# Storage and sharing of large 3D imaging datasets

Richard Boardman[1], **Ian Sinclair**[1], Simon Cox[1], Philippa Reed[1], Kenji Takeda[1,2], Jeremy Frey[1] and Graeme Earl[3]

[1]μ-VIS X-ray Imaging Centre
**www.southampton.ac.uk/muvis**
**muvis@soton.ac.uk**
[2] Microsoft Research Connections, Cambridge, UK
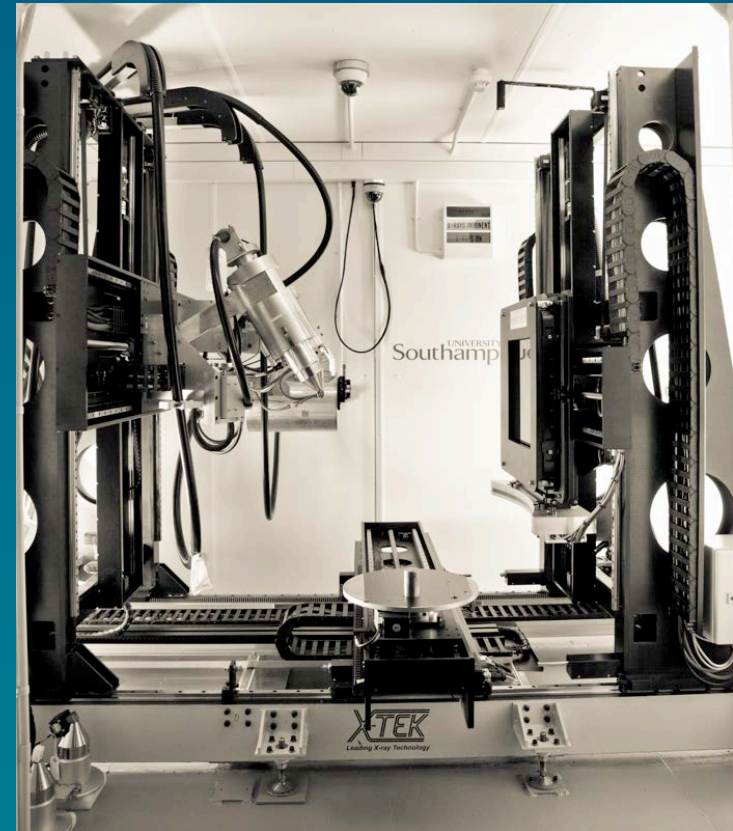[3] Computational Archaeology, Southampton, UK

# Content

- Overview
  - Background
  - Delineation of scope & roles
  - Metadata quality

- Example engineering activities

- Metadata & database strategies

- Archiving practicalities

- Final thoughts

# Motivation

- Significant investments in the generation of large voxel datasets (projects, scanning devices…)
  - High fidelity, large 3D datasets almost inevitably contain more potential than the original researcher/project intended

- To keep value, it is essential to retain the data and record parameters surrounding their acquisition and processing
  - Scientific diligence (e.g. experimental reproducibility…)
  - Sharing: extending the data life cycle
  - Funding body requirements

- Unshared data is a loss to science and engineering
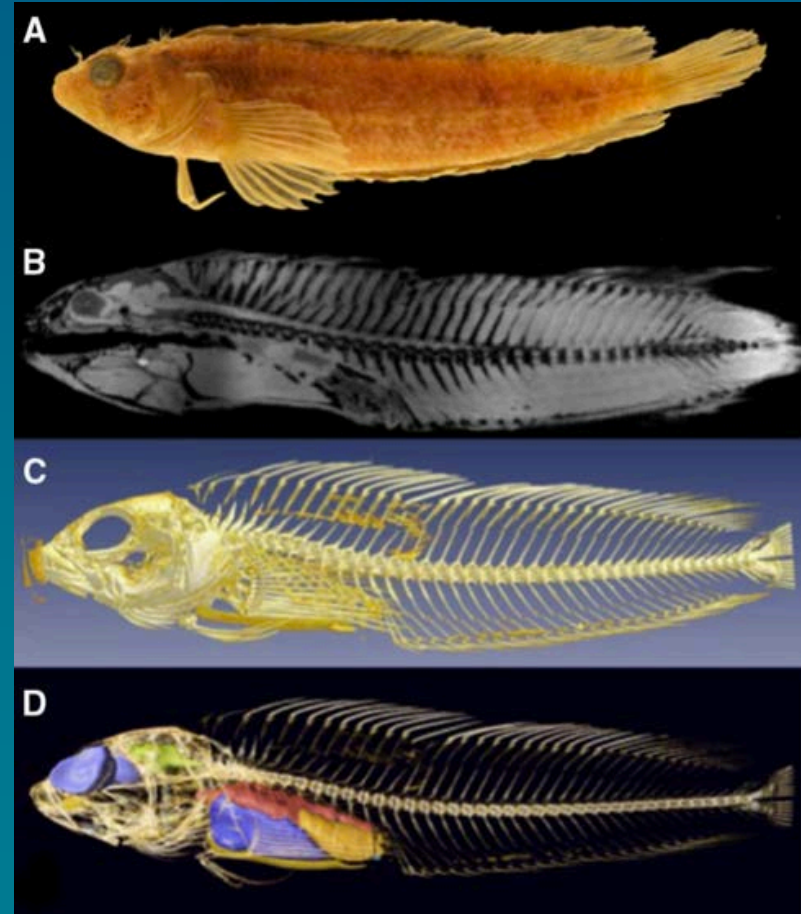
# Data: ownership & management

- The roles that are involved
  - Data authors & users
  - Supervisors
  - Facility managers
  - Computer scientists
  - Institutional leaders
  - Funders (government/others)
  - Open access 'evangelists'
  - Salesmen
  - Legal aspects
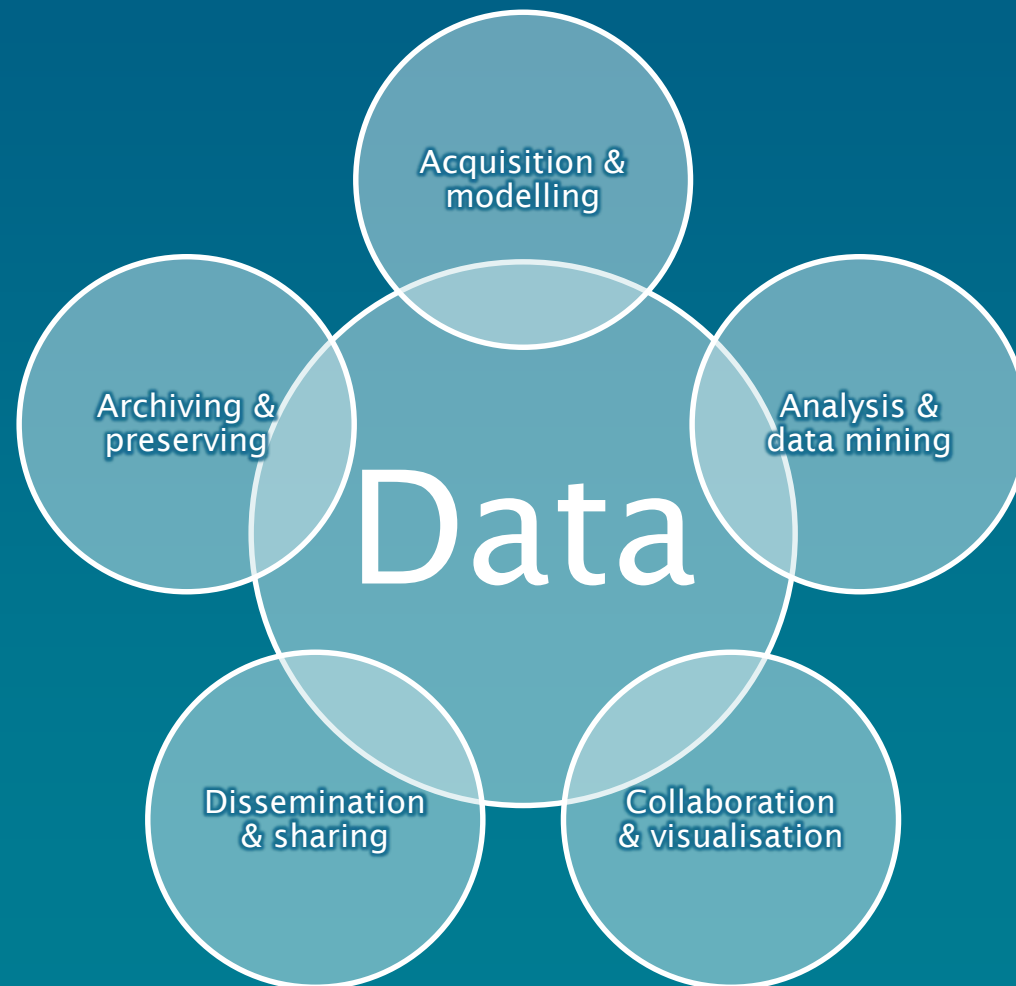


μ-VIS X-ray Imaging Centre
www.southampton.ac.uk/muvis

# Early Developments

- www.digimorph.org
  - XCT data
  - >1000 bio/palaeo samples
- www.digitalfishlibrary.org
  - MRI data
  - >300 samples (fish!)
- Data reduced to 2D and animations
  - <5Mb
- Raw voxels not available as yet
- Recent example
  - 3D Materials Atlas



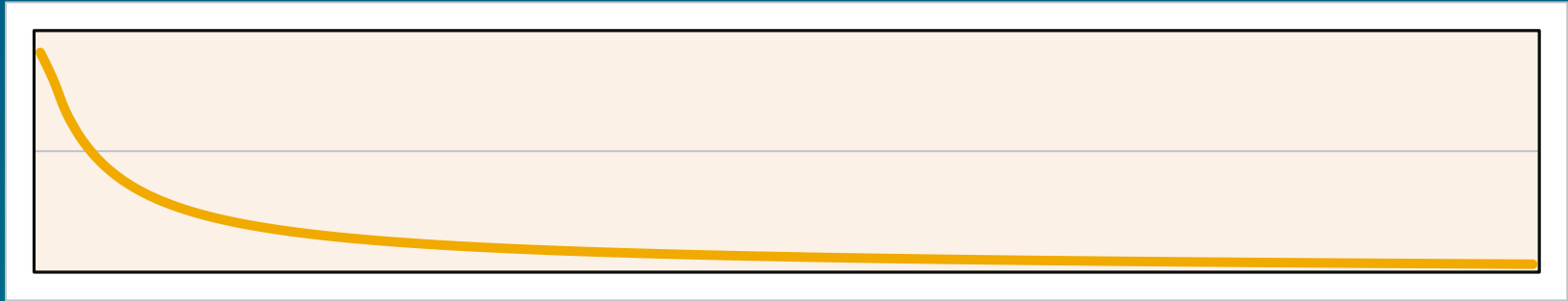Island Kelpfish: MRI & CT data: DigiMorph & Digial Fish Library

# Data-intensive science

# Acquisition: The data deluge

- We can generate data faster than we can consume it
  - Rate of generation now exceeds physical storage capacity (Feb. 2011)*

- Synchrotrons: terabytes per day
  - SLS: ~5TB/day (fast acquisition)
  - AS: 200TB/year (growing to 400TB/year with new beamlines)
  - ESRF (*ca.* 2010): O(100TB) 30-day storage, O(PB) for backups

- µ-VIS lab facility: up to two terabytes per day (robotic operation)
  - 20GB projections + 30GB reconstruction = 50GB in as little as 10-15 minutes
  - Plus O(10MB) metadata

- *LHC ~ 50-100PB/yr, ~20PB stored*

* Hilbert & Lopez, (2011), Science

# Long-tail science



- Small numbers of major projects/facilities responsible for a lot of output
  - Formal data management policies & resources

- Large number of smaller projects/facilities also do a lot!
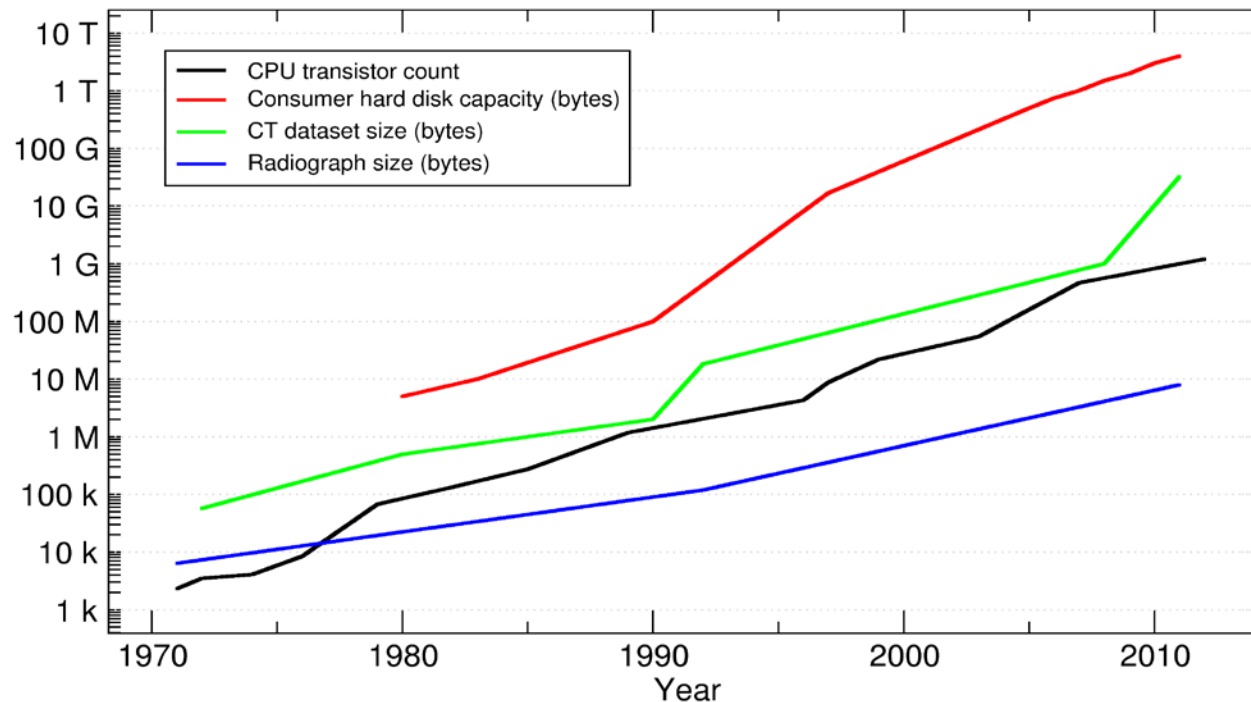  - Data management policies & resources very variable

# What do we mean by "large"

- For our purposes, right now we will say that it is datasets which are O(1GB – 100GB+)

- They won't necessarily fit on CDs, DVDs, Blu-rays

- They will fit on hard drives

- Transferring them around even in isolation may present a challenge (portable drives, institutional networks, FTP/rsync/GridFTP?)
    - → Grid/Cloud capabilities & tools, e.g. Globus, MS Azure, Amazon S3, **Dropbox**
    - → 'Never underestimate the bandwidth of a stationwagon full of tapes hurtling down the highway'?

# What about Moore's Law?

- Integrated circuit transistor count doubles every two years

- Broadly applicable to many areas of technology, including hard disk capacity

- Perhaps we can just wait a while and Moore's Law will help us store data?

- Unfortunately, our large 3D datasets also "obey" Moore's Law
  - This large 3D dataset problem will always be with us

# What about Moore's Law?



Sources: Kalendar, W (2011); μ-VIS X-ray Imaging Centre; IBM; AMD; Intel; DEC; Seagate; Western Digital

- Solution? Greater allocation of resources to storage/archiving – who pays?

# Southampton data sharing projects

- data.gov.uk – "Opening up Government"
  - Founded by Nigel Shadbolt and Tim Berners-Lee

- Open Data Service
  - data.southampton.ac.uk

- Research data access
  - datapool.soton.ac.uk
  - 10 year roadmap
  - *Recognising need for new services, policy framework and data management support*

# Southampton data sharing projects

- Data grades (stars)
  1. Anything/'stuff'
  2. Structured data, e.g. Excel file instead of jpeg of a data table
  3. Open format, e.g. CSV vs. Excel files
  4. Provide persistent link
  5. Links to others data/information, to provide context



- 1* is great, but must aim for 5*

13

# Implementation of e-lab book

- Blog based format

- Purpose built engine

- Fully flexible system with arbitrary metadata

- Full record of changes



**http://chemtools.chem.soton.ac.uk/projects/blog/** "Bio Blogs"

**http://blogs.openwetware.org/scienceintheopen** Discussion

# Implementation of e-lab book

- "Facebook for Scientists" …but different to Facebook!

- A repository of research methods

- A community social network of people and things

- Machinery for coordinating the execution of (scientific) services and linking together (scientific) resources

- Open source (BSD) Ruby on Rails application with HTML, REST and SPARQL interfaces



www.myexperiment.org

# Preserving the record

- Key goal: record the whole experimental process prior to and during, rather than after
  - Ensures we efficiently generate a traceable, complete record of the work
  - Foundation for high quality sharing & reuse of data, extending data life
  - Scalable from a single lab to whole communities

**TAGtivity**

Organizing thoughts and reference material... experience

**Wiki**

Semantic structure, search and reasoning

**BAE Systems**

Software + Services for connecting engineers and experts to users and data

**Rolls-Royce**

Orchestration of gas turbine design calculations

**Airbus**

Robust, reliable and scalable data intensive collaboration

# Centre for Fluid Dynamics Simulation Project

**Concept** → **Computation** → **Data**

Audit trail & Review

Individual review

Small group review

Corporate review

## Data sources

| Tagtivity database Filesystem | Wiki database Workflow database Knowledge database Corporate database | Task database Conversation database Workflow tracking Simulation database | Workflow templates Workflow tracking Filesystem | Sharepoint database Active Directory |

## Technology

| Microsoft Office 2007 SQL Server 2005/ 2008 Windows Presentation Foundation; Matlab | MediaWiki SQL Server 2008 D2R Server ARQ/SPARQL | Windows Server 2008 Hyper-V RC0 Windows HPC Server 2008 Beta2 SQL Server 2008 CTP6 Office Communication Server | Windows Workflow Foundation Windows CCS 2003 Linux (Interop) Visual Studio 2005 SQL Server 2005 | Sharepoint Server Active Directory HP-UX (Interop) Windows Communication Foundation |

18

# µ-VIS X-ray Imaging Centre

- Five CT scanners, including
  - 225/450kV custom "hutch", imaging up to 1x2m, panel shift and line detector
  - 225kV Nikon/Metris HMX with rototic sample exchange
  - Largest single scan >1TB
  - 60TB online data store, 10GbE connectivity
  - Workstations up to 32 CPU cores/128GB RAM/nVidia Tesla GPU rack
  - >100 users/year



µ-VIS X-ray Imaging Centre
www.southampton.ac.uk/muvis

Multidisciplinary, Multiscale, Microtomographic Volume Imaging at Southampton

μ-VIS

UNIVERSITY OF Southampton

| Concept | → | Execution | → | Data |
|---------|---|-----------|---|------|

| Beamtime application | Experiment design (*e.g.* custom mounts, scan condition control) | Scheduling, acquisition and reconstruction | Data analysis | Long term archiving and sharing |
|---------|---------|---------|---------|---------|

**Data sources/targets**

| IMAP email Bugzilla database | Bugzilla database Wiki database Metadata database | Google calendar 10GbE central filestore | 10GbE central filestore Wiki database | Shared 10GbE filesystem Hard disks Dropbox Metadata DB |
|---------|---------|---------|---------|---------|

**Technology**

| HTML/PHP Perl Apache web server Bugzilla | Bugzilla MediaWiki Perl Apache | CTPro (FBP) Digisens (ART) Windows Server 2008R2 Linux | VGStudio MAX Avizo Simpleware ImageJ Matlab IDL | Python MySQL SMB/NFS/FTP |
|---------|---------|---------|---------|---------|

# What data storage and sharing means

Data storage: database, central file store

| Machine acquired | Machine generated | Human generated |
|---|---|---|
| Radiographs | Analyses | |
| Sinograms | Reduced datasets | |
| Shading corrections | Photographs | |
| Scan metadata | | Enquiry metadata |
| Environmental information (*e.g.* radiation levels, temperature) | Volume reconstructions | |
| | Visualisations | |

… many more

# What metadata are relevant?

- In CT, a sensible minimum is:
  - Two projections (at 0 and 90 degrees)
  - A central slice of the reconstructed volume
  - All the available acquisition condition metadata (filters, kV, µA, source to detector distance *&c.*)

- For one CT scan, this might be 30 or 40MB; much more manageable than 50GB

- Once metadata are stored, a web interface provides tools to review and search

# Further metadata

- In our case, if a user adds something "extra" to a CT scan directory, then this is also captured
  - Photographs of the sample, special sample mounting rigs, documentation, charts, videos, anything else

- ... and "Smart Pen" output (operator notes) is added (➜ searchable pdf file)

- E-lab books TBC

# Metadata: browsing and searching

| ID | Data ID | Scan ID | Index date | Name | kV | μA | Exposure (ms) | Projections | 0° | 90° | XY slice | First extra image? |
|----|---------|---------|------------|------|----|----|--------------|-------------|-----|-----|----------|--------------------|
| 2077 | 1104 | 338 | 2012-05-30 23:05:55 | 20120530_HUTCH_338_NS_Giraffe_cranium_1 | 380 | 800 | 125 | 1901 | | | | |
| 2078 | 1104 | 338 | 2012-05-30 23:05:57 | 20120530_HUTCH_338_NS_Giraffe_cranium_1 | 380 | 800 | 125 | 1901 | | | | |

| | |
|---|---|
| XraykV | 380 |
| XrayuA | 800 |
| Stack | 0 |
| Slice | 0 |
| SinogramOffsetX | 0 |
| SinogramBandSampling | 1 |
| SliceThreshold | 0 |
| SliceAreaStartX | 125 |
| SliceAreaEndX | 875 |
| SliceAreaStartY | 125 |
| SliceAreaEndY | 875 |
| Version | V2.2.4182.18577 (Date:1 |
| Product | Product:[XT: CT Pro 3D], |
| Filter_ThicknessMM | 3 |
| Filter_Material | Copper |

Web browser interface
- 2 projections & central slice
- Extensive metadata
- Dataset names, IDs, times
- NetApp/archive location
- Original proposal, emails…

# Metadata standards: DICONDE

- Direct mapping of DICOM to industrial CT
- Firmly established approach, detailed

| ASTM No. | Title | Description | Status |
|---|---|---|---|
| E2339 | Digital Imaging and Communication in Nondestructive Evaluation (DICONDE) | Data and metadata that applies to ALL NDE methods. | Issued 2004 |
| E2663 | Digital Imaging and Communication in Nondestructive Evaluation (DICONDE) for Ultrasonic Test Methods | Data and metadata that are relevent only to ultrasonic test methods | Issued 2008 |
| E2767 | Digital Imaging and Communication in Nondestructive Evaluation (DICONDE) for X-ray Computed Tomography (CT) Test Methods | Data and metadata that are relevent only to x-ray computed tomography test methods | Issued 2010 |
| E2699 | Digital Imaging and Communication in Nondestructive Evaluation (DICONDE) for Digital Radiographic (DR) Test Methods | Data and metadata that are relevent only to digital radiographic test methods | Issued 2010 |
| WK20537 | Digital Imaging and Communication in Nondestructive Evaluation (DICONDE) for Eddy Current Test Methods | Data and metadata that are relevent only to eddy current test methods | Waiting on Public Attributes |
| E2738 | Digital Imaging and Communication Nondestructive Evaluation (DICONDE) for Computed Radiography (CR) Test Methods | Data and metadata that are relevent only to computed radiography test methods | Issued 2009 |

# Heterogeneous Data Centre

- To provide a user-centric software system for users to store and share their data and metadata in a usable way

- The user decides their own metadata structures
  - Stored as name-value pairs and can be hierarchical, providing a flexible approach to data management

- To support a wide variety of data, from small text files to large voxel data files.

- Provide the ability for users to tag data sets with **any** relevant metadata

# CT Dataset browsing in HDC

# Heterogeneous Data Centre

- Data can be uploaded via EPrints with the EP2DC service or directly

# Archiving practicalities

- Many options available: carefully indexed disks or tapes, online NAS, cloud storage

- It is impossible to *100% guarantee* that data will never be lost
  - We can get close (90%, 99%, 99.9%...)

- Cost scales with reliability

# Archiving practicalities *(continued)*

- One copy on one hard disk: ~10-20% chance of data loss over 5 years
  - Approximate cost in 2012: ~$10/TB/year

- Two copies on two separate disks: ~1-4% chance of data loss over 5 years
  - Approximate cost in 2012: ~$20/TB/year

- "Enterprise" class storage (*e.g.* NetApp): <1% chance of data loss over 5 years
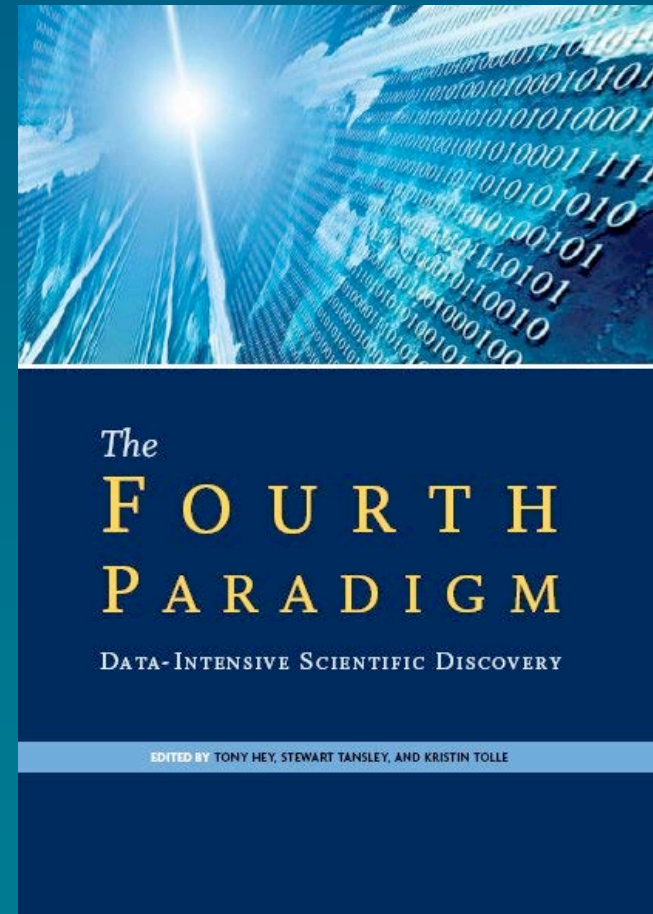  - Approximate cost in 2012: ~$500/TB/year

# Cloud storage

- Provides a scalable and reliable option to store data, e.g. Amazon S3
  - '11 nines' reliability levels

- Typical pricing is O($0.10)/GB/month (2012)
  - around $1200/TB/year; additional charges for uploading and downloading

- Recently, providers have been waiving upload charges
  - 30GB download O($5)
  - May make storing large amounts of data with relatively few downloads more attractive

# Final thoughts

- Local behaviour – what's going on in your lab? Are people carefully looking after their datasets?
  - The generation of quality data and metadata is best done *concurrently*
  - If we look after it, we can make better use of it

- Look for technology that will work well with your, requirements, current systems and budgets
  - *Many* strategies & tools are already in place
  - BioSimGrid, ROOT…

- The first step is, start now…
  - The sooner, the better

# Final thoughts

- Contribution to the '4<sup>th</sup> Paradigm'?
    - Science driven by the capture, curation, analysis of large data
    - *All* data becomes publically **available** and **usable**, like books in the library

- *1<sup>st</sup>: Empirical description of nature (~1000 yrs ago)*

- *2<sup>nd</sup>: Mathematical theory (~100yrs)*

- *3<sup>rd</sup>: Large simulation (~30 yrs)*



The

# FOURTH PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

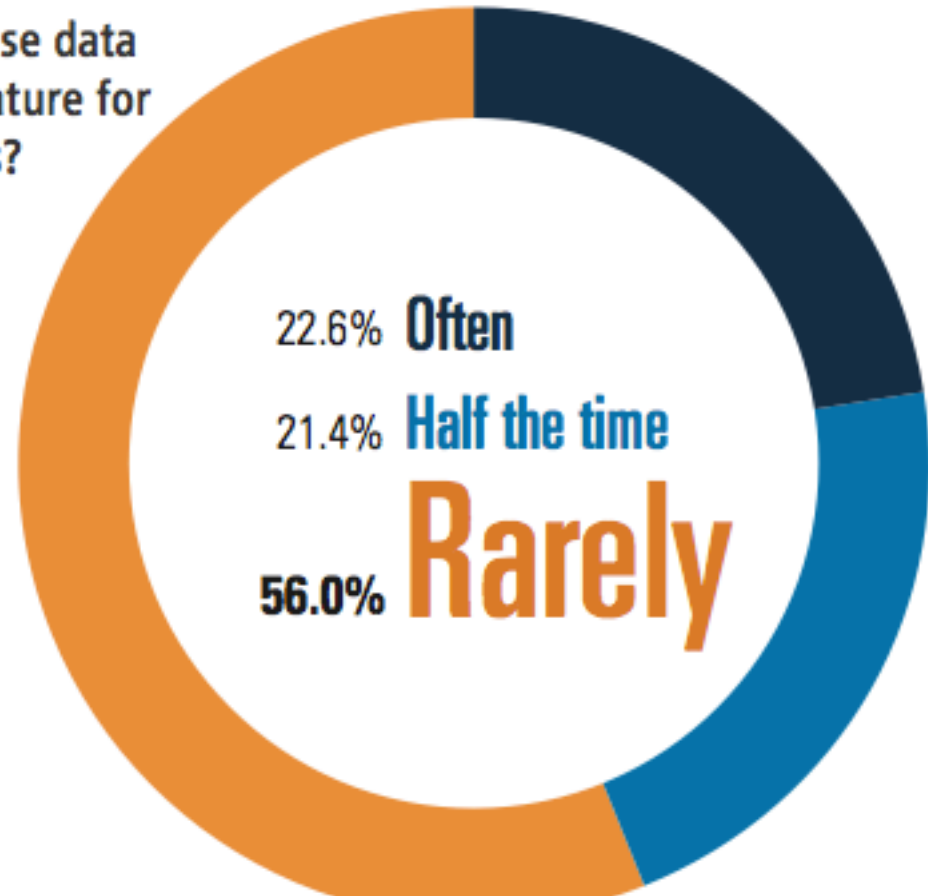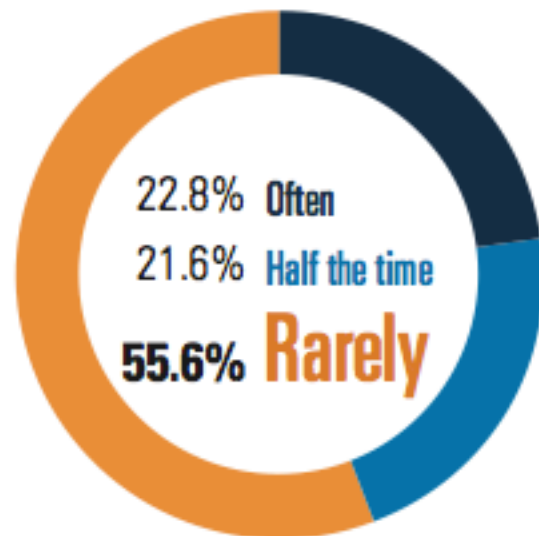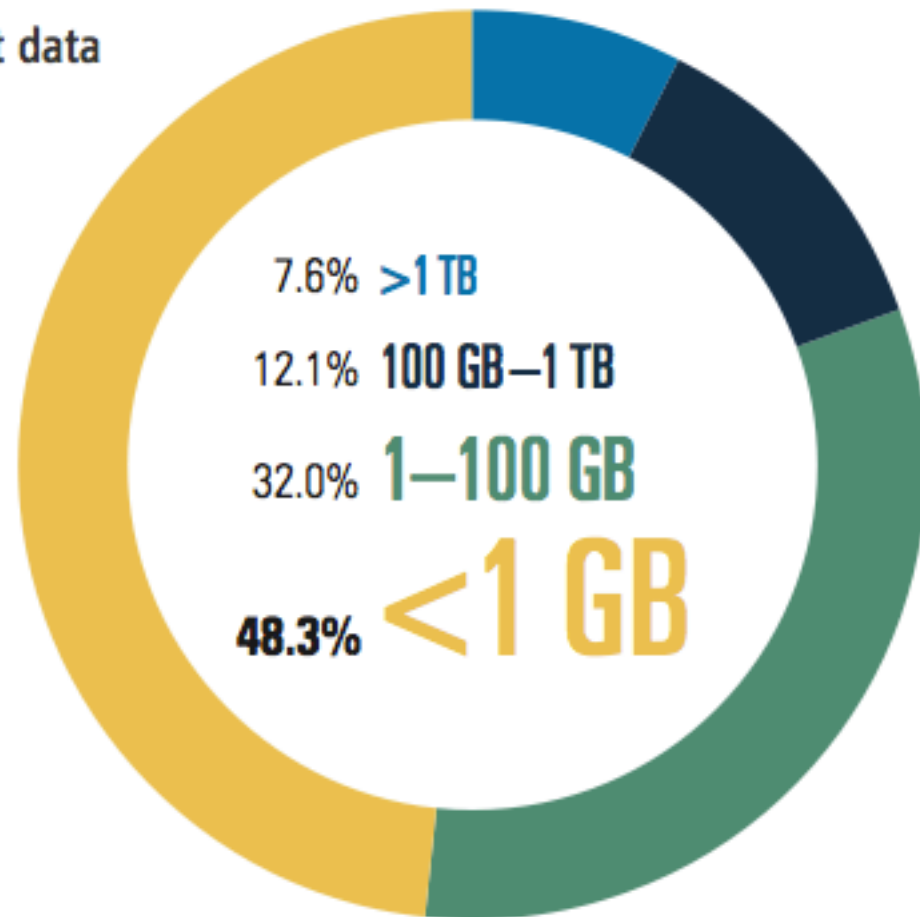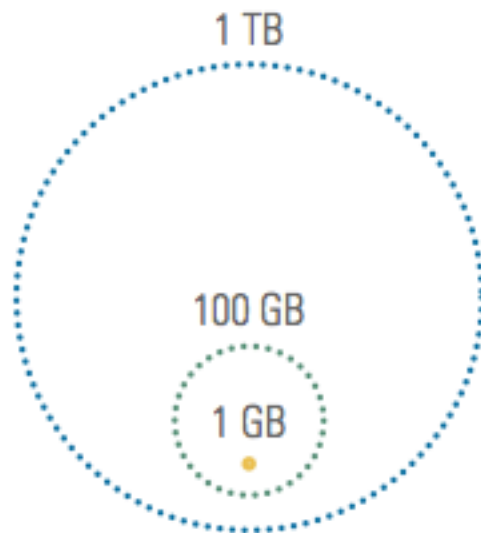EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

# Acknowledgements

- The μ-VIS team

- Mark Scott (HDC)

- Oliver Bunk (Swiss Light Source)

- Uli Felzmann (Australian Synchrotron)

How often do you access or use data sets from the published literature for your original research papers?

From archival databases?

**Archival databases:**
22.8% Often
21.6% Half the time
55.6% Rarely

**Published literature:**
22.6% Often
21.4% Half the time
56.0% Rarely