# Science Informatics

Dan Fay
Director – Earth, Energy and Environment
dan.fay@microsoft.com

# MSR eScience Workshop
# Looking Back 8 yrs to the Beginning

## Scientific Data Intensive Computing Workshop 2004

- Keynote: *20 Questions to a Better Application* – Jim Gray

  ***Online Science the New Computational Science***

- Talk: *Data Explosion: Astrophysics with Terabytes of Data*
  - Alex Szalay
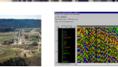
# Online Science the New Computational Science

## Information Avalanche

- In science, industry, government,....
  - better observational instruments and
  - and, better simulations producing a data avalanche
- Examples
  - BaBar: Grows 1TB/day
    - 2/3 simulation Information
    - 1/3 observational Information
  - CERN: LHC will generate 1GB/s .~10 PB/y
  - VLBA (NRAO) generates 1GB/s today
  - Pixar: 100 TB/Movie
- **New emphasis on informatics:**
  - **Capturing, Organizing, Summarizing, Analyzing, Visualizing**

*BaBar, Stanford*
*P&E Gene Sequencer From http://www.pangene.ucl.edu*
*Space Telescope*

## Publishing Data

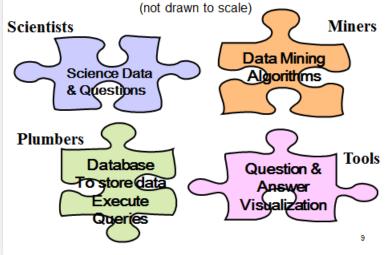| Roles | Traditional | Emerging |
|---|---|---|
| Authors | Scientists | Collaborations |
| Publishers | Journals | Project www site |
| Curators | Libraries | Bigger Archives |
| Consumers | Scientists | Scientists |

- Exponential growth:
  - Projects last at least 3-5 years
  - Data sent upwards only at the end of the project
  - Data will **never** be centralized
- More responsibility on projects
  - Becoming Publishers and Curators
  - Often no explicit funding to do this **(must change)**
- Data will reside with projects
  - Analyses must be close to the data (see later)
- Data cross-correlated with Literature and Metadata [1]

## Global Federations

- Massive datasets live near their owners:
  - Near the instrument's software pipeline
  - Near the applications
  - Near data knowledge and curation
- Each Archive publishes a (web) service
  - Schema: documents the data
  - Methods on objects (queries)
- Scientists get "personalized" extracts
- Uniform access to multiple Archives
  - A common global schema

**Federation**

## What's X-info Needs from us (cs)
### (not drawn to scale)

Scientists — Science Data & Questions

Miners — Data Mining Algorithms

Plumbers — Database To store data Execute Queries

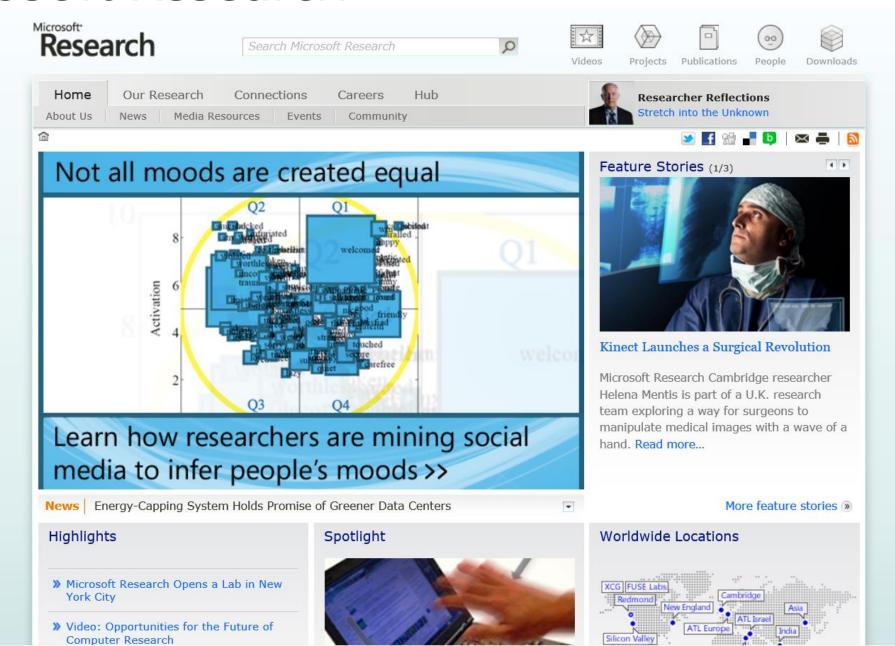Tools — Question & Answer Visualization

9

## How to Help?

- Can't learn the discipline before you start (takes 4 years.)
- Can't go native – you are a CS person not a bio,… person
- Have to learn how to communicate Have to learn the language
- Have to form a working relationship with domain expert(s)
- Have to find problems that leverage your skills

28

## Call to Action
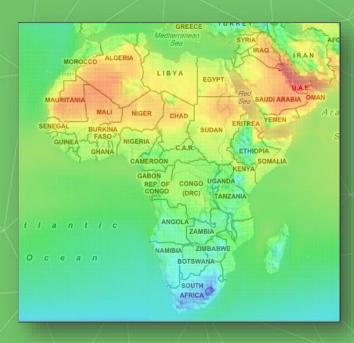
- X-info is emerging.
- Computer Scientists can help in many ways.
  - Tools
  - Concepts
  - Provide technology consulting to the commuity
- There are great CS research problems here
  - Modeling
  - Analysis
  - Visualization
  - Architecture

46

# Microsoft Research
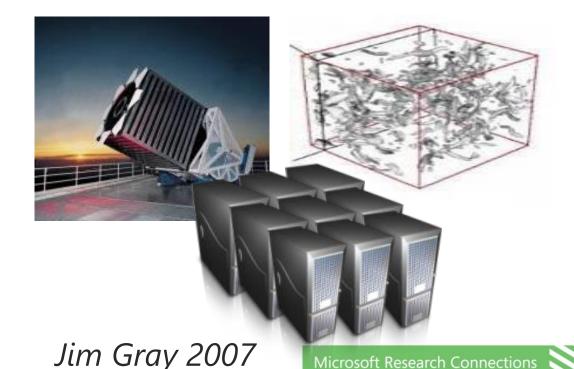
# Earth, Energy & Environment



- **Visualizing and Experiencing E³ Data + Information**: Provide a unique experience to reduce time to insight and knowledge through visualizing data and information
- **Accessible Data:** Ensure $E^3$ data (remote and local sensing) is easily discoverable, accessible and consumable in the scientists domain
- **Enabling Scientific Collaboration**: Look at new ways to enable collaboration in scientific virtual organizations

# Emergence of a Fourth Paradigm

- Thousand years ago – **Experimental Science**
  - Description of natural phenomena
- Last few hundred years – **Theoretical Science**
  - Newton's Laws, Maxwell's Equations...
- Last few decades – **Computational Science**
  - Simulation of complex phenomena
- Today – **Data-Intensive Science**
  Scientists overwhelmed with data sets
         from many different sources
    - Data **captured by instruments**
    - Data **generated by simulations**
    - Data **generated by sensor networks**
- eScience is the set of tools and technologies
         to support data federation and collaboration
    - For analysis and data mining
    - For data visualization and exploration
    - For scholarly communication and dissemination

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - \mathrm{K}\frac{c^2}{a^2}$$

*Jim Gray 2007*

# Changing Nature of Discovery

- Complex models
  - Multidisciplinary interactions
  - Wide temporal and spatial scales
- Large multidisciplinary data
  - Real-time steams
  - Structured and unstructured
- Distributed communities
  - Virtual organizations
  - Socialization and management
- Diverse expectations
  - Client-centric and infrastructure-centric



The FOURTH PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

http://fourthparadigm.org

# The Problem for the e-Scientist

## How to codify and represent our knowledge

Experiments & Instruments → facts →

Simulations → facts →

Literature → facts →

Other Archives → facts →

← Questions

Answers →

### The Generic Problems

- Data ingest
- Managing a petabyte
- Common schema
- How to organize it
- How to *re*organize it
- How to share with others

- Query and Vis tools
- Building and executing models
- Integrating data and Literature
- Documenting experiments
- Curation and long-term preservation

*(With thanks to Jim Gray)*

# EOS Article: *Mountain Hydrology, Snow Color, and the Fourth Paradigm* by Jeff Dozier





Snow is one of nature's most colorful materials
(Landsat Thematic Mapper snow & cloud)



Spatially distributed snow water equivalent

SWE, mm

4500
2500
1900
1300
600
10

(N. Molotch)

04/10/05

# Information about water is more useful as we climb the value ladder



Data >>> Information >>> Insight

>>> Increasing value >>>

Forecasting

Reporting

Analysis

Integration

Distribution

Aggregation

Quality assurance

Collation

Monitoring

**Done poorly,
but a few notable
counter-examples**

**Done poorly to moderately,
not easy to find**

**Sometimes done well,
generally discoverable and available,
but could be improved**

(I. Zaslavsky & CSIRO, BOM, WMO)

Microsoft Research Connections

# Environmental Ecosystem



**Knowledge**

**Action**

**Inform**

# Environmental Ecosystem

# Information ecosystem:

It is chaotic, unstructured and ad hoc

# The Ecological Data Flood



- We're living in a perfect storm of remote sensing, cheap ground-based sensors, internet data access, and commodity computing
- Yet deriving and extracting the variables needed for science remains problematic
  - Specialized knowledge for algorithms, internal file formats, data cleaning, etc, etc
  - Finding the right needle across the distributed heterogeneous and very rapidly growing haystacks

# Data Variety – The Spice of Life

Manual Measurement

Automated Measurement

Sample Collection

Typing

Counting

Historical Photographs

Satellite

Aircraft Surveys

Model Output

Relatively Ubiquitous Motes

Sea ice concentration (%)

Sea surface temperature (deg C)
-2 0 2 4 6 8 10 12 14 16 18 20 22 24 26 28 30 32

# Data Integration Challenges

- Regular rasters, points, and spatial features
- Time series and intermittent
- Vocabulary meanings (ontology)
- Sparse in time, duration, or location
- Science variable derivation
- Gaps
- Spatial/temporal harmonization

# Why Make this Distinction?


PB

Provenance and trust widely varies

Data acquisition, early processing, and reporting ranges from a large government agency to individual scientists.

Smaller data often passed around in email; big data downloads can take days (if at all)

Data sharing concerns and patterns vary

Open access followed by (non-repeatable and tedious) pre-processing


KB


TB

True science ready data set but concerns about misuse, misunderstanding particularly for hard won data.

Computational tools differ.

Not everyone can get an account at a supercomputer center

Very large computations require engineering (error handling)

Space and time aren't always simple dimensions


GB

Complex shared detector                                              Simple instrument (if any)

*Science happens when PBs, TBs, GBs, and KBs can be mashed up simply*

Complex and Heavy process by experts                                 Ad hoc observations and models

# AzureMODIS – Azure Service for Remote Sensing Geoscience

- Science pipeline for download, initial processing, and reduction of satellite imagery. Developed by MSR, UVa, UCB.
- Dramatically lowers resource and complexity barriers to use satellite imagery for terrestrial hydrology and geoscience.
  - Common imagery location determination and upload from diverse sources
  - Optional scientist-provided reduction algorithm (.NET, Java, or MatLab)
  - On-demand scalability beyond local desktop or cluster
- In use now to compute 10 year continental scale water balance for North America. Per year:
  - 500 GB  (~60K files) upload of 9 different source imagery products from 15 different locations
  - 400 GB reprojected harmonized imagery consuming  ~3500 cpu hours
  - 5 GB reduced science result leveraging reported field data aggregates consuming ~60 cpu hour

**Source Imagery Download Sites**

Source Metadata

**Data Collection Stage**

Request Queue

AzureMODIS
Service Web Role Portal

**Reprojection Stage**

Analysis/Reduction Stage
Reduction Queue

Scientist

Scientific Results

Catharine van Ingen (Microsoft Research), Jie Li, Marty Humphreys (UVA), Youngryel Ryu (UCB), Deb Agarwal (BWC/LBL)

Microsoft Research Connections

**Microsoft Codename "Cloud Numerics"**

numerical and data analytics library for data scientists, quantitative analysts, and others

**Microsoft Codename "Data Explorer"**

organize, manage, mash up and gain new insights from your data.

**Microsoft Codename "Data Hub"**

Organizationsto curate and publish its data on a private data marketplace

**Microsoft Codename "Trust Services"**

data encryption services for cloud applications so that they can roam encryption keys in a secure way.

Discoverability

Accessibility

Consumeability

# Fetch Climate



## FetchClimate

Inspiration | Features | Try online | Download | Case studies | People | Acknowledgements

### Retrieve climatic and environmental information with the click of a button or a few lines of code

FetchClimate is a fast, free, intelligent climate information service that operates over the cloud to return exactly the information you need. FetchClimate can be accessed either through a simple web interface, or via a few lines of code inside any .NET program. FetchClimate is intended to make it easy for you to retrieve information for any geographical region, at any grid resolution: from global, through continental, to a few kilometres, and for any range of years (1900 – 2010), days within the year, and / or hours within the day. FetchClimate can also report the uncertainty associated with the values it returns and list data sources used to fulfil the request. When multiple sources could potentially provide information on the same environmental variable, FetchClimate automatically selects the most appropriate data source. Finally, the entire query you ran can be shared as a single URL, enabling others to retrieve the identical information.

### Inspiration

Many environmental science research and applications require information about climate and

### Documentation

() Manual & examples for web version [v

() Manual & examples for programmatic
[pdf]

http://fetchclimate.cloudapp.net/

# Environmental Informatics Framework (EIF)

## Common Problems with Data

➤ To use data from different sources
  - Non-standard formats, scales, and units
  - Lack of data quality control
  - Lack of metadata
  - Difficult to repurpose data for different (my) tools

➤ To share data
  - Lack of incentive (no credit)
  - Need extra resources and tools

➤ Hidden problems, seldom addressed
  - Versioning
  - Provenance
  - Curation

**Data Sources**

SQL

XML

(data)

Data Cube

CSV

# Environmental Informatics Framework (EIF)

## Current State of Data Ecosystem

# Environmental Informatics Framework (EIF)

**Advance data discoverability, accessibility, and consumability**

## Open Data Protocol (OData)
http://www.odata.org

**It allows you to form URLs based on what you know about the underlying data**

➢ A Web protocol for querying and updating data
  ❑ provides a way to unlock your data and free it from data silos
  ❑ does this by building upon Web technologies such as HTTP, Atom Publishing Protocol (AtomPub) and JSON to provide access to information from a variety of applications, services, and stores.

➢ In Open Source/Specifications Promise – being submitted to OASIS

➢ An application of a set of internet standards:
  ❑ HTTP,
  ❑ Atom (RFC 4287),
  ❑ AtomPub (RFC 5023),
  ❑ REST semantics

➢ Existing standards + easy data access API

➢ Added **Geospatial data support** –
  ❑ Feedback from the Community encouraged – www.odata.org

# New ways to analyze and communicate data

# The 'Cosmic Genome' Project

- The Sloan Digital Sky Survey is the first major astronomical survey project:
  - 5 color images of ¼ of the sky
  - Pictures of 300 million celestial objects
  - Distances to the closest 1 million galaxies
- Jim Gray from Microsoft Research worked with astronomer Alex Szalay to build the public 'SkyServer' archive for the survey
- New model of scientific publishing
  - Have to publish the data <u>before</u> astronomers publish their analysis

# Public Use of the SkyServer

- **Posterchild in 21st century data publishing**
  - 380 million web hits in 6 years
  - 930,000 distinct users
    vs 10,000 astronomers
  - 1600 scientific papers
  - Delivered 50,000 hours
    of lectures to high schools
  - Delivered 100B rows of data

- **Citizen Science: GalaxyZoo**
  - Goal of 1 million visual galaxy
    classifications by the public
  - Allows general public to search for
    photographs and classify
    different types of galaxies

# World Wide Telescope

Seamless Rich Social Media Virtual Sky and E...
Web application for science and educatio...

Goals
- Integra...                                    ...ess
- Easy ac...
- Tours f...
- Spatial...

Updat...
- AP...
- Ex...

Not ju...
Being...

We inv...

**www.layerscape.org**

- Community Site for WWT Tours and Layers (Data)
- Sharing by groups/individuals

# Excel to Visualize in seconds

# Natural User Interfaces (NUI)
# Kinect SDK and WWT

- Rethinking ways in which people will interact with computers/technologies of the future

- Re-evaluating everything from their (non-) physical design to the human needs and interaction models

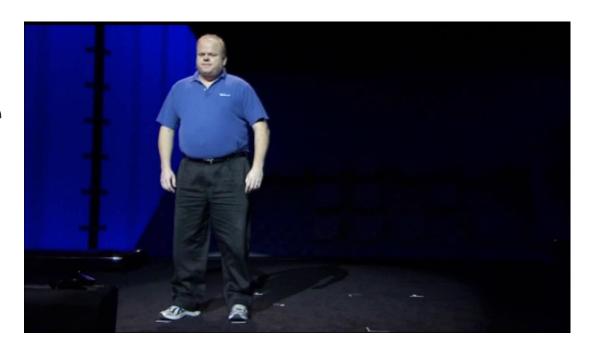- Revolutionize the way we think about technology and what it can do on our behalf

# Our Next Steps

- Continue to work with Scientists
  - Implement features from feedback
    - Ie. Netcdf support
  - New clients and User Interface
  - Make it easier to use/create tours

- Cloud: Windows Azure integration - Write, Run, or Use Software
  - Write: [Platform as a Service](#)
    - Web Sites
    - Storage
    - Languages - .Net, Java, Python, PHP, Ruby, NodeJS, C++
  - Run: [Infrastructure as a Service](#)
    - Virtual Machines – Windows Servers & Linux
    - Connect to on-premises
  - Use: [Software as a Service](#)
    - Azure Media Services
    - HPC and Big Data – Hadoop, etc
    - Windows Azure Marketplace

# eScience in Action

# Microsoft eScience Workshop 2012

**October 8–9, 2012 | Chicago, Illinois**
http://research.microsoft.com/events/escience2012

Microsoft Research Connections

# Data Storage Sustainability?

- Digital Data can be open – who should pay the cost?
- Spinning Disks, Bandwidth, Cooling, etc

# No Silver Bullet - What is needed?

- Algorithms that scale
- Data Management from the Start
- Automatic Ancillary Data capture
- Thinking about the Data, and retention
- Data sharing is natural from the start
- Visualization for everyone
- Best practices – insights and challenges shared amongst domains
    - Ie. eScience Workshop, etc

# Challenges

- Balancing

  Data Acquisition | Bandwidth | Storage/Processing

- Cross Discipline Collaboration – Knowledge sharing
- The data deluge - How to manage and analyze information?
- New types of Scientists:
  - Data Collectors & Data Analysis
- Riding the commodity curve
- Technology/Computing in support of Science

Microsoft®

**Research** Connections