langley Atmospheric Data Center.	
Current solutions	Compute(System)	NASA Center for Climate Simulation (NCCS)
	Storage	The MERRA Analytic Services Hadoop Filesystem (HDFS) is a 36 node Dell cluster, 576 Intel 2,6 GHz SandyBridge cores, 1300 TB raw storage, 1250 GB RAM, 11,7 TF theoretical peak compute capacity.
	Networking	Cluster nodes are connected by an FDR Infiniband network with peak TCP/IP speeds >20 Gbps.
	Software	Cloudera, iRODS, Amazon AWS

<p>Big data characteristics</p>	<p>Data source (distributed/centralized)</p>	<p>MERRA data files are created from the Goddard Earth Observing System version 5 (GEOS-5) model and are stored in HDF-EOS and NetCDF formats. Spatial resolution is 1/2 °latitude × 2/3 °longitude × 72 vertical levels extending through the stratosphere. Temporal resolution is 6-hours for three-dimensional, full spatial resolution, extending from 1979-present, nearly the entire satellite era. Each file contains a single grid with multiple 2D and 3D variables. All data are stored on a longitude latitude grid with a vertical dimension applicable for all 3D variables. The GEOS-5 MERRA products are divided into 25 collections: 18 standard products, 7 chemistry products. The collections comprise monthly means files and daily files at six-hour intervals running from 1979 -2012. MERRA data are typically packaged as multi-dimensional binary data within a self-describing NetCDF file format. Hierarchical metadata in the NetCDF header contain the representation information that allows NetCDF aware software to work with the data. It also contains arbitrary preservation description and policy information that can be used to bring the data into use-specific compliance.</p>
	<p>Volume (size)</p>	<p>480TB</p>
	<p>Velocity (e.g. real time)</p>	<p>Real-time or batch, depending on the analysis. We're developing a set of "canonical ops" -early stage, near-data operations common to many analytic workflows. The goal is for the canonical ops to run in near real-time.</p>

	<p>Variety (multiple datasets, mashup)</p>	<p>There is a need in many types of applications to combine MERRA reanalysis data with other re-analyses and observational data. We are using the Climate Model Inter-comparison Project (CMIP5) Reference standard for ontological alignment across multiple, disparate data sets.</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Variability (rate of change)</p>	<p>The MERRA reanalysis grows by approximately one TB per month.</p>
	<p>Veracity (Robustness Issues, semantics)</p>	<p>Validation provided by data producers, NASA Goddard's Global Modeling and Assimilation Office (GMAO).</p>
<p>Copyrighted document, no reproduction or reproduction for any purpose without the prior written permission of ISO/IEC. For review by ISO/IEC JTC1/SC32, Oct 2024</p>	<p>Visualization</p>	<p>There is a growing need for distributed visualization of analytic outputs.</p>
	<p>Data quality (syntax)</p>	<p>Quality controls applied by data producers, GMAO.</p>
	<p>Data types</p>	<p>See above.</p>
	<p>Data analytics</p>	<p>In our efforts to address the big data challenges of climate science, we are moving toward a notion of climate analytics-as-a-service. We focus on analytics, because it is the knowledge gained from our interactions with big data that ultimately produce societal benefits. We focus on CAaaS because we believe it provides a useful way of thinking about the problem: a specialization of the concept of business process-as-a-service, which is an evolving extension of IaaS, PaaS, and SaaS enabled by Cloud Computing.</p>

<p>Big data specific challenges (Gaps)</p>	<p>A big question is how to use cloud computing to enable better use of climate science's earthbound compute and data resources. Cloud Computing is providing for us a new tier in the data services stack—a cloud-based layer where agile customization occurs and enterprise-level products are transformed to meet the specialized requirements of applications and consumers. It helps us close the gap between the world of traditional, high-performance computing, which, at least for now, resides in a finely-tuned climate modeling environment at the enterprise level and our new customers, whose expectations and manner of work are increasingly influenced by the smart mobility megatrend.</p>
<p>Big data specific challenges in mobility</p>	<p>Most modern smartphones, tablets, etc. actually consist of just the display and user interface components of sophisticated applications that run in cloud data centers. This is a mode of work that CAaaS is intended to accommodate.</p>
<p>Security and privacy technical considerations</p>	<p>No critical issues identified at this time.</p>
<p>Highlight issues for generalizing this Use case (e.g. for ref. architecture)</p>	<p>Map/Reduce and iRODS fundamentally make analytics and data aggregation easier; our approach to software appliance virtualization in makes it easier to transfer capabilities to new users and simplifies their ability to build new applications; the social construction of extended capabilities facilitated by the notion of canonical operations enable adaptability; and the Climate Data Services API that we're developing enables ease of mastery. Taken together, we believe that these core technologies behind CAaaS creates a generative context where inputs from diverse people and groups, who may or may not be working in concert, can contribute capabilities that help address the big data challenges of climate science.</p>
<p>More information (URLs)</p>	<p>Please contact the authors for additional information.</p>

A.8.7 Use case 47: Atmospheric Turbulence—Event Discovery

Use case title	Atmospheric Turbulence - Event Discovery and Predictive Analytics	
Vertical (area)	Scientific Research: Earth Science	
Author/company/email	Michael Seablom, NASA Headquarters, michael.s.seablom@nasa.gov	
Actors/stakeholders and their roles and responsibilities	Researchers with NASA or NSF grants, weather forecasters, aviation interests (for the generalized case, any researcher who has a role in studying phenomena-based events).	
Goals	Enable the discovery of high-impact phenomena contained within voluminous Earth Science data stores and which are difficult to characterize using traditional numerical methods (e.g., turbulence). Correlate such phenomena with global atmospheric re-analysis products to enhance predictive capabilities.	
Use case description	Correlate aircraft reports of turbulence (either from pilot reports or from automated aircraft measurements of eddy dissipation rates) with recently completed atmospheric re-analyses of the entire satellite-observing era. Reanalysis products include the North American Regional Reanalysis (NARR) and the Modern-Era Retrospective-Analysis for Research (MERRA) from NASA.	
Current solutions	Compute(System)	NASA Earth Exchange (NEX) - Pleiades supercomputer.
	Storage	Re-analysis products are on the order of 100 TB each; turbulence data are negligible in size.
	Networking	Re-analysis datasets are likely to be too large to relocate to the supercomputer of choice (in this case NEX), therefore the fastest networking possible would be needed.
	Software	Map/Reduce or the like; SciDB or other scientific database.
Big data characteristics	Data source (distributed/centralized)	Distributed
	Volume (size)	200 TB (current), 500 TB within 5 years
	Velocity (e.g. real time)	Data analyzed incrementally
	Variety (multiple datasets, mashup)	Re-analysis datasets are inconsistent in format, resolution, semantics, and metadata. Likely each of these input streams will have to be interpreted/analyzed into a common product.
	Variability (rate of change)	Turbulence observations would be updated continuously; re-analysis products are released about once every five years.
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Validation would be necessary for the output product (correlations).

	Visualization	Useful for interpretation of results.
	Data quality (syntax)	Input streams would have already been subject to quality control.
	Data types	Gridded output from atmospheric data assimilation systems and textual data from turbulence observations.
	Data analytics	Event-specification language needed to perform data mining / event searches.
Big data specific challenges (Gaps)	Semantics (interpretation of multiple reanalysis products); data movement; database(s) with optimal structuring for 4-dimensional data mining.	
	Development for mobile platforms not essential at this time.	
Security and privacy technical considerations	No critical issues identified.	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Atmospheric turbulence is only one of many phenomena-based events that could be useful for understanding anomalies in the atmosphere or the ocean that are connected over long distances in space and time. However the process has limits to extensibility, i.e., each phenomena may require very different processes for data mining and predictive analysis.	
More information (URLs)	http://oceanworld.tamu.edu/resources/oceanography-book/teleconnections.htm http://www.forbes.com/sites/toddwoody/2012/03/21/meet-the-scientists-mining-big-data-to-predict-the-weather/	

A.8.8 Use case 48: Climate Studies using the Community Earth System Model

Use case title	Climate Studies using the Community Earth System Model at DOE's NERSC center	
Vertical (area)	Research: Climate	
Author/company/email	PI: Warren Washington, NCAR	
Actors/stakeholders and their roles and responsibilities	Climate scientists, U.S. policy makers	
Goals	The goals of the Climate Change Prediction (CCP) group at NCAR are to understand and quantify contributions of natural and anthropogenic-induced patterns of climate variability and change in the 20th and 21st centuries by means of simulations with the Community Earth System Model (CESM).	
Use case description	With these model simulations, researchers are able to investigate mechanisms of climate variability and change, as well as to detect and attribute past climate changes, and to project and predict future changes. The simulations are motivated by broad community interest and are widely used by the national and international research communities.	
Current solutions	Compute(System)	NERSC (24M Hours), DOE LCF (41M), NCAR CSL (17M)
	Storage	1,5 PB at NERSC
	Networking	ESNet
	Software	NCAR PIO library and utilities NCL and NCO, parallel NetCDF
Big data characteristics	Data source (distributed/centralized)	Data is produced at computing centers. The Earth Systems Grid is an open source effort providing a robust, distributed data and computation platform, enabling world wide access to Peta/Exa-scale scientific data. ESGF manages the first-ever decentralized database for handling climate science data, with multiple petabytes of data at dozens of federated sites worldwide. It is recognized as the leading infrastructure for the management and access of large distributed data volumes for climate change research. It supports the Coupled Model Intercomparison Project (CMIP), whose protocols enable the periodic assessments carried out by the Intergovernmental Panel on Climate Change (IPCC).

	Volume (size)	30 PB at NERSC (assuming 15 end-to-end climate change experiments) in 2017; many times more worldwide
	Velocity (e.g. real time)	42 GBytes/sec are produced by the simulations
	Variety (multiple datasets, mashup)	Data must be compared among those from observations, historical reanalysis, and a number of independently produced simulations. The Program for Climate Model Diagnosis and Intercomparison develops methods and tools for the diagnosis and intercomparison of general circulation models (GCMs) that simulate the global climate. The need for innovative analysis of GCM climate simulations is apparent, as increasingly more complex models are developed, while the disagreements among these simulations and relative to climate observations remain significant and poorly understood. The nature and causes of these disagreements must be accounted for in a systematic fashion in order to confidently use GCMs for simulation of putative global climate change.
	Variability (rate of change)	Data is produced by codes running at supercomputer centers. During runtime, intense periods of data i/o occur regularly, but typically consume only a few percent of the total run time. Runs are carried out routinely, but spike as deadlines for reports approach.

Copyrighted document, no reproduction or circulation allowed. For review by ISO/IEC JTC1/SC29/WG2. Oct 2024

STANDARDSISO.COM: Click to view the full PDF of ISO/IEC TR 20547-2:2018

Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Data produced by climate simulations is plays a large role in informing discussion of climate change simulations. Therefore it must be robust, both from the standpoint of providing a scientifically valid representation of processes that influence climate, but also as that data is stored long term and transferred world-wide to collaborators and other scientists.
	Visualization	Visualization is crucial to understanding a system as complex as the Earth ecosystem.
	Data quality (syntax)	Earth system scientists are being inundated by an explosion of data generated by ever-increasing resolution in both global models and remote sensors.
	Data types	There is a need to provide data reduction and analysis web services through the Earth System Grid (ESG). A pressing need is emerging for data analysis capabilities closely linked to data archives.
Big data Specific Challenges (Gaps)	Data analytics	
Big data specific challenges (Gaps)	Data from simulations and observations must be shared among a large widely distributed community.	
Big data specific challenges in mobility		
Security and privacy technical considerations	ESGF is in the early stages of being adapted for use in two additional domains: biology (to accelerate drug design and development) and energy (infrastructure for California Energy Systems for the 21st Century (CES21)).	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	http://esgf.org/ http://www-pcmdi.llnl.gov/ http://www.nersc.gov/ http://science.energy.gov/ber/research/cesd/ http://www2.cisl.ucar.edu/	
More information (URLs)		

A.8.9 Use case 49: Subsurface Biogeochemistry

Use case title	DOE-BER Subsurface Biogeochemistry Scientific Focus Area	
Vertical (area)	Research: Earth Science	
Author/company/email	Deb Agarwal, Lawrence Berkeley Lab. daagarwal@lbl.gov	
Actors/stakeholders and their roles and responsibilities	LBNL Sustainable Systems SFA 2.0, Subsurface Scientists, Hydrologists, Geophysicists, Genomics Experts, JGI, Climate scientists, and DOE SBR.	
Goals	The Sustainable Systems Scientific Focus Area 2.0 Science Plan ("SFA 2.0") has been developed to advance predictive understanding of complex and multiscale terrestrial environments relevant to the DOE mission through specifically considering the scientific gaps defined above.	
Use case description	Development of a Genome-Enabled Watershed Simulation Capability (GEWaSC) that will provide a predictive framework for understanding how genomic information stored in a subsurface microbiome affects biogeochemical watershed functioning, how watershed-scale processes affect microbial functioning, and how these interactions co-evolve. While modeling capabilities developed by our team and others in the community have represented processes occurring over an impressive range of scales (ranging from a single bacterial cell to that of a contaminant plume), to date little effort has been devoted to developing a framework for systematically connecting scales, as is needed to identify key controls and to simulate important feedbacks. A simulation framework that formally scales from genomes to watersheds is the primary focus of this GEWaSC deliverable.	
Current solutions	Compute(System)	NERSC
	Storage	NERSC
	Networking	ESNet
	Software	PFLOWTran, postgres, HDF5, Akuna, NEWT, etc.
Big data characteristics	Data source (distributed/centralized)	Terabase-scale sequencing data from JGI, subsurface and surface hydrological and biogeochemical data from a variety of sensors (including dense geophysical datasets) experimental data from field and lab analysis
	Volume (size)	
	Velocity (e.g. real time)	

	<p>Variety (multiple datasets, mashup)</p>	<p>Data crosses all scales from genomics of the microbes in the soil to watershed hydro-biogeochemistry. The SFA requires the synthesis of diverse and disparate field, laboratory, and simulation datasets across different semantic, spatial, and temporal scales through GE-WaSC. Such datasets will be generated by the different research areas and include simulation data, field data (hydrological, geochemical, geophysical), 'omics data, and data from laboratory experiments.</p>
	<p>Variability (rate of change)</p>	<p>Simulations and experiments</p>
<p>Big data science (collection, curation, analysis, action)</p>	<p>Veracity (Robustness Issues, semantics)</p>	<p>Each of the sources samples different properties with different footprints – extremely heterogeneous. Each of the sources has different levels of uncertainty and precision associated with it. In addition, the translation across scales and domains introduces uncertainty as does the data mining. Data quality is critical.</p>
	<p>Visualization</p>	<p>Visualization is crucial to understanding the data.</p>
	<p>Data quality (syntax)</p>	<p>Described in “Variety” above.</p>
	<p>Data types</p>	<p>Data mining, data quality assessment, cross-correlation across datasets, reduced model development, statistics, quality assessment, data fusion, etc.</p>
<p>Big data Specific Challenges (Gaps)</p>	<p>Data analytics</p>	
<p>Big data specific challenges (Gaps)</p>	<p>Field experiment data taking would be improved by access to existing data and automated entry of new data via mobile devices.</p>	
<p>Big data specific challenges in mobility</p>		
<p>Security and privacy technical considerations</p>	<p>A wide array of programs in the earth sciences are working on challenges that cross the same domains as this project.</p>	
<p>Highlight issues for generalizing this Use case (e.g. for ref. architecture)</p>	<p>Under development</p>	
<p>More information (URLs)</p>		

A.8.10 Use case 50: AmeriFlux and FLUXNET

Use case title	DOE-BER AmeriFlux and FLUXNET Networks	
Vertical (area)	Research: Earth Science	
Author/company/email	Deb Agarwal, Lawrence Berkeley Lab. daagarwal@lbl.gov	
Actors/stakeholders and their roles and responsibilities	AmeriFlux scientists, Data Management Team, ICOS, DOE TES, USDA, NSF, and Climate modelers.	
Goals	AmeriFlux Network and FLUXNET measurements provide the crucial linkage between organisms, ecosystems, and process-scale studies at climate-relevant scales of landscapes, regions, and continents, which can be incorporated into biogeochemical and climate models. Results from individual flux sites provide the foundation for a growing body of synthesis and modeling analyses.	
Use case description	AmeriFlux network observations enable scaling of trace gas fluxes (CO ₂ , water vapor) across a broad spectrum of times (hours, days, seasons, years, and decades) and space. Moreover, AmeriFlux and FLUXNET datasets provide the crucial linkages among organisms, ecosystems, and process-scale studies—at climate-relevant scales of landscapes, regions, and continents—for incorporation into biogeochemical and climate models	
Current solutions	Compute(System)	NERSC
	Storage	NERSC
	Networking	ESNet
	Software	EddyPro, Custom analysis software, R, python, neural networks, Matlab.
Big data characteristics	Data source (distributed/centralized)	~150 towers in AmeriFlux and over 500 towers distributed globally collecting flux measurements.
	Volume (size)	
	Velocity (e.g. real time)	
	Variety (multiple datasets, mashup)	The flux data is relatively uniform, however, the biological, disturbance, and other ancillary data needed to process and to interpret the data is extensive and varies widely. Merging this data with the flux data is challenging in today's systems.
	Variability (rate of change)	

Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Each site has unique measurement and data processing techniques. The network brings this data together and performs a common processing, gap-filling, and quality assessment. Thousands of users
	Visualization	Graphs and 3D surfaces are used to visualize the data.
	Data quality (syntax)	Described in "Variety" above.
	Data types	Data mining, data quality assessment, cross-correlation across datasets, data assimilation, data interpolation, statistics, quality assessment, data fusion, etc.
Big data Specific Challenges (Gaps)	Data analytics	
Big data specific challenges (Gaps)	Field experiment data taking would be improved by access to existing data and automated entry of new data via mobile devices.	
Big data specific challenges in mobility		
Security and privacy technical considerations		
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	http://Ameriflux.lbl.gov http://www.fluxdata.org	
More information (URLs)		

A.9 Energy

A.9.1 Use case 51: Consumption Forecasting in Smart Grids

Use case title	Consumption forecasting in Smart Grids	
Vertical (area)	Energy Informatics	
Author/company/email	Yogesh Simmhan, University of Southern California, simmhan@usc.edu	
Actors/stakeholders and their roles and responsibilities	Electric Utilities, Campus MicroGrids, Building Managers, Power Consumers, Energy Markets	
Goals	Develop scalable and accurate forecasting models to predict the energy consumption (kWh) within the utility service area under different spatial and temporal granularities to help improve grid reliability and efficiency.	
Use case description	<p>Deployment of smart meters are making available near-realtime energy usage data (kWh) every 15 min at the granularity individual consumers within the service area of smart power utilities. This unprecedented and growing access to fine-grained energy consumption information allows novel analytics capabilities to be developed for predicting energy consumption for customers, transformers, sub-stations and the utility service area. Near-term forecast can be used by utilities and microgrid managers to take preventive action before consumption spikes cause brown/blackouts through demand-response optimization by engaging consumers, bringing peaker units online, or purchasing power from the energy markets. These form an OODA feedback loop. Customers can also use them for energy use planning and budgeting. Medium- to long-term predictions can help utilities and building managers plan generation capacity, renewable portfolio, energy purchasing contracts and sustainable building improvements.</p> <p>Steps involved include 1) <i>Data Collection and Storage</i>: time-series data from (potentially) millions of smart meters in near-realtime, features on consumers, facilities and regions, weather forecasts, archival of data for training, testing and validating models; 2) <i>Data Cleaning and Normalization</i>: Spatio-temporal normalization, gap filling/Interpolation, outlier detection, semantic annotation; 3) <i>Training Forecast Models</i>: Using univariate timeseries models like ARIMA, and data-driven machine learning models like regression tree, ANN, for different spatial (consumer, transformer) and temporal (15-min, 24-hour) granularities; 4) <i>Prediction</i>: Predict consumption for different spatio-temporal granularities and prediction horizons using near-realtime and historic data fed to the forecast model with thresholds on prediction latencies.</p>	
Current solutions	Compute(System)	Many-core servers, Commodity Cluster, Workstations
	Storage	SQL Databases, CSV Files, HDFS, Meter Data Management
	Networking	Gigabit Ethernet
	Software	R/Matlab, Weka, Hadoop
Big data characteristics	Data source (distributed/centralized)	Head-end of smart meters (distributed), Utility databases (Customer Information, Network topology; centralized), US Census data (distributed), NOAA weather data (distributed), Microgrid building information system (centralized), Microgrid sensor network (distributed)
	Volume (size)	10 GB/day; 4 TB/year (<i>City scale</i>)

	Velocity (e.g. real time)	Los Angeles: Once every 15 min (~100 k streams); Once every 8 h (~1,4 M streams) with finer grain data aggregated to 8-hour interval
	Variety (multiple datasets, mashup)	Tuple-based: Timeseries, database rows; Graph-based: Network topology, customer connectivity; Some semantic data for normalization.
	Variability (rate of change)	Meter and weather data change, and are collected/used, on hourly basis. Customer/building/grid topology information is slow changing on a weekly basis
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	Versioning and reproducibility is necessary to validate/compare past and current models. Resilience of storage and analytics is important for operational needs. Semantic normalization can help with inter-disciplinary analysis (e.g. utility operators, building managers, power engineers, behavioral scientists)
	Visualization	Map-based visualization of grid service topology, stress; Energy heat-maps; Plots of demand forecasts vs. capacity, what-if analysis; Realtime information display; Apps with push notification of alerts
	Data quality (syntax)	Gaps in smart meters and weather data; Quality issues in sensor data; Rigorous checks done for "billing quality" meter data;
	Data types	Timeseries (CSV, SQL tuples), Static information (RDF, XML), topology (shape files)
	Data analytics	Forecasting models, machine learning models, time series analysis, clustering, motif detection, complex event processing, visual network analysis,
Big data specific challenges (Gaps)	Scalable realtime analytics over large data streams Low-latency analytics for operational needs Federated analytics at utility and microgrid levels Robust time series analytics over millions of customer consumption data Customer behavior modeling, targeted curtailment requests	
Big data specific challenges in mobility	Apps for engaging with customers: Data collection from customers/premises for behavior modeling, feature extraction; Notification of curtailment requests by utility/building managers; Suggestions on energy efficiency; Geo-localized display of energy footprint.	

Security and privacy technical considerations	Personally identifiable customer data requires careful handling. Customer energy usage data can reveal behavior patterns. Anonymization of information. Data aggregation to avoid customer identification. Data sharing restrictions by federal and state energy regulators. Surveys by behavioral scientists may have IRB restrictions.
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	Realtime data-driven analytics for cyber physical systems
More information (URLs)	http://smartgrid.usc.edu http://ganges.usc.edu/wiki/Smart_Grid https://www.ladwp.com/ladwp/faces/ladwp/aboutus/a-power/a-p-smartgridla http://ieeexplore.ieee.org/xpl/articleDetails.jsp?arnumber=6475927

A.9.2 Use case 52: Home Energy Management System

Use case title	Homes Energy Management System	
Vertical (area)	Service Business	
Author/company/email	Japan National Body	
Actors/stakeholders and their roles and responsibilities	Users are general people who live in private homes. Energy manager is the company who provides several sensors and devices into private homes. Energy manager may be a home builder, or energy company. Information manager is the company who gathers data from private homes and responds to privacy and security of users. Servicer is the company who analyzes those data and provides valuable information to users as a service.	
Goals	Provide useful information services to users by combining and analyzing power usage data with other available data.	
Use case description	HEMS (Home Energy Management System) is useful system for energy conservation in private homes. In the HEMS, many kinds of sensors and devices are introduced into private homes, such as, smart meter, electric vehicle, solar power panel, light, air conditioner, fuel cell, water heater, storage battery. Energy manager gathers those data generated at private homes and stores them into cloud database named the large HEMS information platform. Information manager operates the large HEMS information platform and manages data. Privacy and security of users are responsible to Information manager. Servicer analyzes data and provides valuable information to users as a service.	
Current solutions	Compute(System)	n/a
	Storage	n/a
	Networking	n/a
	Software	n/a
Big data characteristics	Data source (distributed/centralized)	Data sources are distributed into individual private homes.
	Volume (size)	About 14,000 households. n/a about data size

	Velocity (e.g. real time)	Real time, streaming data from sensors
	Variety (multiple datasets, mashup)	smart meter, electric vehicle, solar power panel, light, air conditioner, fuel cell, water heater, storage battery
	Variability (rate of change)	Data varies in seconds, minutes or hours
Big data science (collection, curation, analysis, action)	Veracity (Robustness Issues, semantics)	
	Visualization	Visualization of sensor data is basic service.
	Data quality (syntax)	Data quality affects directly to the quality of services.
	Data types	Timeseries
	Data analytics	Prediction, time series analysis
Big data specific challenges (Gaps)		
Big data specific challenges in mobility		
Security and privacy technical considerations	<p>Personally identifiable data should be carefully handled to protect user's privacy.</p> <p>Service needs to inform how to handle the data and what kind of service is provided. Users need to select services which they want to obtain by considering service's information.</p>	
Highlight issues for generalizing this Use case (e.g. for ref. architecture)	<p>Several players act as supply chain in the big data flow. One of the player A receives data from Data source (Data provider). The Player A provides data to another player B. Both Player A and B analyze data. Thus data provision may chain through several players.</p>	
More information (URLs)	<p>https://www.ntt-east.co.jp/release/detail/20140828_01.html (in Japanese)</p>	
NOTE <additional comments>		

Annex B

Summary of Key Properties

Information related to five key properties was extracted from each use case. The five key properties were three big data characteristics (volume, velocity, and variety), software related information, and associated analytics. The extracted information is presented in [Table B.1](#).

Table B.1 — Use case Specific Information by Key Properties

Use case	Volume	Velocity	Variety	Software	Analytics
A.1.1 UC #001 Census 2000 and 2010	380 TB	Static for 75 years	Scanned documents	Robust archival storage	None for 75 years
A.1.2 UC #002 NARA: Search, Retrieve, Preservation	Hundreds of terabytes, and growing	Data loaded in batches, so bursty	Unstructured and structured data: textual documents, emails, photos, scanned documents, multimedia, social networks, web sites, databases, etc.	Custom software, commercial products, commercial databases	Crawl/index, search, ranking, predictive search; data categorization (sensitive, confidential, etc.); personally identifiable information (PII) detection and flagging
A.1.3 UC #003 Statistical Survey Response Improvement	Approximately 1 PB	Variable, field data streamed continuously, Census was ~150 million records transmitted	Strings and numerical data	Hadoop, Spark, Hive, R, SAS, Mahout, Allegrograph, MySQL, Oracle, Storm, BigMemory, Cassandra, Pig	Recommendation systems, continued monitoring
A.1.4 UC #004 Non-Traditional Data in Statistical Survey Response Improvement	—	—	Survey data, other government administrative data, web-scraped data, wireless data, e-transaction data, (potentially) social media data and positioning data from various sources	Hadoop, Spark, Hive, R, SAS, Mahout, Allegrograph, MySQL, Oracle, Storm, BigMemory, Cassandra, Pig	New analytics to create reliable information from non-traditional disparate sources
A.2.1 UC #005 Cloud Eco-System for Finance	—	Real time	—	Hadoop RDBMS XBRL	Fraud detection
A.2.2 UC #006 Mendeley	15 TB presently, growing about 1 TB per month	Currently Hadoop batch jobs scheduled daily, real-time recommended in future	PDF documents and log files of social network and client activities	Hadoop, Scribe, Hive, Mahout, Python	Standard libraries for machine learning and analytics, LDA, custom-built reporting tools for aggregating readership and social activities per document

Table B.1 (continued)

Use case	Volume	Velocity	Variety	Software	Analytics
A.2.3 UC #007 Netflix Movie Service	Summer 2012 – 25 million subscribers, 4 million ratings per day, 3 million searches per day, 1 billion hours streamed in June 2012; Cloud storage – 2 petabytes in June 2013	Media (video and properties) and rankings continually updated	Data vary from digital media to user rankings, user profiles, and media properties for content-based recommendations	Hadoop and Pig; Cassandra; Teradata	Personalized recommender systems using logistic/linear regression, elastic nets, matrix factorization, clustering, LDA, association rules, gradient-boosted decision trees, and others; streaming video delivery
A.2.4 UC #008 Web Search	45 billion web pages total, 500 million photos uploaded each day, 100 hours of video uploaded to YouTube each minute	Real-time updating and real-time responses to queries	Multiple media	Map/Reduce + Bigtable; Dryad + Cosmos; PageRank; final step essentially a recommender engine	Crawling; searching, including topic-based searches; ranking; recommending
A.2.5 UC #009 Business Continuity and Disaster Recovery Within a Cloud Eco-System	Terabytes up to petabytes	Can be real time for recent changes	Must work for all data	Hadoop, Map/Reduce, open source, and/or vendor proprietary such as AWS, Google Cloud Services, and Microsoft	Robust backup
A.2.6 UC #010 Cargo Shipping	—	Needs to become real time, currently updated at events	Event-based	—	Distributed event analysis identifying problems
A.2.7 UC #011 Materials Data for Manufacturing	500,000 material types in 1980s, much growth since then	Ongoing increase in new materials	Many datasets with no standards	National programs (Japan, Korea, and China), application areas (EU nuclear program), proprietary systems (Granta, etc.)	No broadly applicable analytics
A.2.8 UC #012 Simulation-Driven Materials Genomics	100 TB (current), 500 TB within five years, scalable key-value and object store databases needed	Regular data added from simulations	Varied data and simulation results	MongoDB, GFFS, PyMatGen, FireWorks, VASP, ABINIT, NWChem, BerkeleyGW, varied community codes	Map/Reduce and search that join simulation and experimental data
A.3.1 UC #013 Large-Scale Geospatial Analysis and Visualization	Imagery – hundreds of terabytes; vector data – tens of GBs but billions of points	Vectors transmitted in near real time	Imagery, vector (various formats such as shape files, KML, text streams) and many object structures	Geospatially enabled RDBMS, Esri ArcServer, Geoserver	Closest point of approach, deviation from route, point density over time, PCA and ICA

Table B.1 (continued)

Use case	Volume	Velocity	Variety	Software	Analytics
A.3.2 UC #014 Object Identification and Tracking	FMV – 30–60 frames per second at full-color 1080P resolution; WALF – 1–10 frames per second at 10,000 × 10,000 full-color resolution	Real time	A few standard imagery or video formats	Custom software and tools including traditional RDBMS and display tools	Visualization as overlays on a GIS, basic object detection analytics and integration with sophisticated situation awareness tools with data fusion
A.3.3 UC #015 Intelligence Data Processing and Analysis	Tens of terabytes to hundreds of petabytes, individual warfighters (first responders) would have at most one to hundreds of GBs	Much real-time, imagery intelligence devices that gather a petabyte of data in a few hours	Text files, raw media, imagery, video, audio, electronic data, human-generated data	Hadoop, Accumulo (BigTable), Solr, NLP, Puppet (for deployment and security) and Storm; GIS	Near real-time alerts based on patterns and baseline changes, link analysis, geo-spatial analysis, text analytics (sentiment, entity extraction, etc.)
A.4.1 UC #016 EMR Data	12 million patients, more than 4 billion discrete clinical observations, >20 TB raw data	0.5 – 1.5 million new real-time clinical transactions added per day	Broad variety of data from doctors, nurses, laboratories and instruments	Teradata, PostgreSQL, MongoDB, Hadoop, Hive, R	Information retrieval methods (tf-idf), NLP, maximum likelihood estimators, Bayesian networks
A.4.2 UC #017 Pathology Imaging	1 GB raw image data + 1.5 GB analytical results per 2D image, 1 TB raw image data + 1 TB analytical results per 3D image, 1 PB data per moderated hospital per year	Once generated, data will not be changed	Images	MPI for image analysis, Map/Reduce + Hive with spatial extension	Image analysis, spatial queries and analytics, feature clustering and classification
A.4.3 UC #018 Computational Bioimaging	Medical diagnostic imaging around 70 PB annually, 32 TB on emerging machines for a single scan	Volume of data acquisition requires HPC back end	Multi-modal imaging with disparate channels of data	Scalable key-value and object store databases; ImageJ, Omero, VolRover, advanced segmentation and feature detection methods	Machine learning (support vector machine [SVM] and random forest [RF]) for classification and recommendation services
A.4.4 UC #019 Genomic Measurements	>100 TB in 1 to 2 years at NIST, many PBs in healthcare community	~300 GB of compressed data/day generated by DNA sequencers	File formats not well-standardized, though some standards exist; generally structured data	Open-source sequencing bioinformatics software from academic groups	Processing of raw data to produce variant calls, clinical interpretation of variants

Table B.1 (continued)

Use case	Volume	Velocity	Variety	Software	Analytics
A.4.5 UC #020 Comparative Analysis for Metagenomes and Genomes	50 TB	New sequencers stream in data at growing rate	Biological data that are inherently heterogeneous, complex, structural, and hierarchical; besides core genomic data, new types of omics data such as transcriptomics, methylomics, and proteomics	Standard bioinformatics tools (BLAST, HMMER, multiple alignment and phylogenetic tools, gene callers, sequence feature predictors), Perl/Python wrapper scripts	Descriptive statistics, statistical significance in hypothesis testing, data clustering and classification
A.4.6 UC #021 Individualized Diabetes Management	5 million patients	Not real time but updated periodically	100 controlled vocabulary values and 1,000 continuous values per patient, mostly time-stamped values	HDFS supplementing Mayo internal data warehouse (EDT)	Integration of data into semantic graphs, using graph traverse to replace SQL join; development of semantic graph-mining algorithms to identify graph patterns, index graph, and search graph; indexed Hbase; custom code to develop new patient properties from stored data
A.4.7 UC #022 Statistical Relational Artificial Intelligence for Health Care	Hundreds of GBs for a single cohort of a few hundred people; possibly on the order of 1 PB when dealing with millions of patients	Constant updates to EHRs; in other controlled studies, data often in batches at regular intervals	Critical features – data typically in multiple tables, need to be merged to perform analysis	Mainly Java-based, in-house tools to process the data	Relational probabilistic models (Statistical Relational AI) learned from multiple data types
A.4.8 UC #023 World Population-Scale Epidemiological Study	100 TB	Low number of data feeding into the simulation, massive amounts of real-time data generated by simulation	Can be rich with various population activities, geographical, socio-economic, cultural variations	Charm++, MPI	Simulations on a synthetic population
A.4.9 UC #024 Social Contagion Modeling for Planning	Tens of terabytes per year	During social unrest events, human interactions and mobility leads to rapid changes in data; e.g., who follows whom in Twitter	Big issues – data fusion, combining data from different sources, dealing with missing or incomplete data	Specialized simulators, open source software, proprietary modeling environments; databases	Models of behavior of humans and hard infrastructures, models of their interactions, visualization of results

Table B.1 (continued)

Use case	Volume	Velocity	Variety	Software	Analytics
A.4.10 UC #025 Biodiversity and LifeWatch	N/A	Real-time processing and analysis in case of natural or industrial disaster	Rich variety and number of involved databases and observation data	RDBMS	Requires advanced and rich visualization
A.5.1 UC #026 Large-Scale Deep Learning	Current datasets typically 1 to 10 TB, possibly 100 million images to train a self-driving car	Much faster than real-time processing, for autonomous driving, need to process thousands of high-resolution (six megapixels or more) images per second	Neural net very heterogeneous as it learns many different features	In-house GPU kernels and MPI-based communication developed by Stanford, C++/Python source	Small degree of batch statistical pre-processing, all other data analysis performed by the learning algorithm itself
A.5.2 UC #027 Organizing Large-Scale Unstructured Collections of Consumer Photos	500+ billion photos on Facebook, 5+ billion photos on Flickr	Over 500 million images uploaded to Facebook each day	Images and metadata including EXIF (Exchangeable Image File type, etc.)	Hadoop Map/Reduce, simple hand-written multi-threaded tools (Secure Shell [SSH] and sockets for communication)	Robust non-linear least squares optimization problem, SVM
A.5.3 UC #028 Truthy Twitter Data	30 TB/year compressed data	Near real-time data storage, querying and analysis	Schema provided by social media data source; currently using Twitter only; plans to expand, incorporating Google+ and Facebook	Hadoop IndexedHBase and HDFS; Hadoop, Hive, Redis for data management; Python: SciPy NumPy and MPI for data analysis	Anomaly detection, stream clustering, signal classification, online learning; information diffusion, clustering, dynamic network visualization
A.5.4 UC #029 Crowd Sourcing in Humanities	GBs (text, surveys, experiment values) to hundreds of terabytes (multimedia)	Data continuously updated and analyzed incrementally	So far mostly homogeneous small data sets; expected large distributed heterogeneous datasets	XML technology, traditional relational databases	Pattern recognition (e.g., speech recognition, automatic audio-visual analysis, cultural structures (lexical units, linguistic rules, etc.)
A.5.5 UC #030 CINET for Network Science	Can be hundreds of GBs for a single network, 1,000–5,000 networks and methods	Dynamic networks, network collection growing	Many types of networks	Graph libraries (Galib, NetworkX); distributed workflow management (Simfracture, databases, semantic web tools)	Network visualization