



OGSA-DAI Introduction and Overview

OGSA-DAI Tutorial
GGF15, Boston, USA
6 October 2005

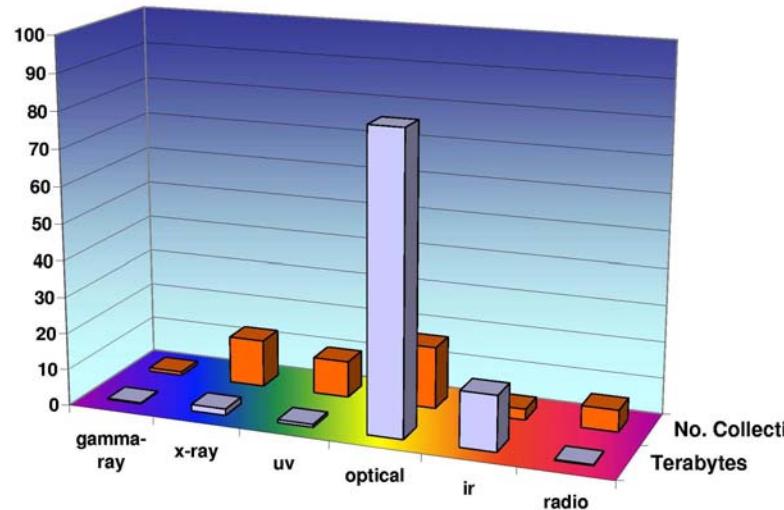
Neil Chue Hong
Project Manager, EPCC
N.ChueHong@epcc.ed.ac.uk
+44 131 650 5957

Malcolm Atkinson
Director, National e-Science Centre
mpa@nesc.ac.uk
+44 131 651 4040

- Growing volumes
 - Growing diversity
 - Growing complexity
- } **⇒ Rich resource**
- How do we mine its riches for nuggets of information?
 - Find & Access
 - Understand
 - Extract, Combine & Digest
 - Test hypotheses
 - Bingo!
- } **⇒ Use OGSA-DAI**

OGSA-DAI

Composing Observations in Astronomy



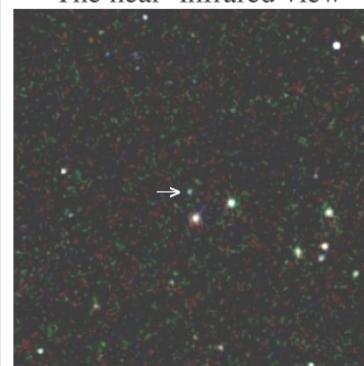
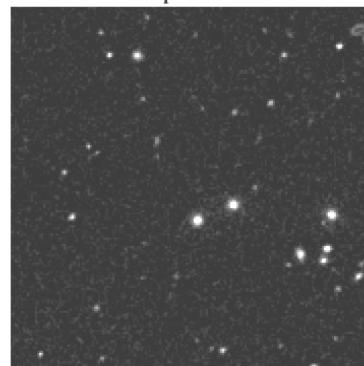
No. & sizes of data sets as of mid-2002,
grouped by wavelength

- 12 waveband coverage of large areas of the sky
- Total about 200 TB data
- Doubling every 12 months
- Largest catalogues near 1B objects

2MASSW J1217-03

A methane (T-type) dwarf in the constellation Virgo

The near-infrared view
The optical view

2MASS Composite JHK_s Atlas Image
Palomar Digitized Sky Survey



2MASS
MICRON ALL-SKY SURVEY

A.J.Burgasser (Caltech), J.D.Kirkpatrick (IPAC/Caltech), M.E.Brown (Caltech), I.N.Reid (U.Penn), J.E.Gizis (U.Mass), C.C.Dahn & D.G.Monet (USNO, Flagstaff), C.A.Bechman (JPL), J.I.Lebert (Arizona), R.M.Cutri (IPAC/Caltech), M.E.Skrutskie (U.Mass)

The 2MASS Project is a collaboration between the University of Massachusetts and IPAC

IB 1, 1999

Astronomers Detect New Category of Elusive 'Brown Dwarfs'

By JOHN NOBLE WILFORD CHICAGO, May 31 — Ambitious Apache Point, N.M. Dr. Michael Strauss and a graduate student, Xiaohui Fan, were able to detect a brown dwarf, but was not associated with a star companion. An estimate of their mass have been possible in hotter, younger objects. An estimate of their mass

Data and images courtesy Alex Szalay, John Hopkins



Biomedical data – making connections

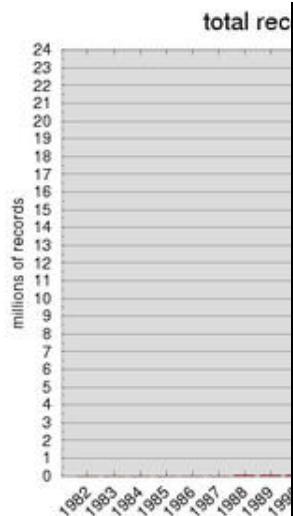


Slide provided by Carole Goble: University of Manchester

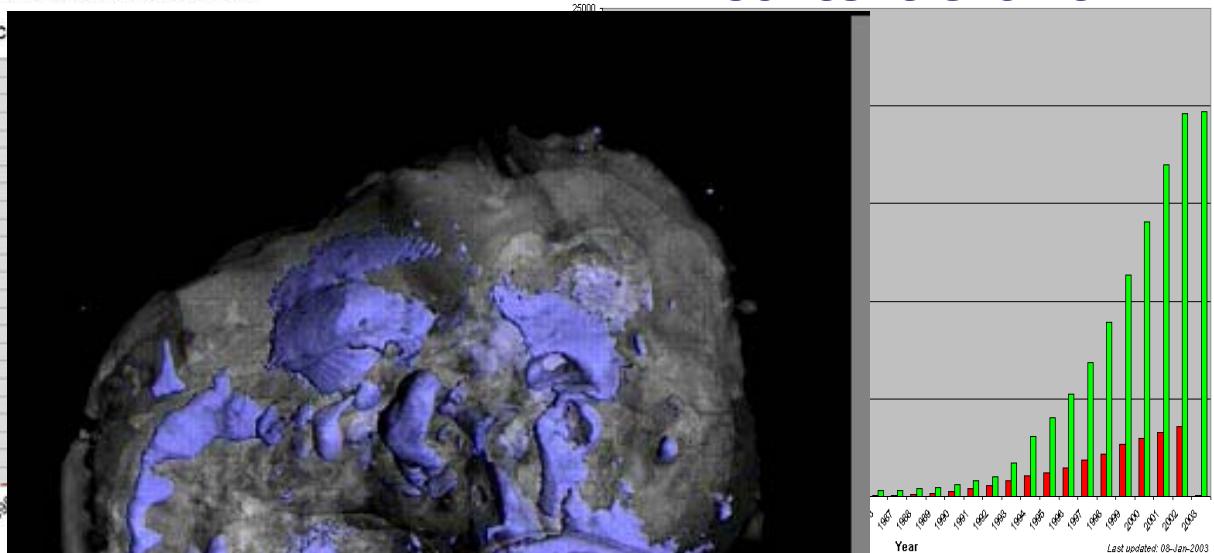


Database Growth

EMBL Database Growth



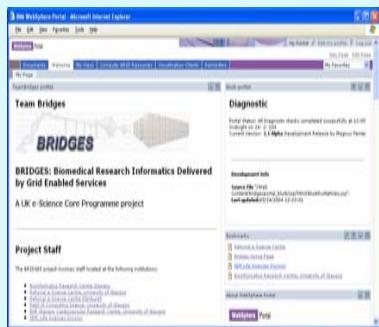
PDB Content Growth



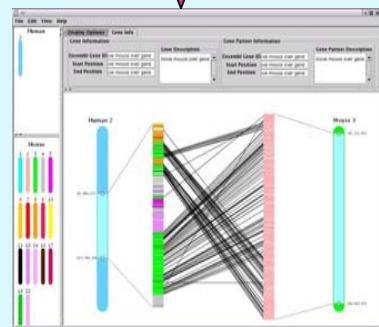
Slide provided by Richard Baldock: MRC HGU Edinburgh

BRIDGES

VO Authorisation



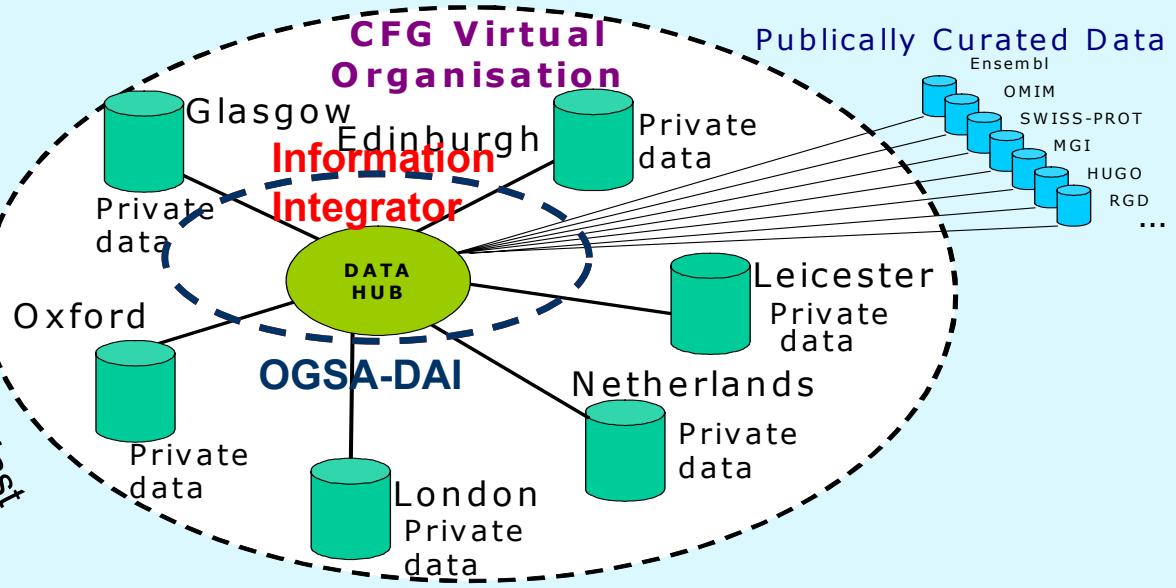
Synteny
Grid
Service



blast

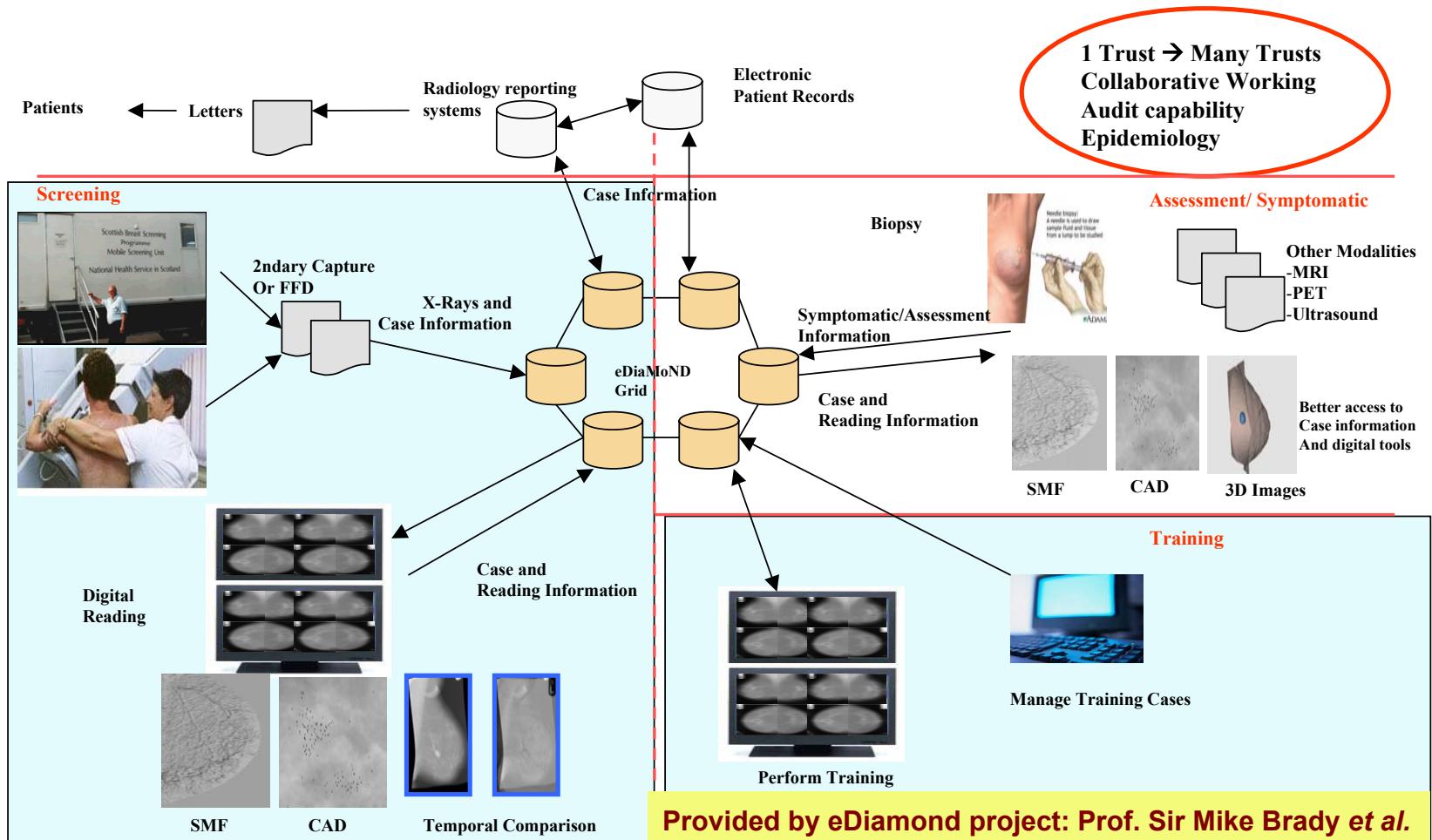


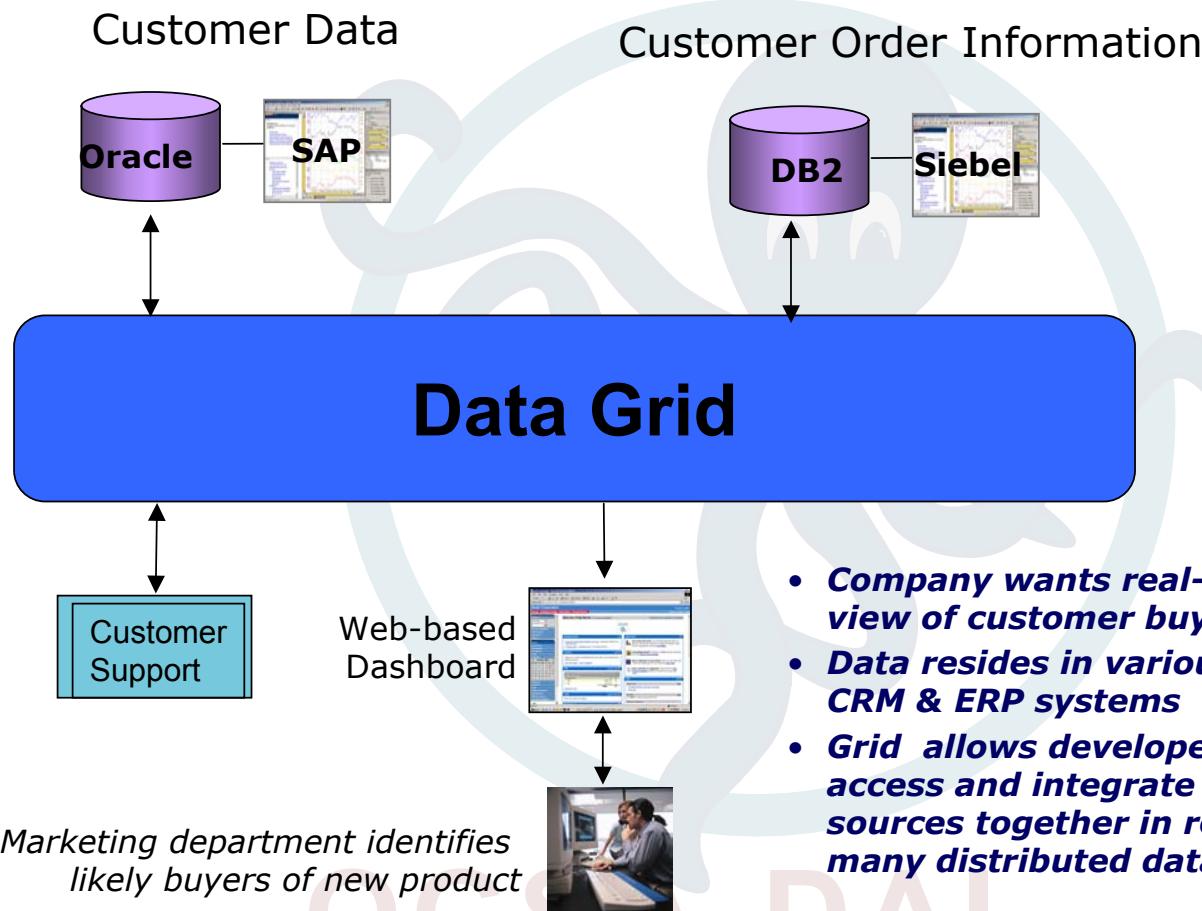
+



Slide provided by Richard Sinnott: University of Glasgow

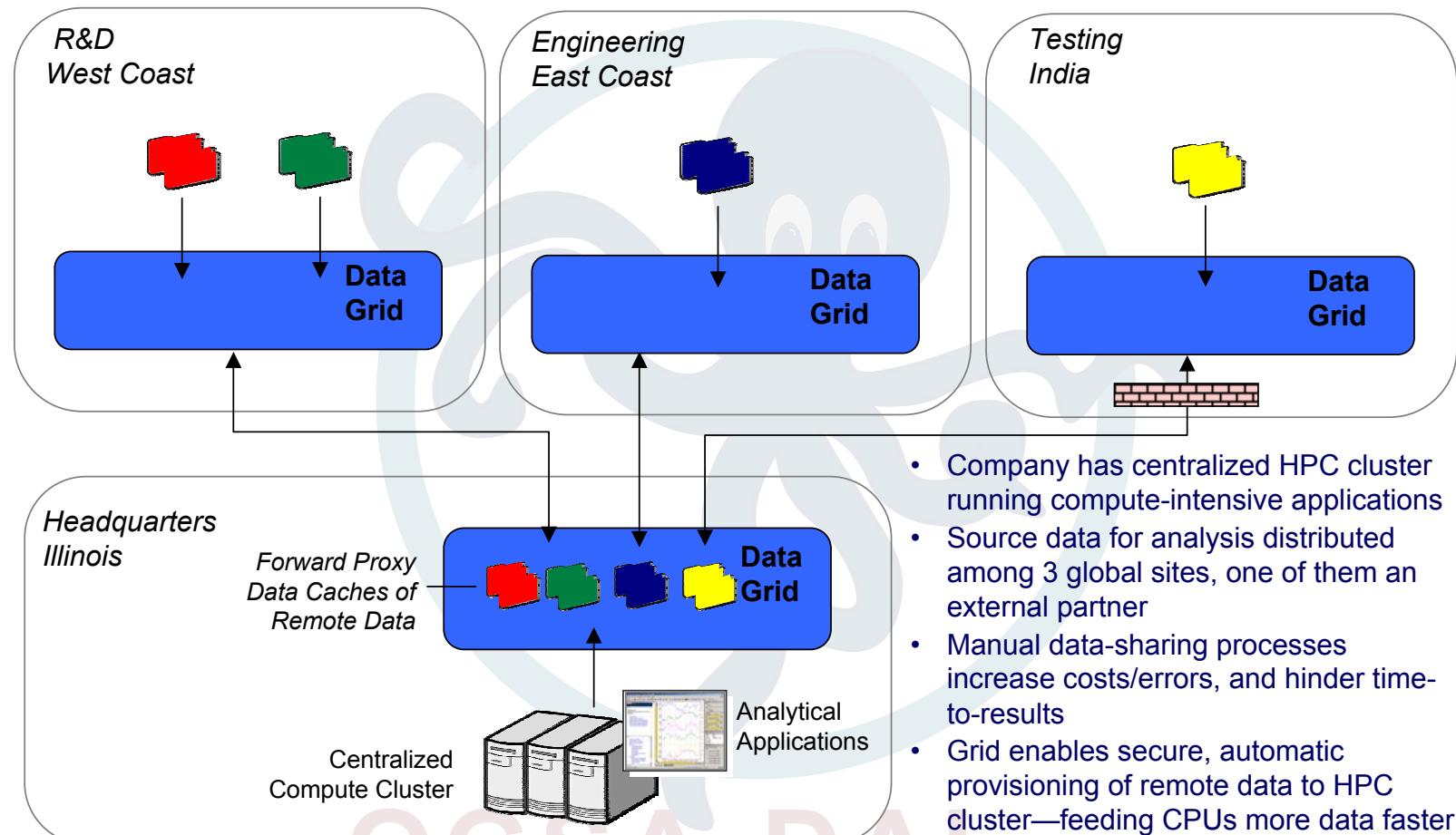
eDiaMoND: Screening for Breast Cancer





- **Company wants real-time integrated view of customer buying behavior**
 - **Data resides in various distributed CRM & ERP systems**
 - **Grid allows developers and apps to access and integrate customer data sources together in real time--across many distributed databases**

Providing Data to Cluster-Based Analytical Application



- There is a lot of data
 - growing in every dimension
 - Distributed
 - Many different producers & owners
 - Heterogeneous
 - High value resource
 - Combined it is more valuable
- There are many requirements for integration
 - Takes many forms
 - Driven by insights
 - Enable conversion of insights to tested hypothesis
- There are many data owners
 - Their autonomy and policies must be respected

Generic Repeatable Solutions Required

OGSA-DAI



Changing the way we manage Data

Neil Chue Hong
Project Manager, EPCC
N.ChueHong@epcc.ed.ac.uk
+44 131 650 5957

Malcolm Atkinson
Director, National e-Science Centre
mpa@nesc.ac.uk
+44 131 651 4040

	Terabyte	Petabyte
RAM time to move	15 minutes	2 months
1GB WAN move time	10 hours (\$1000)	14 months (\$1 million)
Disk cost	7 disks = \$5000 (SCSI)	6800 Disks + 490 units + 32 racks = \$7 million
Disk power	100 Watts	100 Kilowatts
Disk weight	5.6 Kg	33 Tonnes
Disk footprint	Inside machine	60 m ²

Approximately Correct in May 2003 *Distributed Computing Economics*
Jim Gray, Microsoft Research, MSR-TR-2003-24

- Petabytes of Data cannot be moved
 - It stays where it is produced or curated
 - Hospitals, observatories, European Bioinformatics Institute
 - A few caches and a small proportion cached
- Distributed collaborating communities
 - Expertise in curation, simulation & analysis
- Diverse data collections
 - Discovery depends on insights
 - Unpredictable or unexpected use of data

OGSA-DAI

- Assumption: code size << data size
 - Minimise data transport
- Provision combined storage & compute resources
- Develop the database philosophy for this?
- Develop the storage architecture for this?
- Develop experiment, sensor & simulation architectures
 - That take code to select and digest data as an output control
 - That attach the provenance & metadata
- Data Cutter a step in this direction
 - Sub-setting and aggregation of datasets using filters executed close to data
 - <http://www.cs.umd.edu/projects/hpsl/ResearchAreas/DataCutter.htm>

UGDA-DHI

- Choosing data sources
 - How do you find them?
 - How are they described and advertised?
 - Is the equivalent of Google possible?
- Meta-data is required describing
 - Content
 - Provenance
 - Structure
 - Types, Formats & Ontologies
 - Operations available
 - Access requirements
 - Quality of service
- No established standards for heterogeneous data sources

OGSA-DAI

- Changing the way we work?
- Publication and sharing of result data
 - Increased volume and diversity = increased opportunity?
 - Allows independent validation of methods and derivatives
 - Responsibility, ownership, credit, citation
- Many distributed data resources
 - Data collected from observation, simulation & experiment
 - Independently owned & managed
 - No common goals or design
 - Work hard for agreements on foundation types and ontologies
 - Autonomous decisions change data, structure, policy, etc
 - Requires negotiations and patience!
- Diversity
 - No “one size fits all” solutions will work

OGSA-DAI

- Data production, publication & management
 - Many researchers contributing increments of data
 - Who pays for storage, transport, management and curation?
- Data longevity
 - Research requirements may outlive technical decisions
 - Data does not preserve itself indefinitely!
- When a community is dependent on a data resource
 - Who pays or decides to switch it off?

OGSA-DAI

- Obtaining access to that data
 - Overcoming administrative barriers
 - Overcoming technical barriers
- Understanding that data
 - The parts you care about for your research
- Combing them using sophisticated models
 - The picture of reality in your head
- Analysis on scales required by statistics
 - Coupling data access with computation
- Repeated Processes
 - Examining variations, covering a set of candidates
 - Monitoring the emerging details
 - Coupling with scientific workflows

Three communities

epcc1

Users: Individual & Organisations

**Data,
Information &
Knowledge
Providers**

**Computer
Storage &
Communications
Providers**

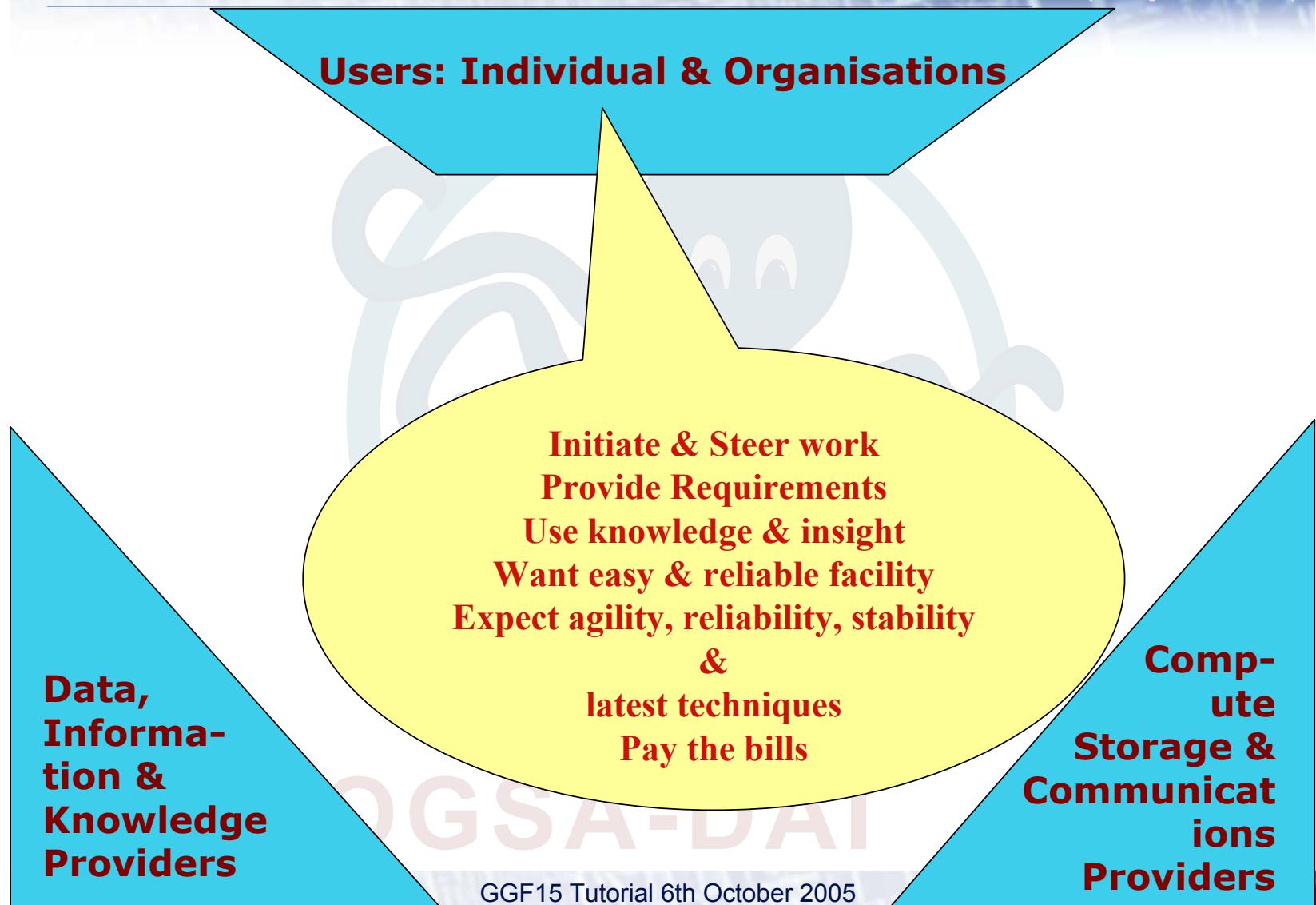


OGSA-DAI

GGF15 Tutorial 6th October 2005

Three communities

|epcc|



Three communities

epcc1

Users: Individual & Organisations

Provide & operate resources
Storage centres, Data centres,
DBMS & File Systems
Computation environment,
Processing & Communications
Need to *change* facilities & policies
Prefer consolidated requirements
Must be paid

Data,
Information &
Knowledge
Providers

Computer
Storage &
Communications
Providers

Three communities

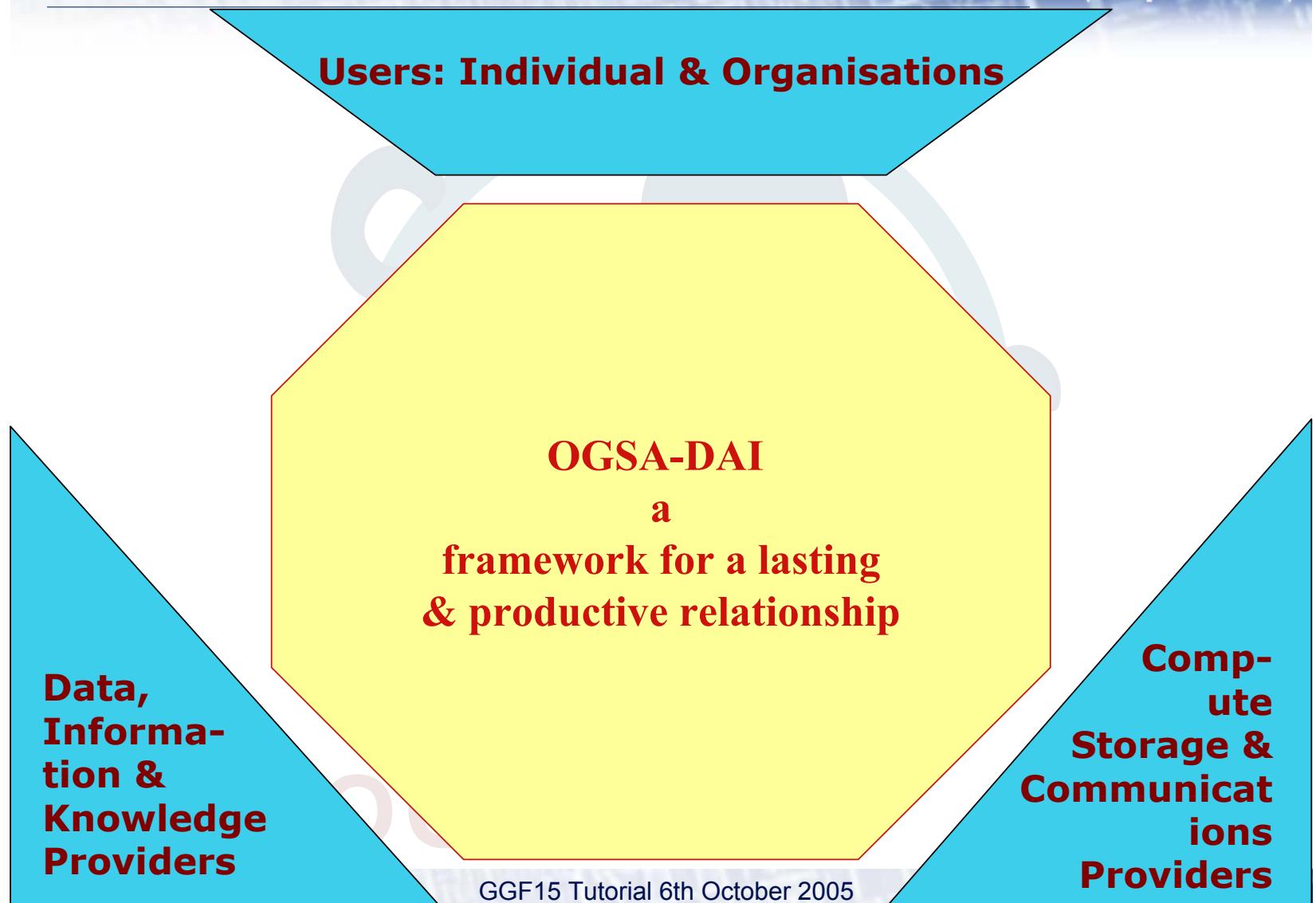
epcc1

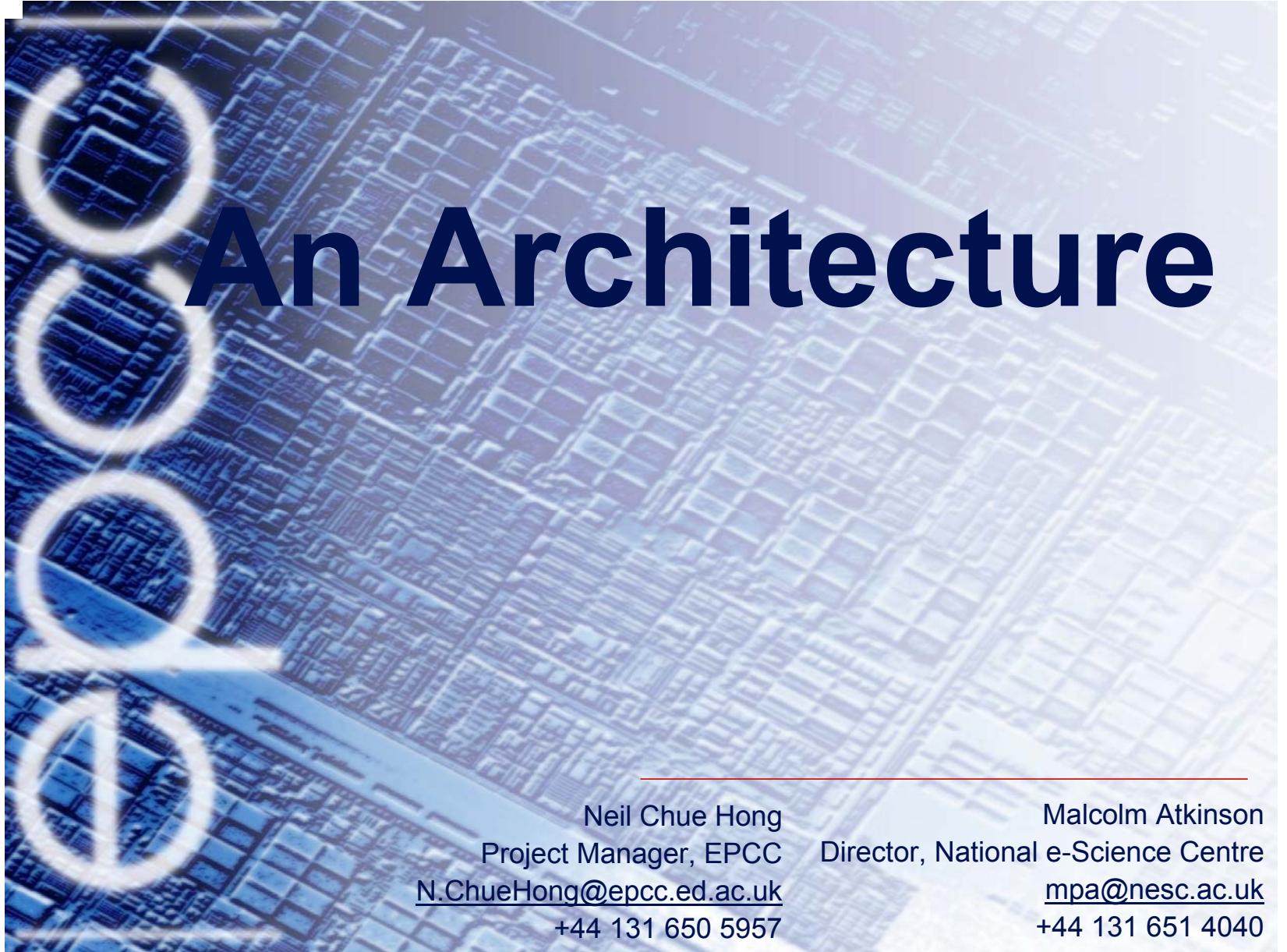
Users: Individual & Organisations

Create & Collect data
Organise & Structure data
Provide, organise & maintain metadata
Offer access and domain specific services
Establish use policies
Will *change* structure, services & policies
May pay or be paid

Data,
Information &
Knowledge
Providers

Computer
Storage &
Communications
Providers





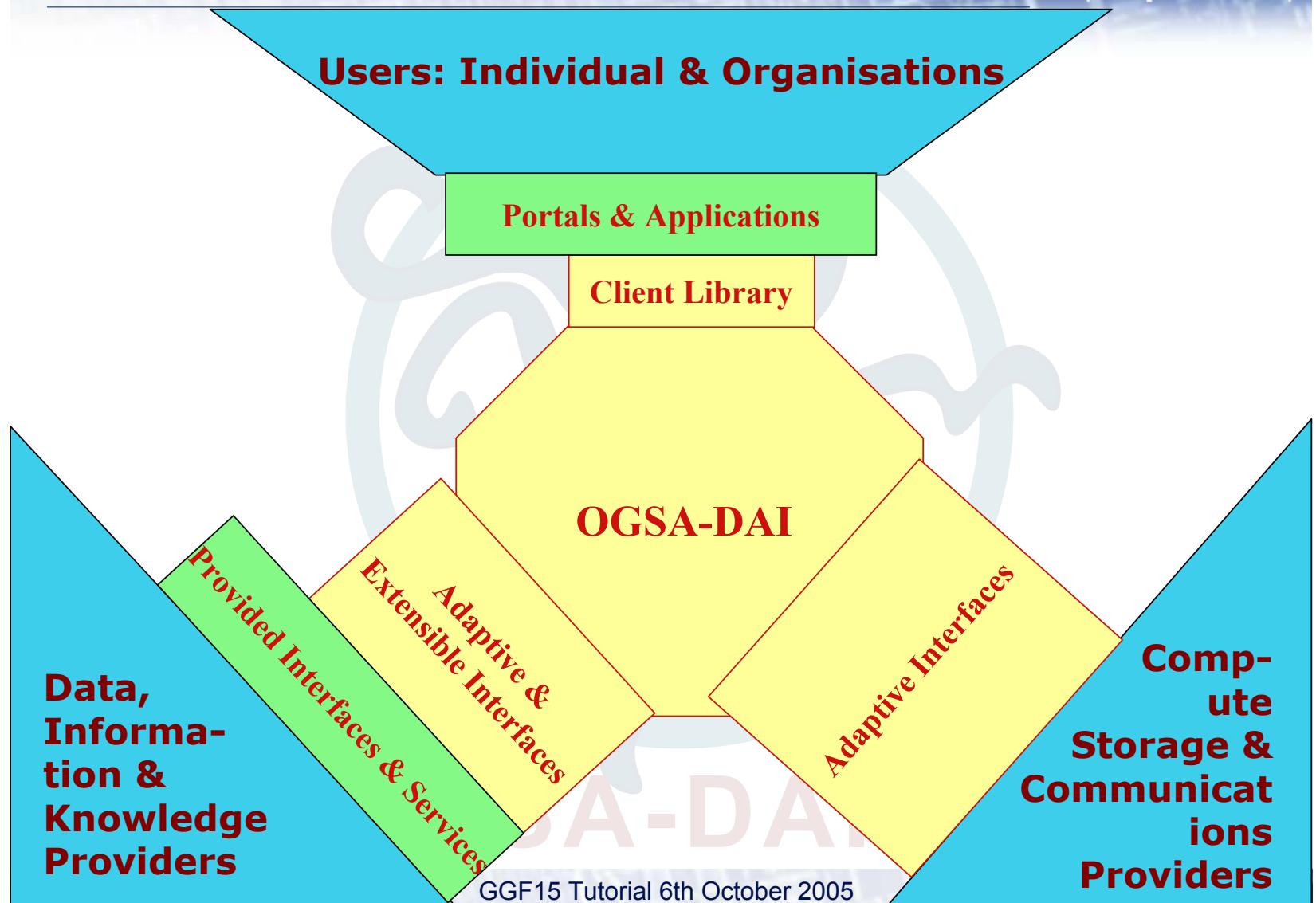
An Architecture

Neil Chue Hong
Project Manager, EPCC
N.ChueHong@epcc.ed.ac.uk
+44 131 650 5957

Malcolm Atkinson
Director, National e-Science Centre
mpa@nesc.ac.uk
+44 131 651 4040

Three communities

|epcc|





DAIS Working Group

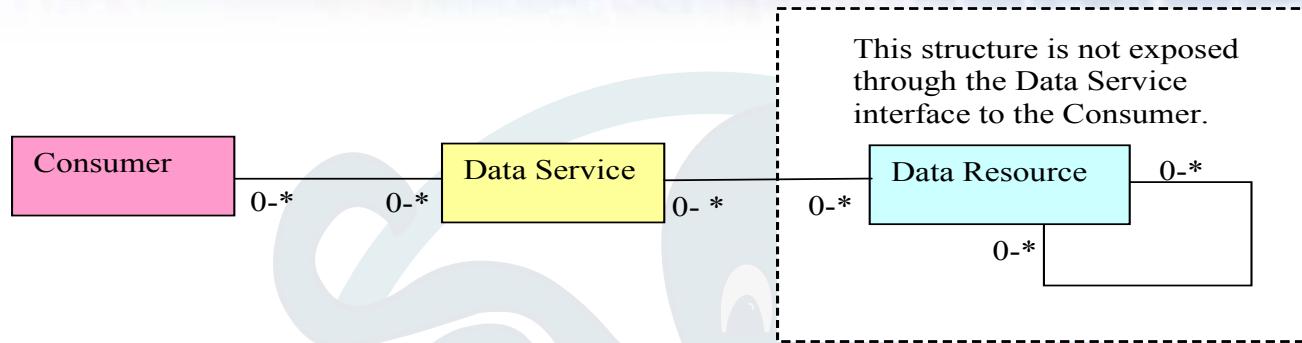
Neil Chue Hong
Project Manager, EPCC
N.ChueHong@epcc.ed.ac.uk
+44 131 650 5957

Malcolm Atkinson
Director, National e-Science Centre
mpa@nesc.ac.uk
+44 131 651 4040

- Provide service-based access to structured data resources as part of OGSA architecture
- Specify a selection of interfaces tailored to various styles of data access starting with relational and XML
- Interact well with other GGF OGSA specs

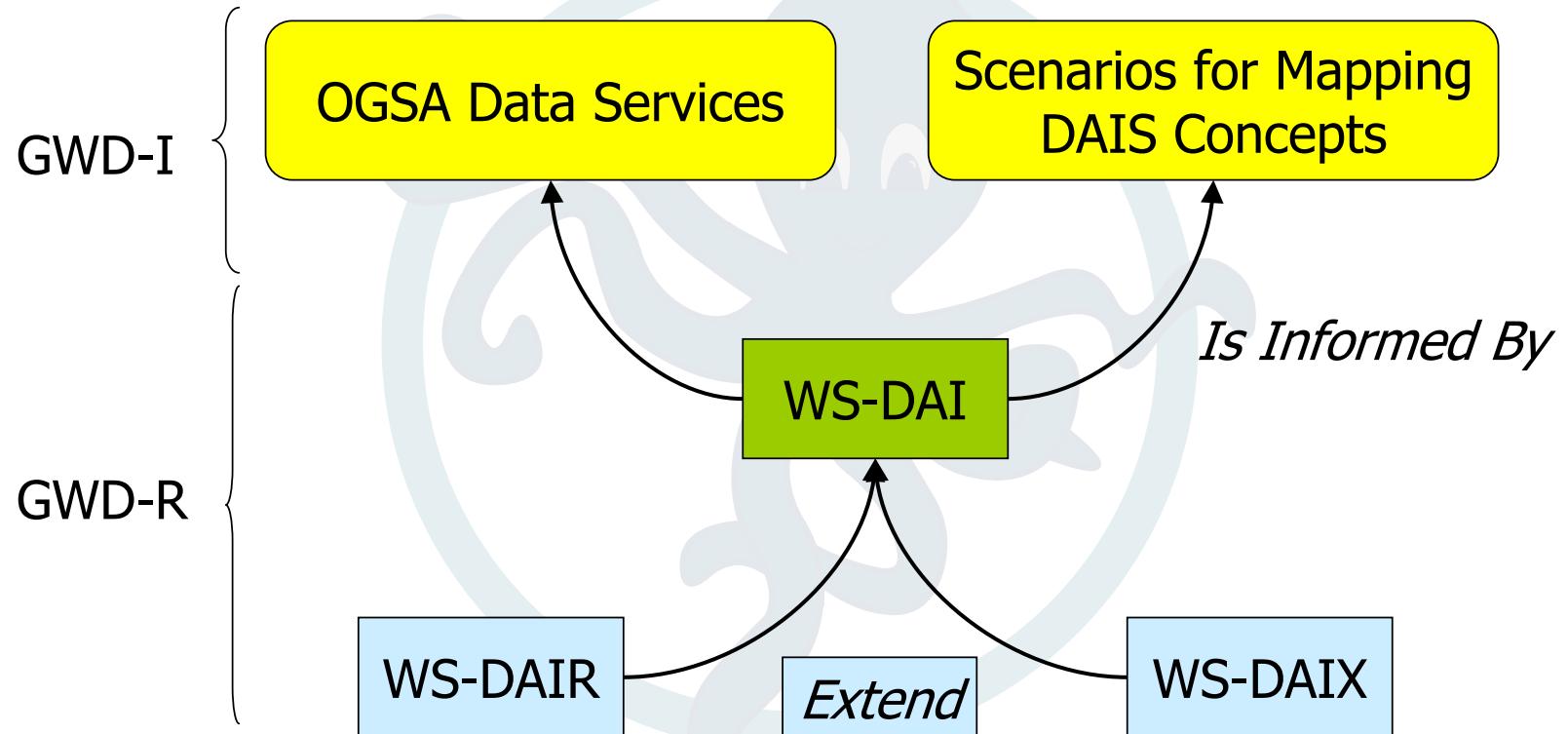
OGSA-DAI

- No new common query language
- No schema integration or common data model
- No common namespace or naming scheme
- No data resource management
 - e.g. starting/stopping database managers
- No push based delivery
 - Information Dissemination WG?



- A Data Service presents a Consumer with an interface to a Data Resource.
- A Data Resource can have arbitrary complexity, for example, a file on an NFS mounted file system or a federation of relational databases.
- A Consumer is not typically exposed to this complexity and operates within the bounds and semantics of the interface provided by the Data Service

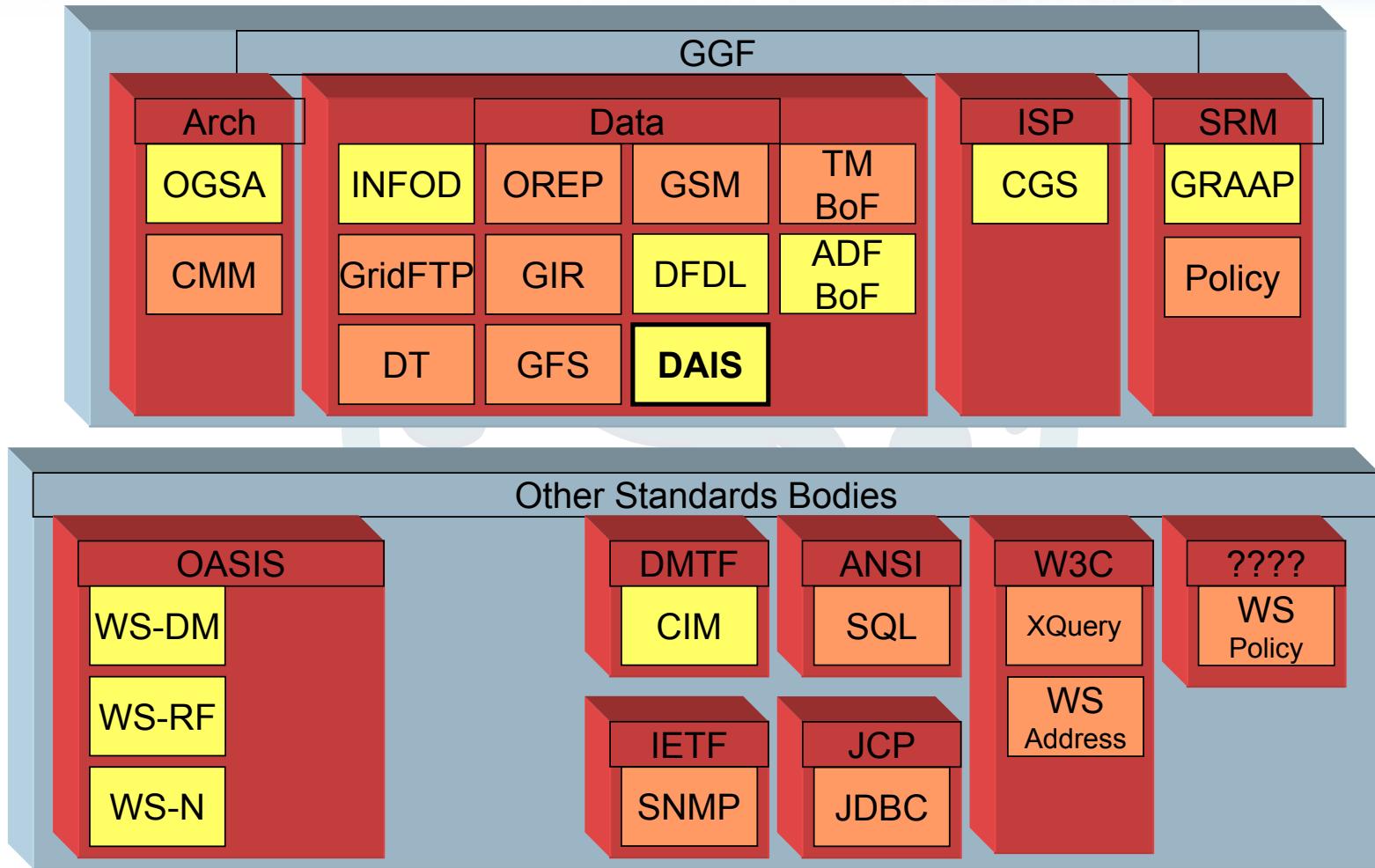
OGSA-DAI

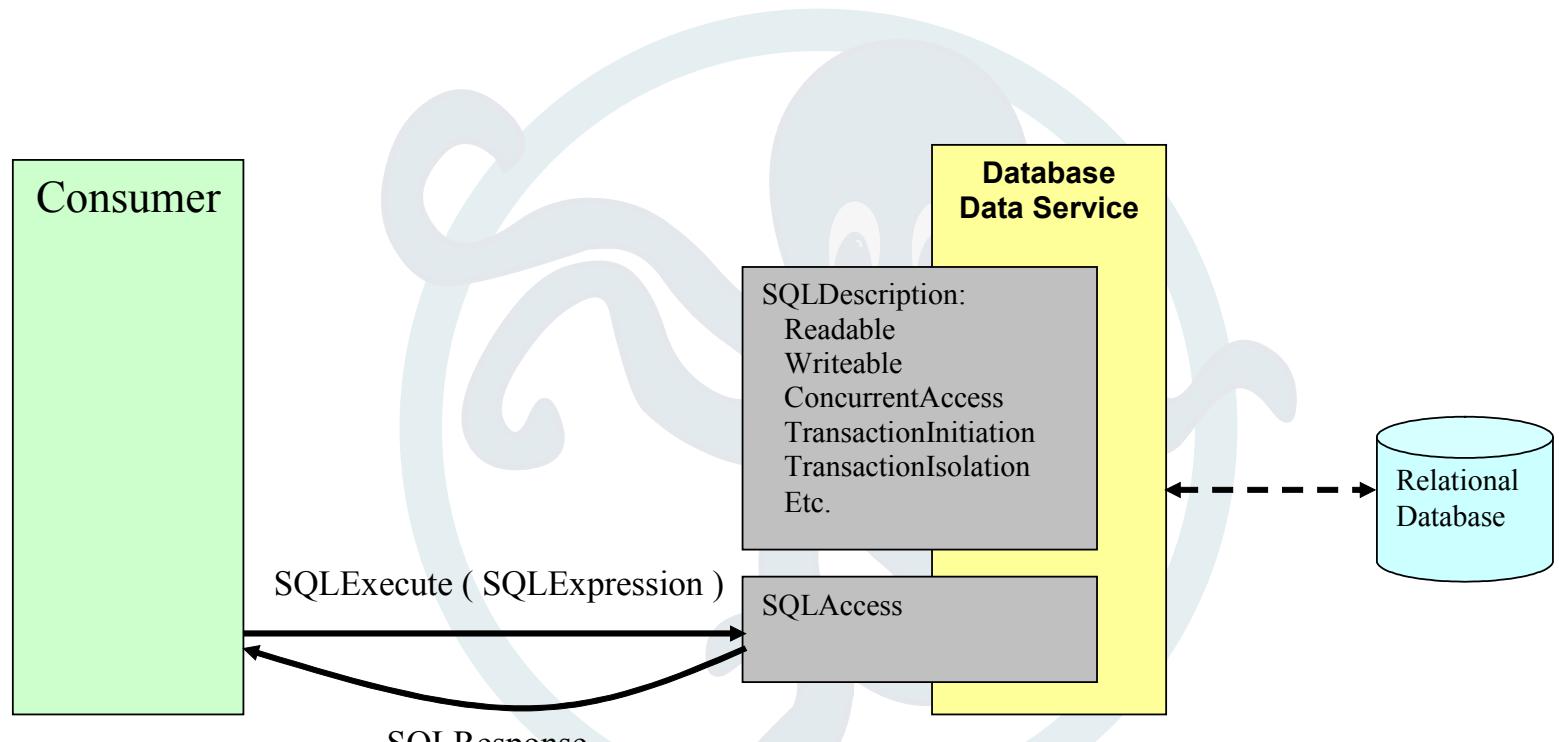


OGSA-DAI

DAIS and Other Standards/Specifications

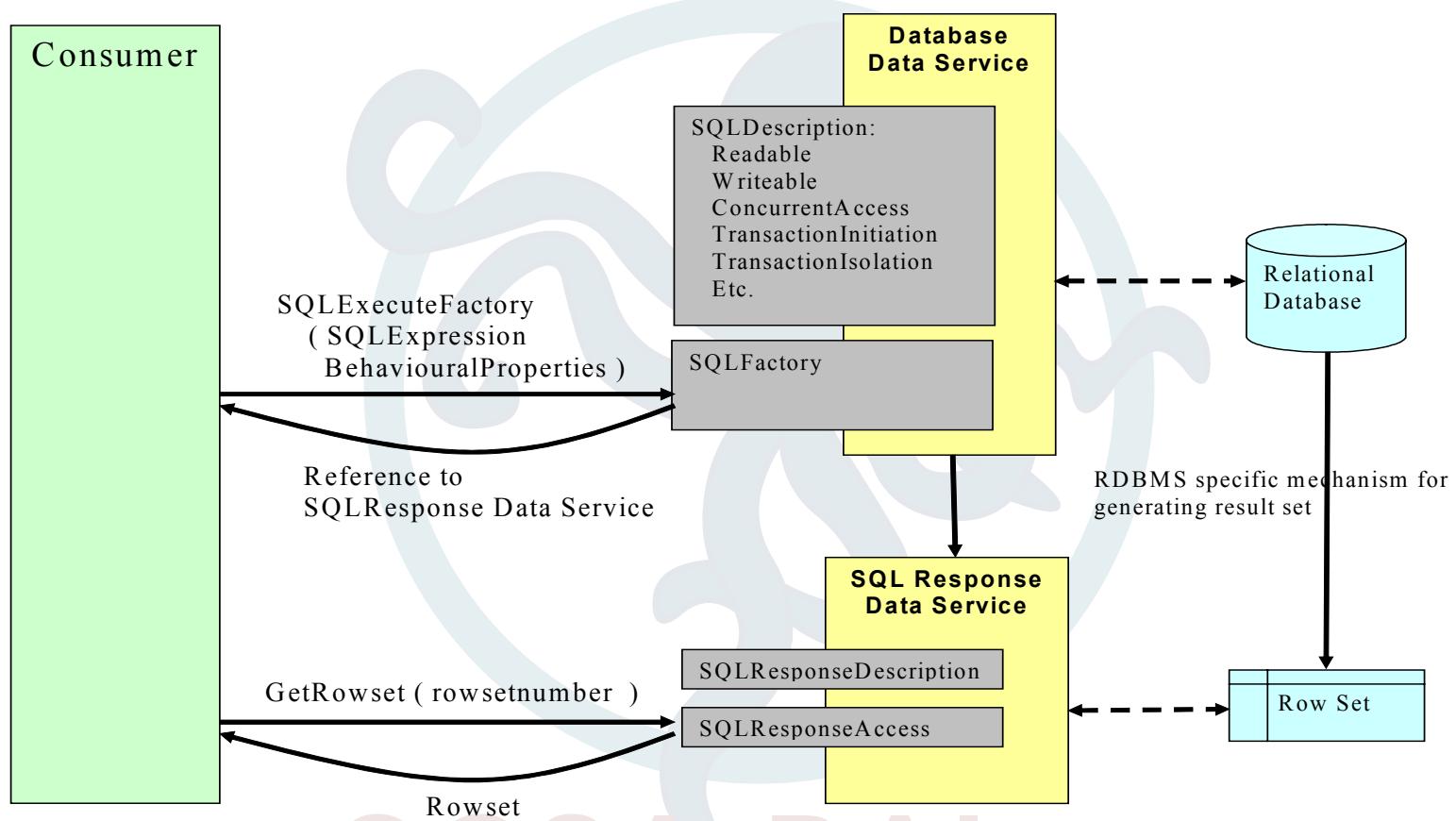
epccI





OGSA-DAI

DAIS Derived Data Access



OGSA-DAI

- Data Resource
 - Any object that can source/sink data
 - Currently databases in scope
- Data Service
 - Common interface to a data resource
 - Exposes capabilities of data resource
 - SQL Queries, X-Path Queries
 - May provide additional capabilities
 - Data transformations, 3rd party data delivery
- OGSA-DAI
 - Open Grid Services Architecture Data Access and Integration

OGSA-DAI



And now OGSA-DAI

Neil Chue Hong
Project Manager, EPCC
N.ChueHong@epcc.ed.ac.uk
+44 131 650 5957

Malcolm Atkinson
Director, National e-Science Centre
mpa@nesc.ac.uk
+44 131 651 4040

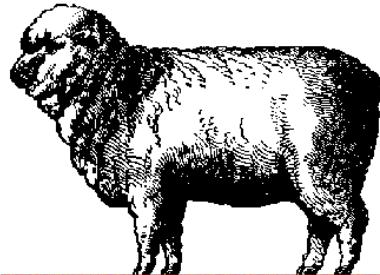
- Develop a component library
 - Access and manipulate data in a grid
 - Serve UK and International e-Science communities
- Provide an *extensible* framework
- Provide
 - Common interface to data resources
 - Simple integration of distributed queries to multiple data resources
- Contribute to standardisation efforts
 - Input into GGF DAIS WG and other groups
 - Provide a reference implementation of DAIS spec
- Based on Open Grid Services Architecture (OGSA)
 - WSRF (GT4) & WS-I (OMII_2) versions
- Support Application Developers & Contributors

OGSA-DAI

- Scale
 - Many sites, large collections, many uses
- Longevity
 - Research requirements outlive technical decisions
- Diversity
 - No “one size fits all” solutions will work
 - Primary Data, Data Products, Meta Data, Administrative data, ...
- Many Data Resources
 - Independently owned & managed
 - No common goals
 - No common design
 - Work hard for agreements on foundation types and ontologies
 - Autonomous decisions change data, structure, policy, ...
 - Geographically distributed
- and I haven't even mentioned security yet!



Slide from Neil Chue Hong



OGSA-DAI IN A NUTSHELL

A Desktop Quick Reference

With apologies to
O'REILLY®

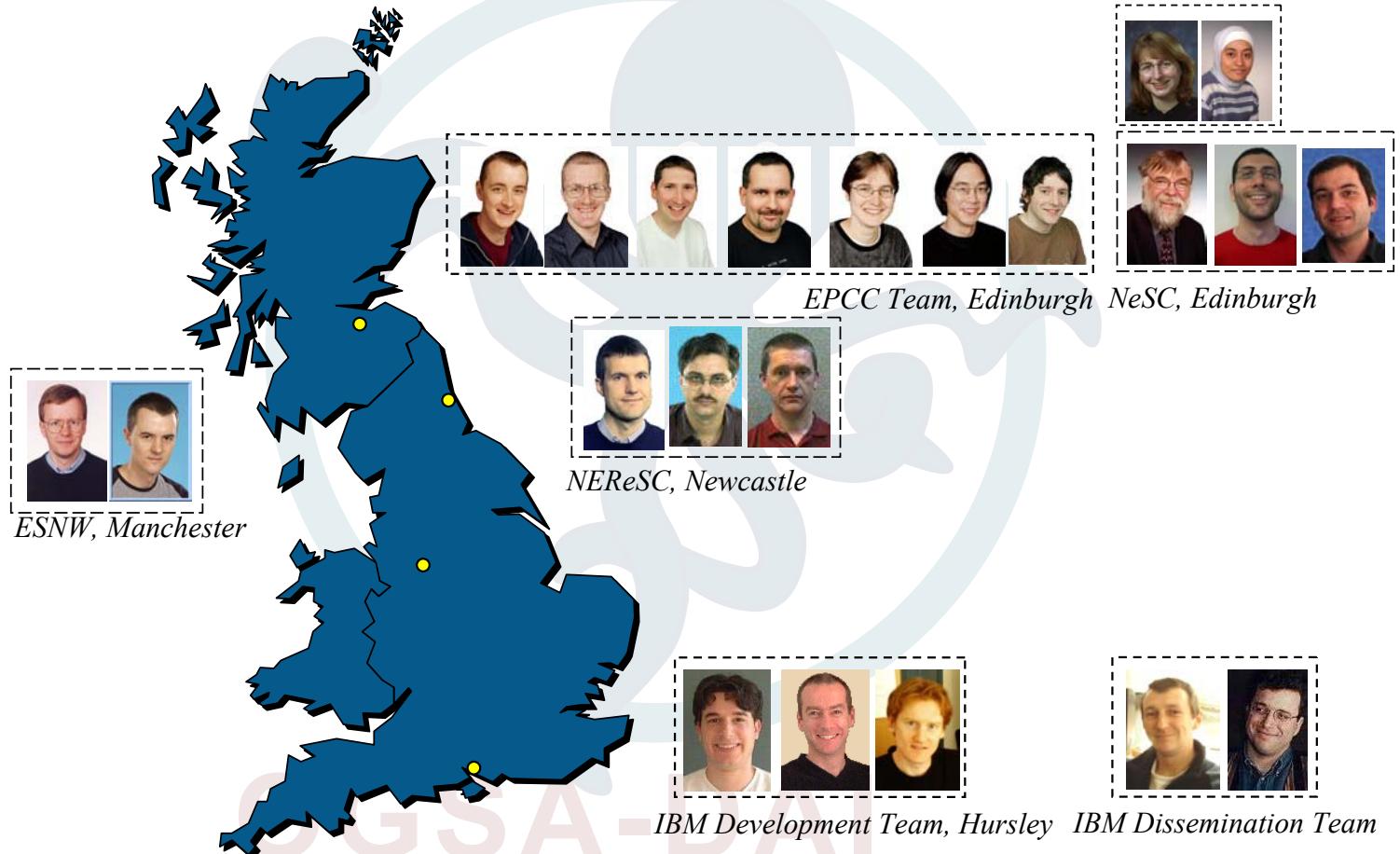
Neil Chue Hong

Slide from Neil Chue Hong

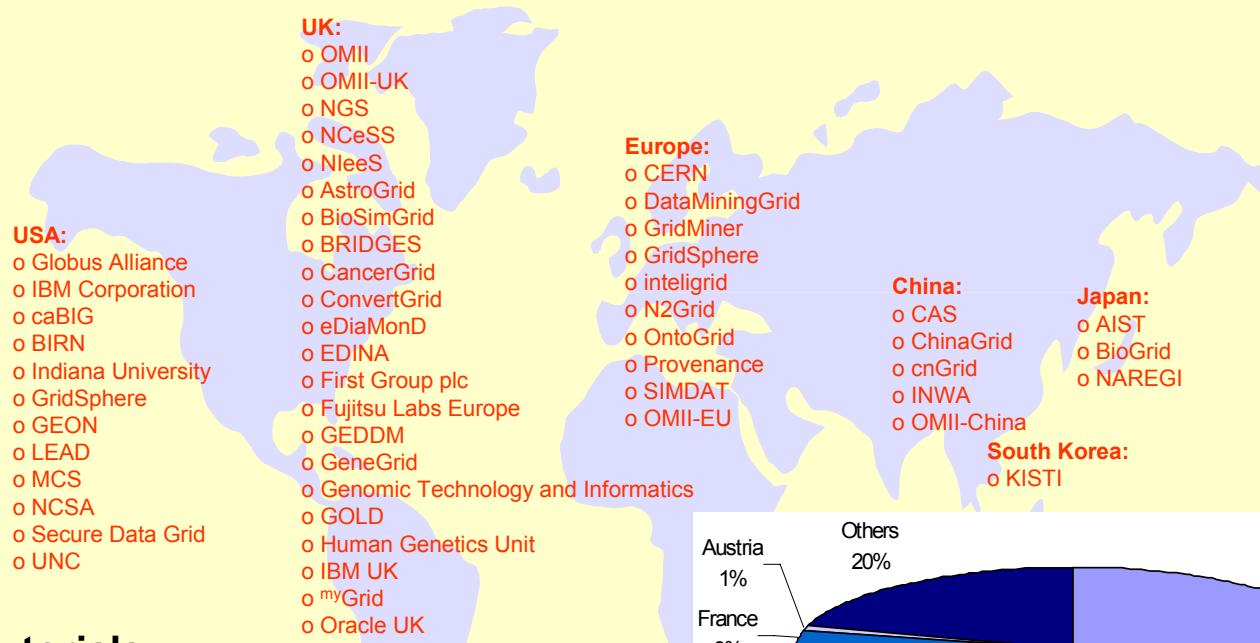
- An *extensible framework* for data access and integration.
- Expose heterogeneous data resources to a grid through web services.
- Interact with data resources:
 - Queries and updates.
 - Data transformation / compression
 - Data delivery.
- Customise for your project using
 - Additional Activities
 - Client Toolkit APIs
 - Data Resource handlers
- A base for higher-level services
 - federation, mining, visualisation,...

OGSA-DAI team

epcc|



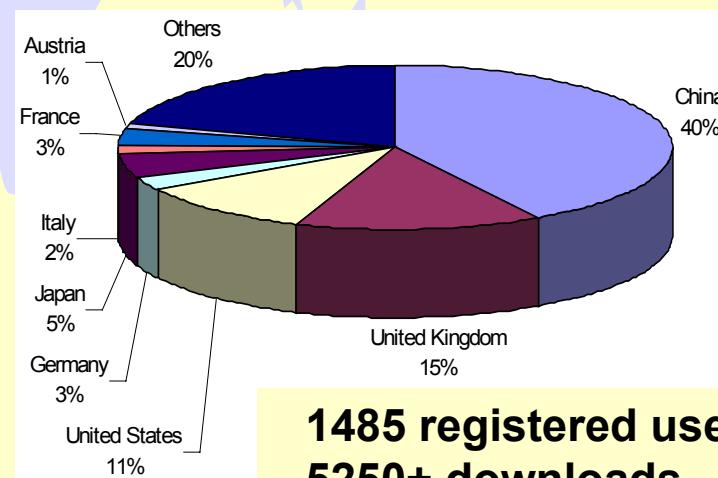
International Collaboration & Use



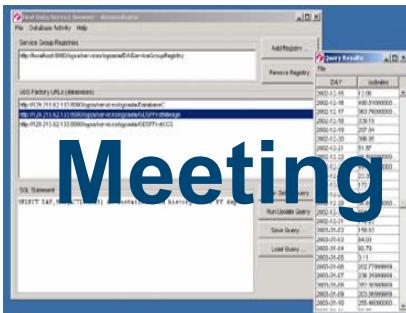
Tutorials

Boston
CERN
Edinburgh
San Francisco
Seoul
Tokyo

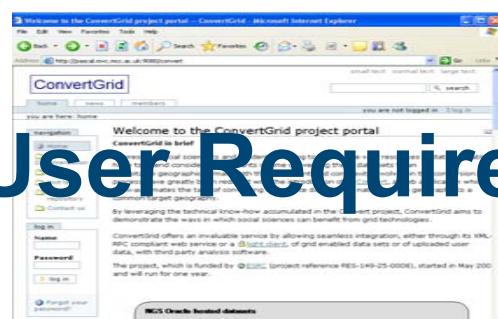
Cambridge
Chicago
London
Seattle
Singapore
ISSGC 03 to 05



1485 registered users
5250+ downloads



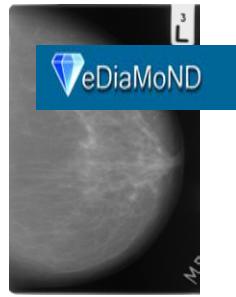
FirstDIG



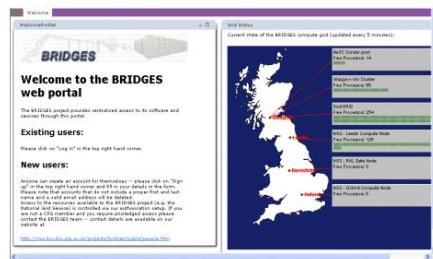
ConvertGrid



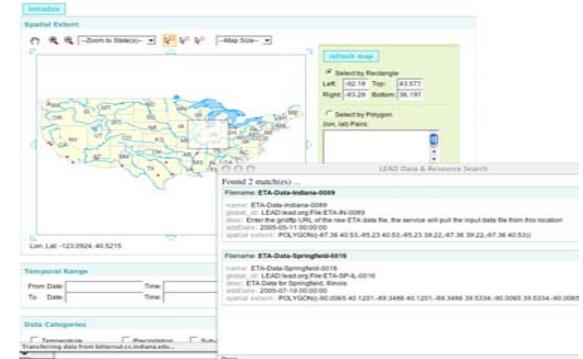
GeneGrid



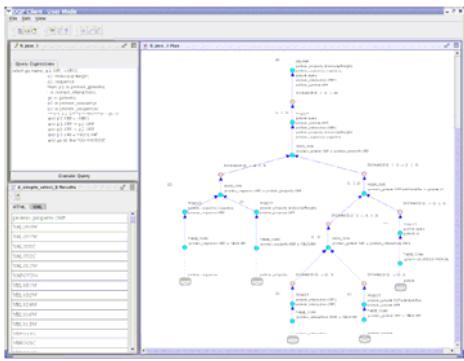
eDiaMoND



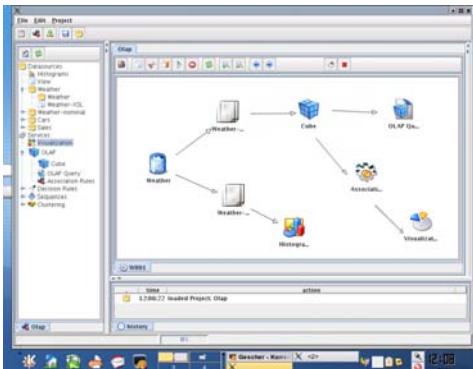
BRIDGES



LEAD



OGSA-DQP



Grid Miner

OGSA WebDB



caBIG

Project Partners

epcc|

Powered by



Funded by the Grid Core Programme

OGSA-DAI

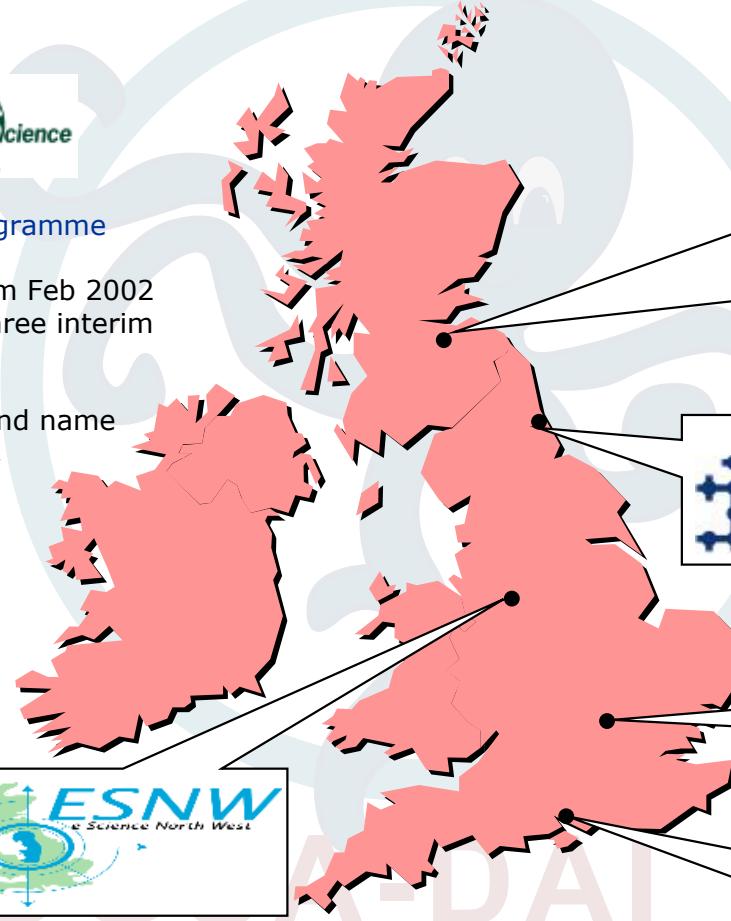
£3 million, 18 months, from Feb 2002
Three major releases, three interim releases

DAIT (DAI-Two)

Keep the OGSA-DAI brand name
£1.5 million, 24 months,
from Oct 2003
Four major releases

OMII-UK

To October 2008

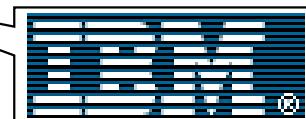


National
e-Science
Centre

epcc|



ORACLE®





Thanks for attending!

Neil Chue Hong
Project Manager, EPCC
N.ChueHong@epcc.ed.ac.uk
+44 131 650 5957

Malcolm Atkinson
Director, National e-Science Centre
mpa@nesc.ac.uk
+44 131 651 4040