# Advancing Environmental Understanding:
## the Role of eScience

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

# MSR eScience Workshop 2011

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**

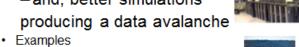- Talk: *Data Explosion: Astrophysics with Terabytes of Data*

  - 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**

*Image courtesy C. Meneveau & A. Szalay @ JHU*

*BaBar, Stanford*

*P&E Gene Sequencer From http://www.genome.ucl.edu*

*Space Telescope*

## Publishing Data

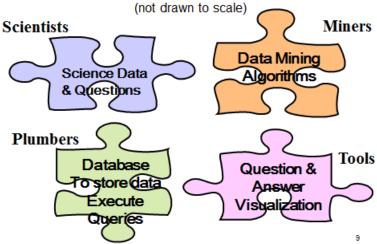| Roles | Traditional | Emerging |
|-------|-------------|----------|
| Authors | Scientists | Collaborations |
| Publishers | Journals | Project www site |
| Curators | Libraries | Bigger Archives |
| Consumers | Scientists | Scientists |

- 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**

**Miners**

Science Data & Questions

Data Mining Algorithms

**Plumbers**

Database To store data Execute Queries

Question & Answer Visualization

**Tools**

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 community
- There are great CS research problems here
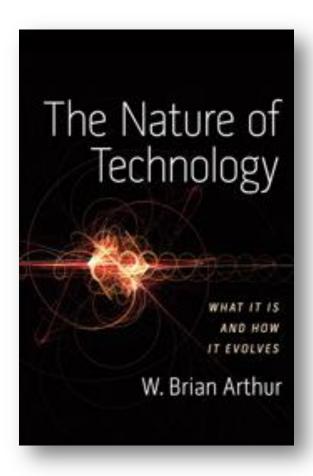  - Modeling
  - Analysis
  - Visualization
  - Architecture

46
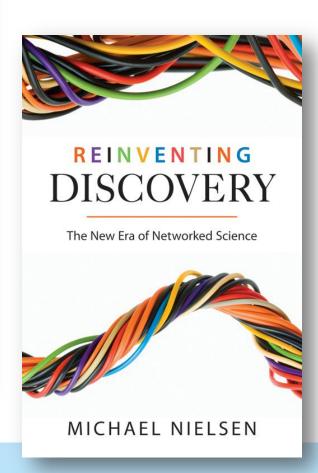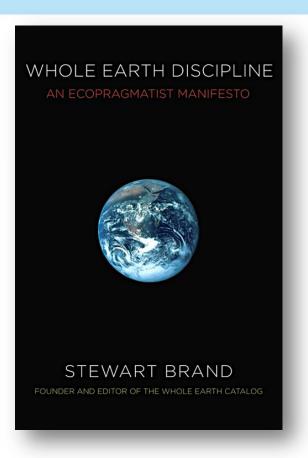
# A Tidal Wave of Scientific Data

# Interesting Thinking



The Nature of Technology
WHAT IT IS AND HOW IT EVOLVES
W. Brian Arthur

REINVENTING DISCOVERY
The New Era of Networked Science
MICHAEL NIELSEN

WHOLE EARTH DISCIPLINE
AN ECOPRAGMATIST MANIFESTO
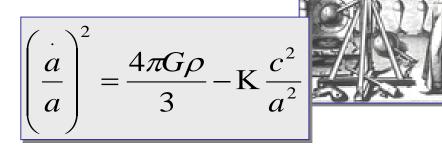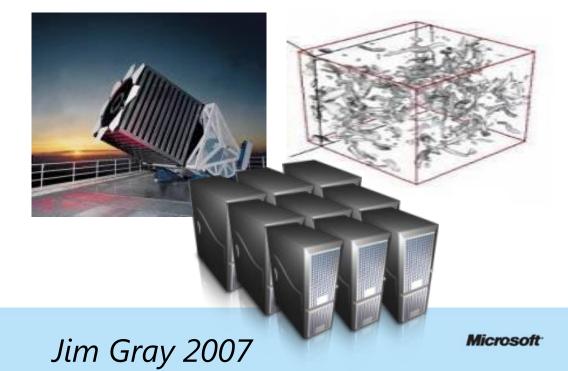STEWART BRAND
FOUNDER AND EDITOR OF THE WHOLE EARTH CATALOG

# 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
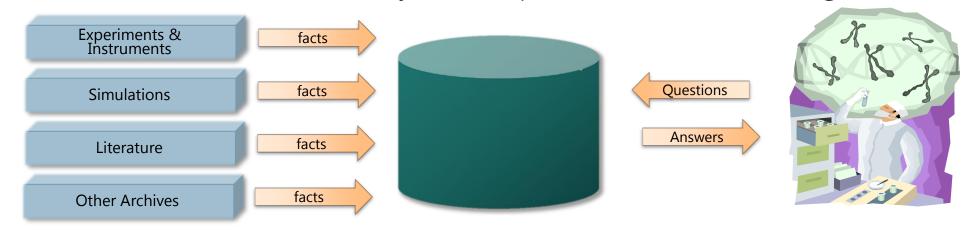  - Client-centric and infrastructure-centric

The FOURTH PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

http://fourthparadigm.org

Microsoft

# 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)*

Microsoft

# What is eScience?

Definition of eScience (Wikipedia)

**E-Science** (or **eScience**) is computationally intensive <u>science</u> that is carried out in highly distributed <u>network</u> environments, or science that uses immense <u>data</u> sets that require <u>grid computing</u>; the term sometimes includes technologies that enable distributed collaboration

Definition from Microsoft Research

How computing technologies can help address scientific challenges
eScience efforts at Microsoft Research seek to further the understanding of these challenges, support the developing community, develop computational tools that will enable the advancement of scientific research, and catalyze discovery through funded collaborative research.

What it really means

How can current and future products and technologies can be applied to scientific challenges to help with scientific insight in a easy to use system

Technology in support of Science

# 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



Forecasting

Reporting

Analysis

Integration

Distribution

Aggregation

Quality assurance

Collation

Monitoring

Data >>> Information >>> Insight

>>> Increasing value >>>

**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**

# Environmental Ecosystem



Knowledge

Action

Inform

# Environmental 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

Historical Photographs

Typing

Counting

Model Output

Relatively Ubiquitous Motes

Aircraft Surveys

Satellite

# 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

# Tiling: Do Scientists Have to be Computer Scientists?

- Reprojection
  - Converts one geo-spatial representation to another.
  - Example is converting from latitude-longitude swaths to sinusoidal cells.
- Spatial resampling
  - Converts one spatial resolution to another.
  - Example is converting from 1 KM to 5 KB pixels.
- Temporal resampling
  - Converts one temporal resolution to another.
  - Example is converting from daily observation to 8 day averages.
- Gap filling
  - Assigns values to pixels without data either due to inherent data issues such as clouds or missing pixels introduced by one of the above.
- Masking
  - Eliminates uninteresting or unneeded pixels.
  - Examples are eliminating pixels over the ocean when computing a land product or eliminating pixels outside a spatial feature such as a watershed.

**Source Data (Swath format)**

Surface Temperature (MODIS Swaths) DOY 190, 2003

**Reprojected Data (Sinusoidal format)**

# Why Make this Distinction?

**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

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

**PB**

**KB**     **TB**

**GB**

Complex shared detector                                Simple instrument (if any)

*Science happens when PBs, TBs, GBs, and KBs can be mashed up simply*
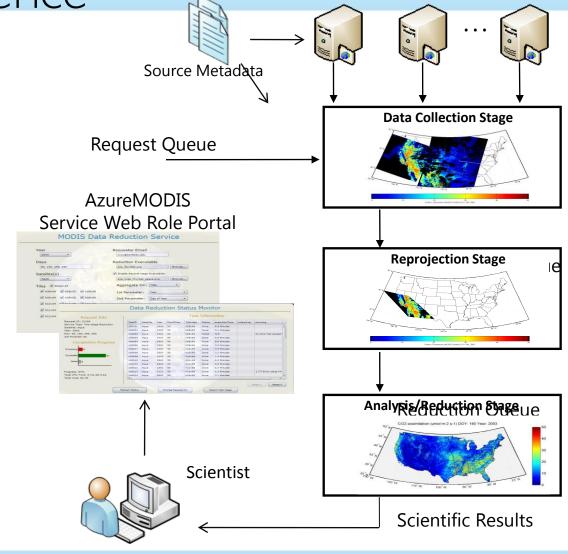
Complex and Heavy process by experts                   Ad hoc observations and models

# 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

Request Queue

AzureMODIS
Service Web Role Portal

**Data Collection Stage**

**Reprojection Stage**

**Analysis/Reduction Stage**

Scientist

Scientific Results

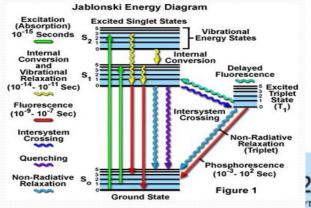# Micrometeorology

Pilot study "R1": 2009
- 20 million observations
- Engineering success

Collaborators:
- Humberto da Rocha (USP)
- Andreas Terzis (JHU)
- Juliana Salles, Rob Fatland  (MSR)
- Brito Cruz (FAPESP)

$$\frac{\partial \bar{c}}{\partial t} + \bar{u}\frac{\partial \bar{c}}{\partial x} + \bar{w}\frac{\partial \bar{c}}{\partial z} + \bar{c}\left(\frac{\partial \bar{u}}{\partial x} + \frac{\partial \bar{w}}{\partial z}\right) + \frac{\partial \overline{u'c'}}{\partial x} + \frac{\partial \overline{w'c'}}{\partial z}$$

Continuation "R2": 2012+
- Science and engineering objectives
- *"Solve the carbon balance problem"*
- *"Build an interoperable data system"*

# Biogeochemistry: Carbon+

First Objective: Characterize fate of terrigenous carbon
- Multiple spectral analysis methods
- Data reduction: From correlation to machine learning

- Second objective: Library
  - Follow Environmental Information Framework
  - Contribute merit to Data Publishers
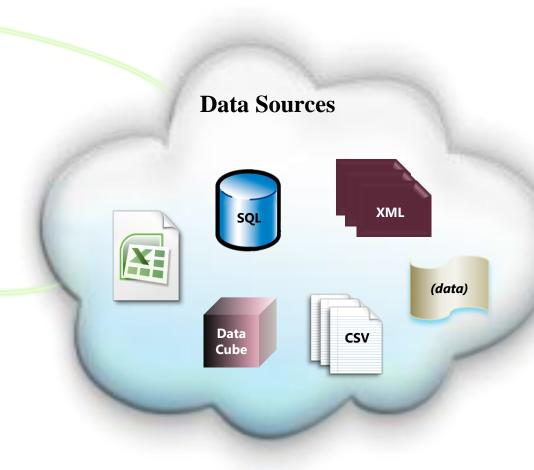  - Address { curation, versioning, provenance }

# Environmental Informatics Framework (EIF)

## Common Problems with Data

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

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

➢ Hidden problems, seldom addressed
- o Versioning
- o Provenance
- o Curation

**Data Sources**

SQL

XML

(data)

Data Cube

CSV

# Environmental Informatics Framework (EIF)

## Current State of Data Ecosystem

# Environmental Informatics Framework (EIF)

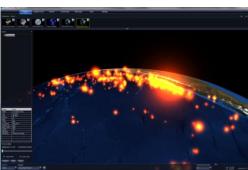## Advance data discoverability, accessibility, and consumability

### 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

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

➢ Existing standards + easy data access API

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

# "Data Explorer"

- A new SQL Azure Lab that provides a new way to organize, manage, mashup and gain new insights from your data.

- Data Explorer provides capabilities for data curation, collaboration, classification and mashup, opening new capabilities and opportunities around the data that you own or want to work with.

**DISCOVER**

**ENRICH**

**PUBLISH**

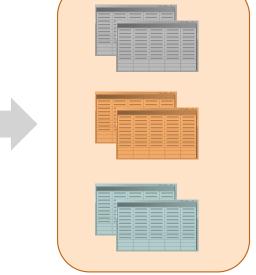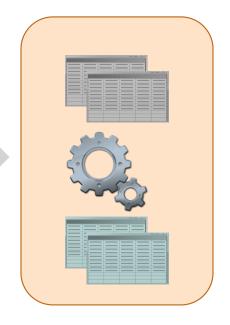# Microsoft Codename "Data Explorer" capabilities

**DISCOVER** → **ENRICH** → **PUBLISH**



| | | |
|---|---|---|
| "OTHER" DATA | | |
| PUBLIC REF. DATA | | |
| LICENSED REF. DATA | | |
| FILES & SHARES | | |
| ON-PREM RDBMS's | | |

**Add & Manage Data Sources**

**Classify
Understand
Recommend**
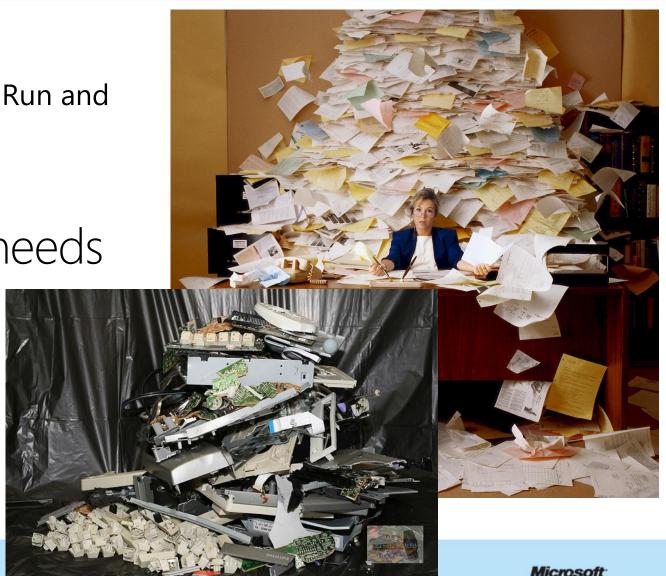
**Transform
Mash up
Cleanse**

**Snapshot
Publish
Sell**

# Data Tsunami or Scientific Data Hording

- ## Technology Trends
  - Compute and Storage make it easy to Run and Keep
  - Does it get used?  Could it be mined?

- ## Data Collaboration reuse needs
- ## Discoverability
  - Catalogs, etc
- ## Accessibility
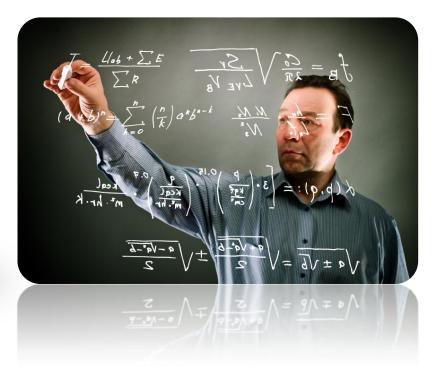  - Protocols
- ## Consumability
  - Tooling support
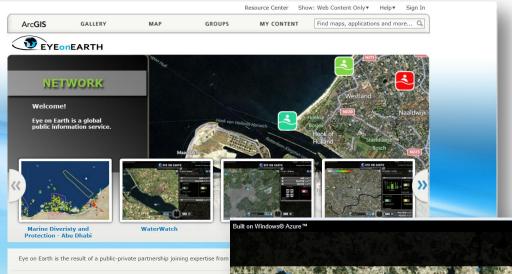
# New ways to analyze and communicate data

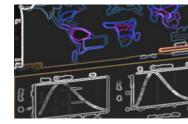# EYEonEARTH Network
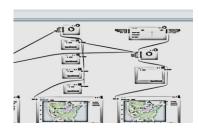
http://eyeonearth.org/
http://network.eyeonearth.org

# Computational Ecology and Environmental Science Group

- Developing new concepts, methods and tools to enable better information and predictions about our planet

- Joining up theory, with data, via statistics, to produce useful predictive models

- Targeting predictive models at key policy requirements to help develop better informed solutions

http://research.microsoft.com/ecology/



Global Carbon Model



Ecosystem function



Food security



Forest dynamics



Threats to biodiversity



Species distributions

# Global Carbon-Climate Feedback Model



Carbon in Atmosphere

Carbon in Ocean

Carbon on Land

Climate

e.g. temperature, winds

e.g. temperature, precipitation

# Current models are limited

# Multi-Component Model Framework

# Bringing Together Data for Models

# Complex event processing

- Event-driven computations
- Reason about time
- Detect interesting patterns
- Connect to and correlate heterogeneous sources
- Process late data
- Process lots of data
- Re-use existing functions and algorithms

# Microsoft StreamInsight

- API to build CEP applications
- Development in .NET & LINQ
- Maturing tool support
  - Event Flow Debugger
  - LINQPad
- Flexible integration of data sources / sinks
- Extensibility model
- Needs a SQL Server 2008 R2 License



StreamInsight Platform

# NodeXL

## Network graph visualization



Binary and source code:
http://nodexl.codeplex.com

# NodeXL - Network Overview Discovery and Exploration add-in for Excel 2007/2010

# World Wide Telescope

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

Goals
- Integration of data sets and one-click contextual access
- Easy access and use

Up...

Not...

We...

www.layerscape.org

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

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

**KINECT**™
for Windows®

- 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

# 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

# eScience in Action

# Microsoft eScience Workshop 2012

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

**Microsoft**®

Microsoft® Research Connections

# Presentation Fonts/Typography

## What are the font choices and sizes?

Any font size 28pt or larger should use Segoe UI Light
Any type that is less than 28pt should use Segoe UI
On one slide, try to use a maximum of 3 font sizes

## Where to start?

Control focus by using scale versus bullets
Use color to draw focus when necessary
Start with main topic at 40pt and subtopics at 20pt

## Title/sentence casing and periods

Title text should be "title caps" including a, is, of, & and
All supporting slide text should be sentence case
Periods should only be used on the title & main topics for complete sentences

# Slide Palette Info

The PowerPoint palette for this template has been built for you and is shown below. Avoid using too many colors in your presentation.

Color text. Select the 4th color from the left for subheads and 1st level non-bulleted text color, or wherever "color" text is preferred over the default black/white text

**Primary colors**

**Accent colors**

| Text Light 2 | Primary 1 | Primary 2 | Primary 3 | | Accent 4 | Accent 5 | Accent 6 |

Use **Primary 1** as the main color. Use **Primary 2** and **Primary 3** when additional colors are needed.

Use **Accents 4-6** sparingly – only when more colors are necessary.

# Preferred Text Layout (No Bullets)

## Main topic 1: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 2: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 3: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

# Same Color Text Layout (No Bullets)

## Main topic 1: size 40pt
Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 2: size 40pt
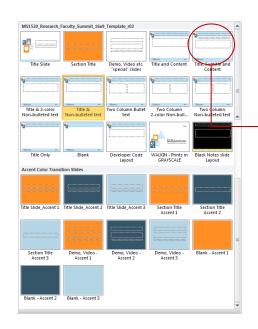Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 3: size 40pt
Size 20pt for the subtopics
Size 20pt for the subtopics

This is an almost identical layout to the previous slide, except that it has **all white text**. Sometimes you may prefer not to use colored text – for example if your list is only top level points, all white might look better. You can choose the layout you prefer.

Here's how to select different layouts:

1. Click on the Home tab at the top (if not already selected)
2. Click on Layout. A drop down list similar to the one shown on the left will appear. Notice that the layout for the slide you are on is **highlighted**. This slide uses a layout called "Title & Non-bulleted text"
3. Try clicking on the Layout to the left of it, called "Title & 2-color Non-bulleted text". Notice how the 1st level subheads change to a color.
4. Next try clicking on the layout called "Title and Content" layout. This is a the bulleted layout used on the next slide.
5. Use Layouts to set up new slides or to change existing slide layouts.

# Adjusting List Levels

## Main topic 1: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 2: size 40pt

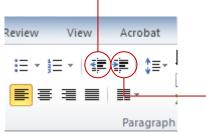Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 3: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

Use the "Increase List Level" and "Decrease List Level" tools on the Home Menu to change text levels.

Try this:
1. Place your cursor in any row of text to the left that says "Size 20pt for subtopics"
2. Next click the Home tab, and then on the "Decrease List level" tool. Notice how the line jumps up a level in size.
3. Now try placing your cursor in one of the "Main topic..." lines of text. Click the "Increase List Level" tool and see how the text is pushed down one level

Use these 2 tools to adjust your text levels as you work

# Preferred Two Column Layout

## Main topic 1: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 2: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 3: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 1: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 2: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 3: size 40pt

Size 20pt for the subtopics
Size 20pt for the subtopics

# Same Color Two Column Layout

## Main topic 1: size 40pt
Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 2: size 40pt
Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 3: size 40pt
Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 1: size 40pt
Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 2: size 40pt
Size 20pt for the subtopics
Size 20pt for the subtopics

## Main topic 3: size 40pt
Size 20pt for the subtopics
Size 20pt for the subtopics
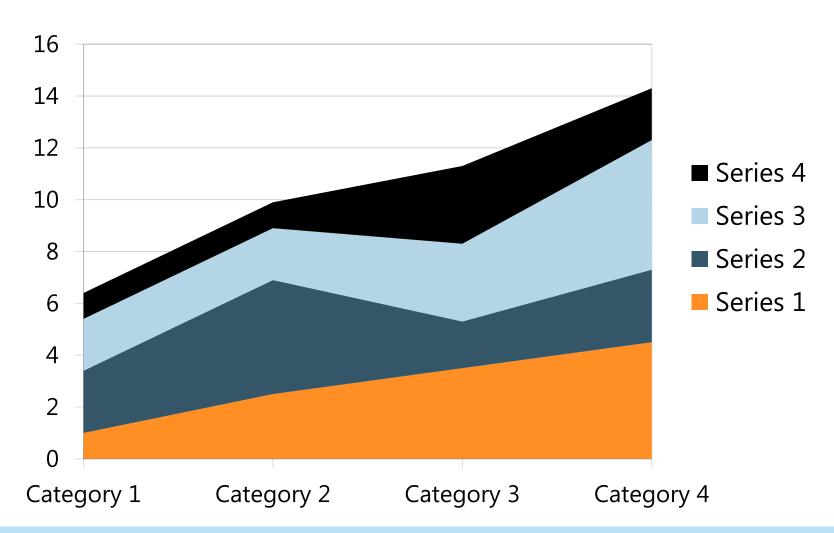
# Demo Title

# demo

Name
Title

# Video Title

video

# Chart Example

# Slide for Showing Software Code

Add code here