# A PROVENANCE-BASED INFRASTRUCTURE FOR CREATING REPRODUCIBLE PAPERS

Juliana Freire

*Juliana.freire@nyu.edu*
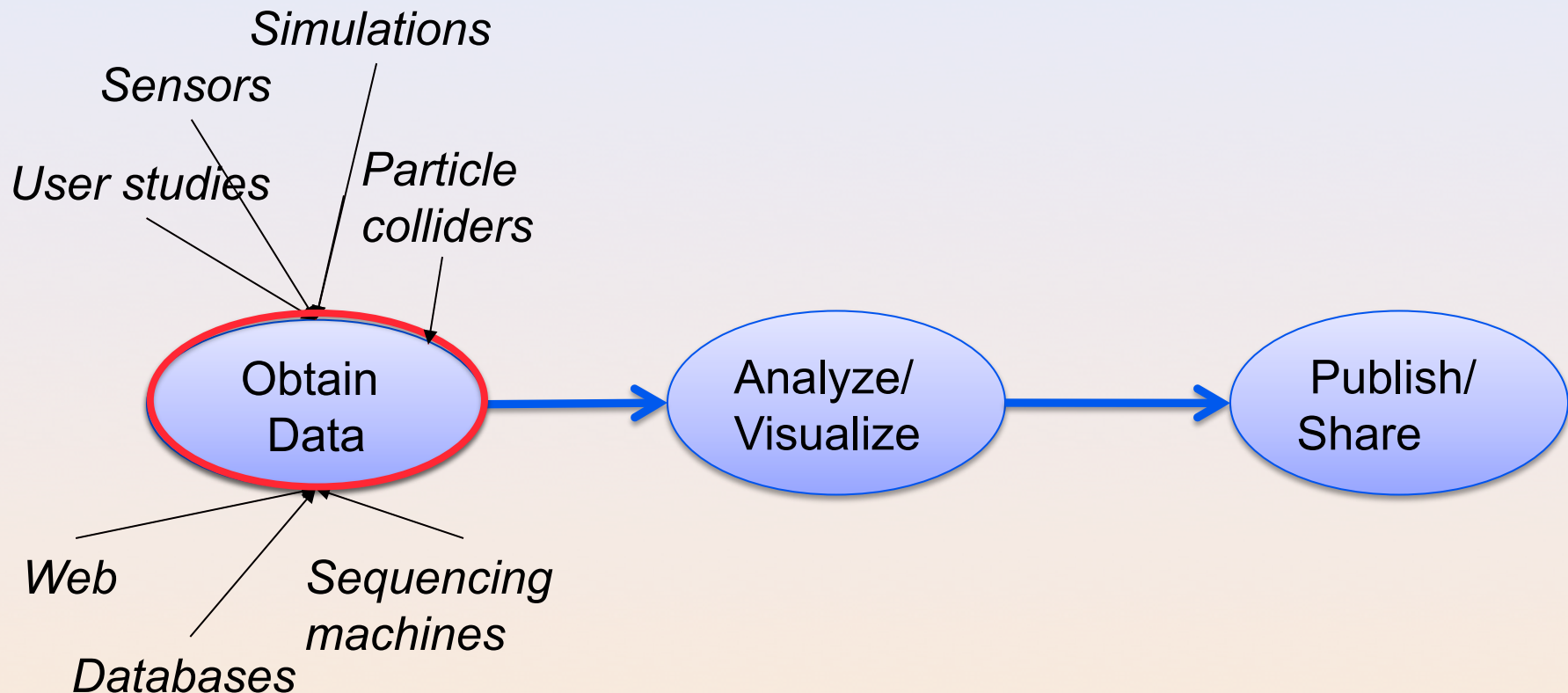
VisTrails Group & Web and Databases Lab
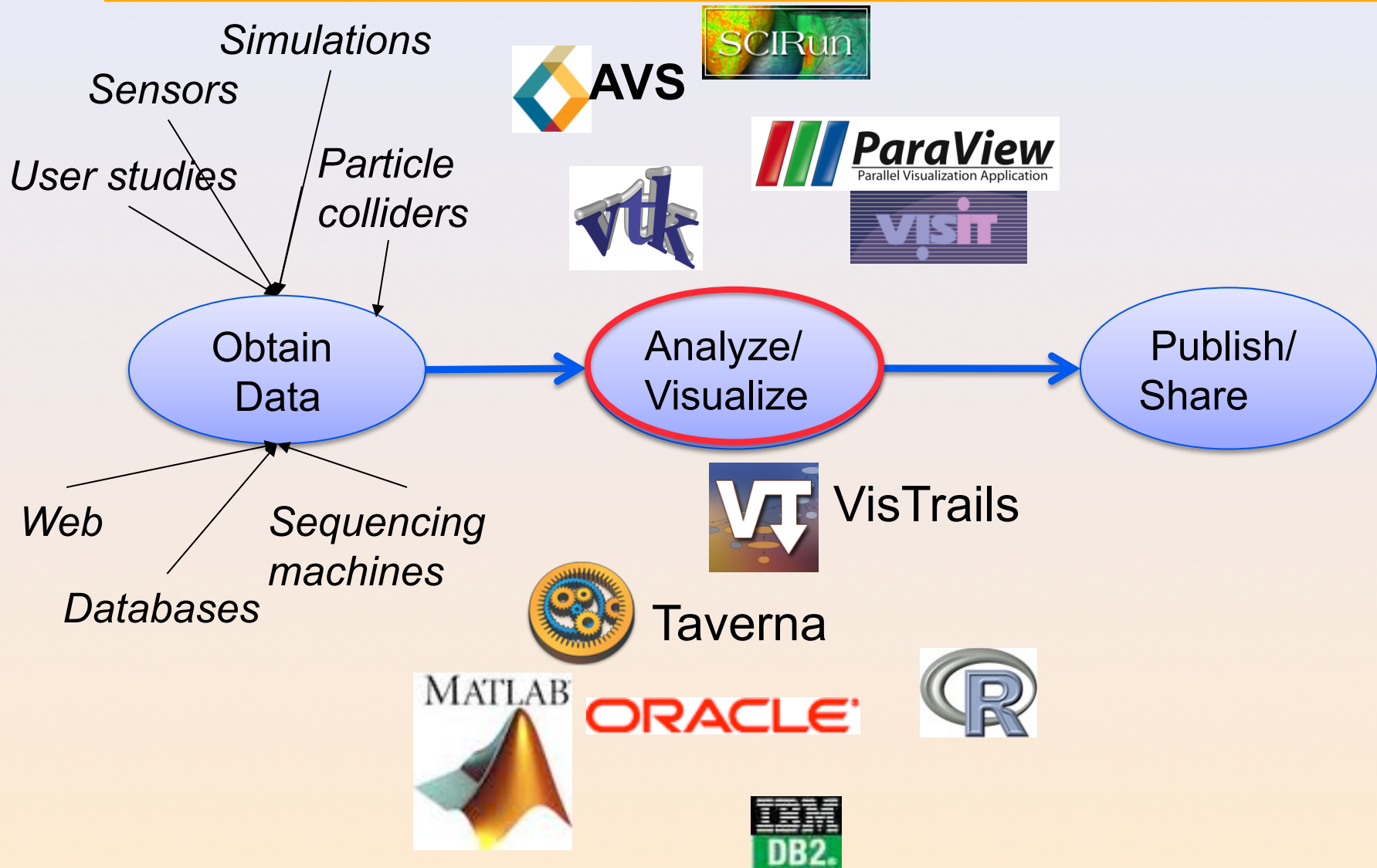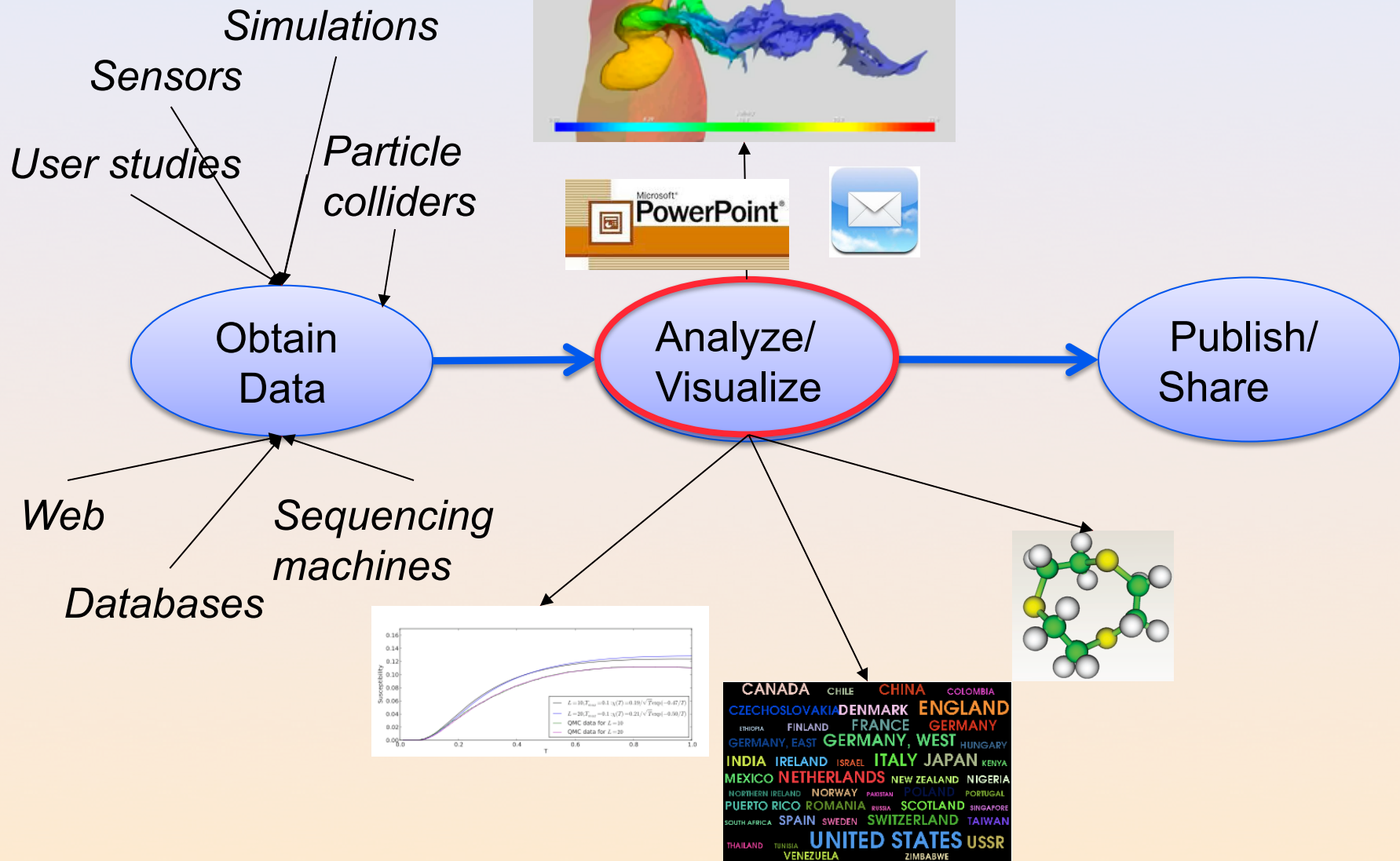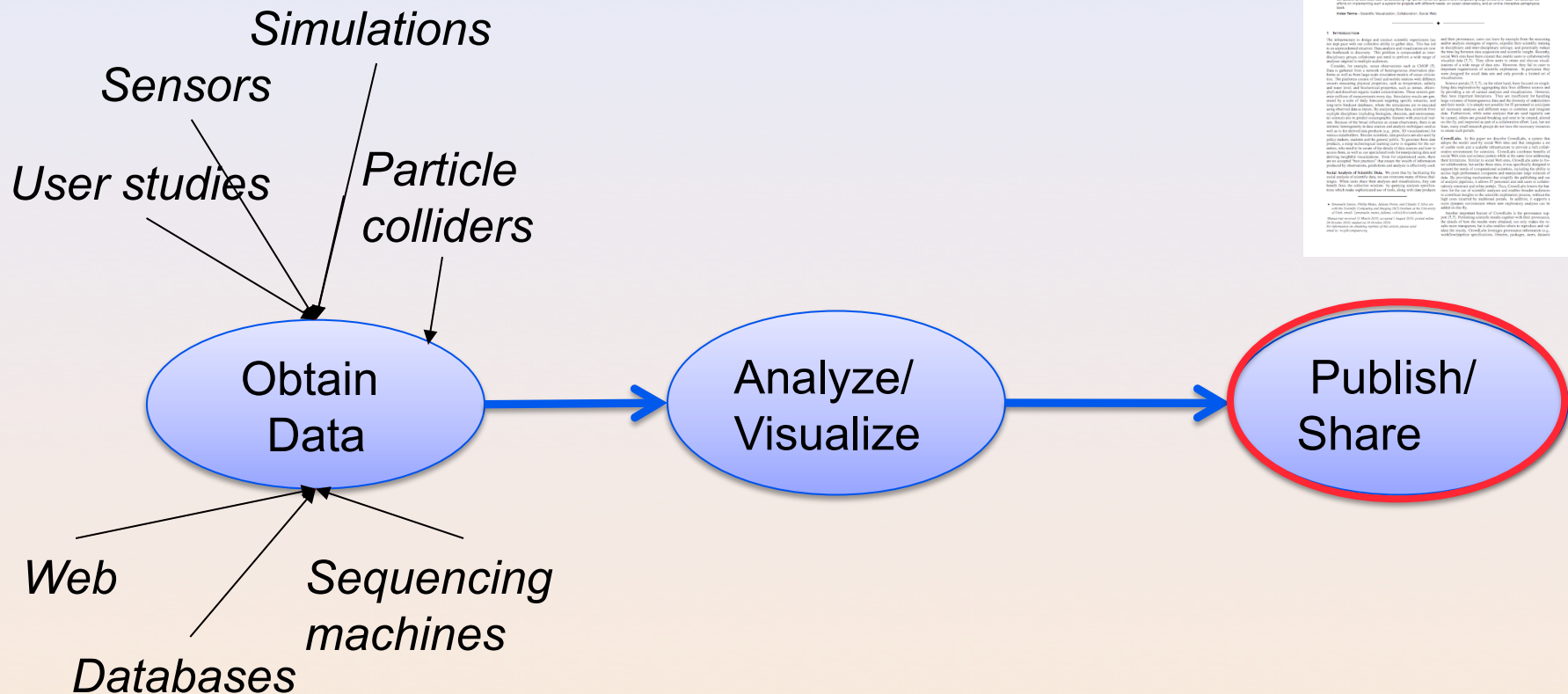
NYU Poly

# Science Today: Data Intensive



*Simulations*

*Sensors*

*User studies*

*Particle colliders*

**Obtain Data**

**Analyze/ Visualize**

**Publish/ Share**

*Web*

*Sequencing machines*

*Databases*

# Science Today: Data + Computing Intensive



*Simulations*

*Sensors*

*User studies*

*Particle colliders*

**AVS**

**SCIRun**

**vtk**

**ParaView** — Parallel Visualization Application

**VISIT**

Obtain Data → Analyze/Visualize → Publish/Share

*Web*

*Sequencing machines*

*Databases*

**VT** VisTrails

Taverna

**MATLAB**

**ORACLE**

**R**

**IBM DB2**

# Science Today: Data + Computing Intensive

*Simulations*

*Sensors*

*User studies*

*Particle colliders*

*Web*

*Sequencing machines*

*Databases*

**Obtain Data** → **Analyze/Visualize** → **Publish/Share**

# Science Today: Data + Computing Inte

*Simulations*

*Sensors*

*User studies*

*Particle colliders*

**Obtain Data** → **Analyze/ Visualize** → **Publish/ Share**

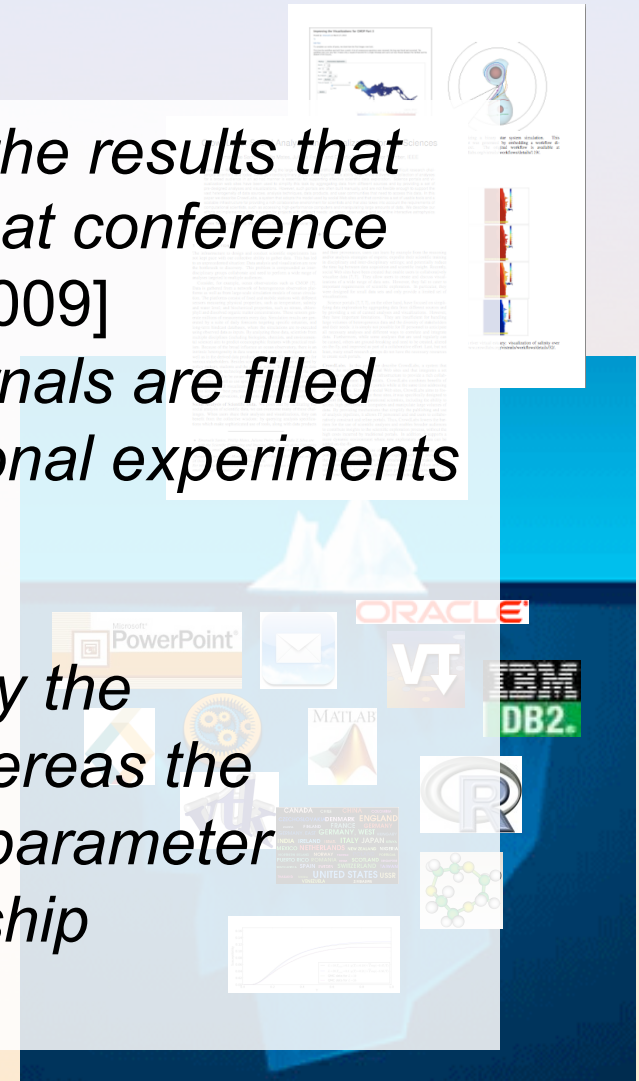*Web*

*Sequencing machines*

*Databases*

# Science Today: Incomplete Publications

- ◆ Publications are just the tip of the iceberg
  - – Scientific record is incomplete--- to large to fit in a paper
  - – Large volumes of data
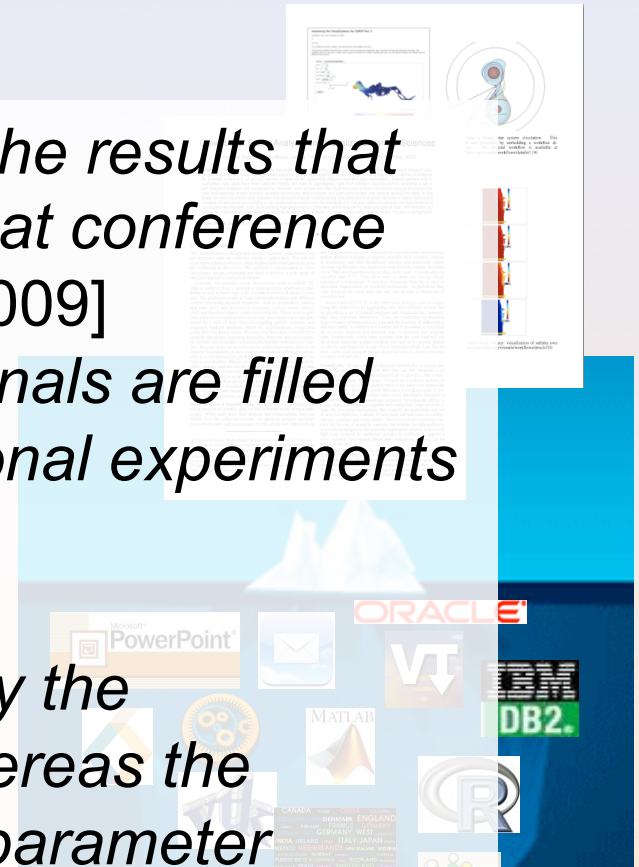  - – Complex processes
- ◆ Can't (easily) reproduce results

# Science Today: Incomplete Publications

◆ Publications are just the tip of the
iceberg

  – Scientific record is incomplete---
    to large to fit in a paper
  – Large volumes of data
  – Complex processes

◆ Can't (easily) reproduce results

*"It's impossible to verify most of the results that computational scientists present at conference and in papers."* [Donoho et al., 2009]

*"Scientific and mathematical journals are filled with pretty pictures of computational experiments that the reader has no hope of repeating."* [LeVeque, 2009]

*"Published documents are merely the advertisement of scholarship whereas the computer programs, input data, parameter values, etc. embody the scholarship itself."* [Schwab et al., 2007]

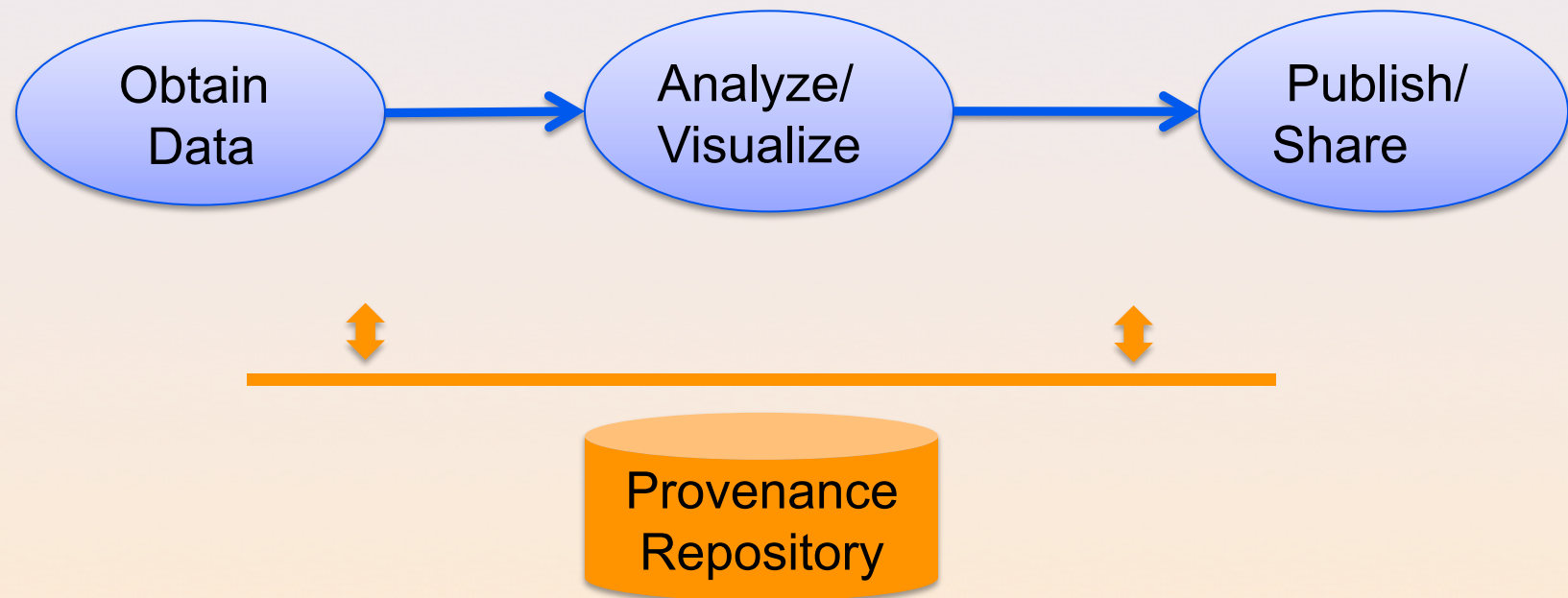# Science Today: Incomplete Publications

◆ Publications are just the tip of the iceberg

– Scientific record is incomplete--- to large to fit in a paper

– Large volumes of data

– Complex processes

◆ Can't (easily) reproduce results

*"It's impossible to verify most of the results that computational scientists present at conference and in papers."* [Donoho et al., 2009]

*"Scientific and mathematical journals are filled with pretty pictures of computational experiments that the reader has no hope of repeating."* [LeVeque, 2009]

*"Published documents are merely the advertisement of scholarship whereas the computer programs, input data, parameter values... itself."*

http://en.wikipedia.org/wiki/Scientific_misconduct

http://ori.dhhs.gov/misconduct/cases/

Nobel Laureate Retracts Two Papers, NYTimes 09/24/2010

# Vision: Provenance-Rich Science

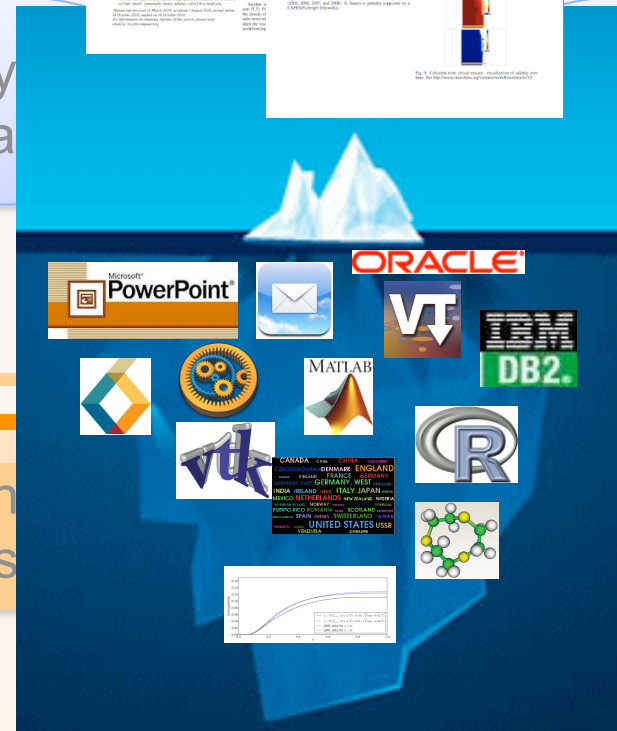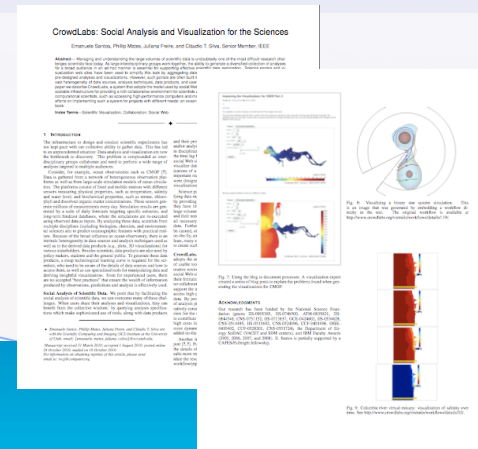# Vision: Provenance-Rich Science



*Provenance is the scientific record*

Provenance Repository

# Provenance-Rich Publications

◆ Bridge the gap between the scientific process and publications
  – Papers with *deep* captions and a *complete and trustworthy* scientific record

◆ Show me the proof: results that can be reproduced and validated

◆ Encouraged by ACM SIGMOD, a number of journals, funding agencies, academic institutions
  – E.g., ETH
    http://www.vpf.ethz.ch/services/researchethics/Broschure

◆ Several workshops, different communities
  – Beyond The PDF, SIAM Symposium on Reproducible Research, AMP Workshop on Reproducible Research, Workshop on Archiving Experiments

# Provenance-Rich Publications: Benefits

- Produce more knowledge---not just text
- Allow scientists to stand on the shoulders of giants
- Science can move faster
  - http://www.nytimes.com/2011/06/26/opinion/sunday/26ideas.html?_r=1
- Allow scientists to stand on their own shoulders!
- Higher-quality publications
  - Authors will be more careful
  - Many eyes to check results
- Describe more of the discovery process: people only describe successes, can we learn from mistakes?
- Expose scientific community to different techniques and tools: expedite their training; and potentially reduce their time to insight

# Provenance-Rich Publications: Challenges

- ◆ It is too hard, time-consuming for authors to prepare compendia of reproducible results
  - – Data, computations, parameter settings, environment, etc.
- ◆ It is too hard for reviewers (and readers) to install, compile, and reproduce experiments
  - – Different OSes, library versions, hardware, large data, incompatible data formats…
- ◆ Need to simplify the process of sharing, reviewing and re-using scientific experiments and results

# Our Approach: An Infrastructure to Support Provenance-Rich Papers [Koop et al., ICCS 2011]

◆ Tools for *authors* to create reproducible papers

– Specifications that encode the computational processes

– Package the results    *Support different approaches*

– Link from publications

◆ Tools for testers to repeat and validate results

– Explore different parameters, data sets, algorithms

◆ Interfaces for searching, comparing and analyzing experiments and results

– Can we discover better approaches to a given problem?

– Or discover relationships among workflows and the problems?

– How to describe experiments?

# An *Provenance-Rich* Paper: ALPS2.0



The ALPS project release 2.0:
Open source software for strongly correlated systems

B. Bauer[1] L. D. Carr[2] H.G. Evertz[3] A. Feiguin[4] J. Freire[5]
S. Fuchs[6] L. Gamper[1] J. Gukelberger[1] E. Gull[7] S. Guertler[8]
A. Hehn[1] R. Igarashi[9,10] S.V. Isakov[1] D. Koop[5] P.N. Ma[1]
P. Mates[1,5] H. Matsuo[11] O. Parcollet[12] G. Pawłowski[13]
J.D. Picon[14] L. Pollet[1,15] E. Santos[5] V.W. Scarola[16]
U. Schollwöck[17] C. Silva[5] B. Surer[1] S. Todo[10,11] S. Trebst[18]
M. Troyer[1]‡ M. L. Wall[2] P. Werner[1] S. Wessel[19,20]

[1]Theoretische Physik, ETH Zurich, 8093 Zurich, Switzerland
[2]Department of Physics, Colorado School of Mines, Golden, CO 80401, USA
[3]Institut für Theoretische Physik, Technische Universität Graz, A-8010 Graz, Austria
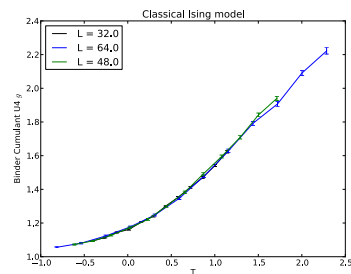[4]Department of Physics and Astronomy, University of Wyoming, Laramie, Wyoming 82071, USA
[5]Scientific Computing and Imaging Institute, University of Utah, Salt Lake City, Utah 84112, USA
[6]Institut für Theoretische Physik, Georg-August-Universität Göttingen, Göttingen, Germany
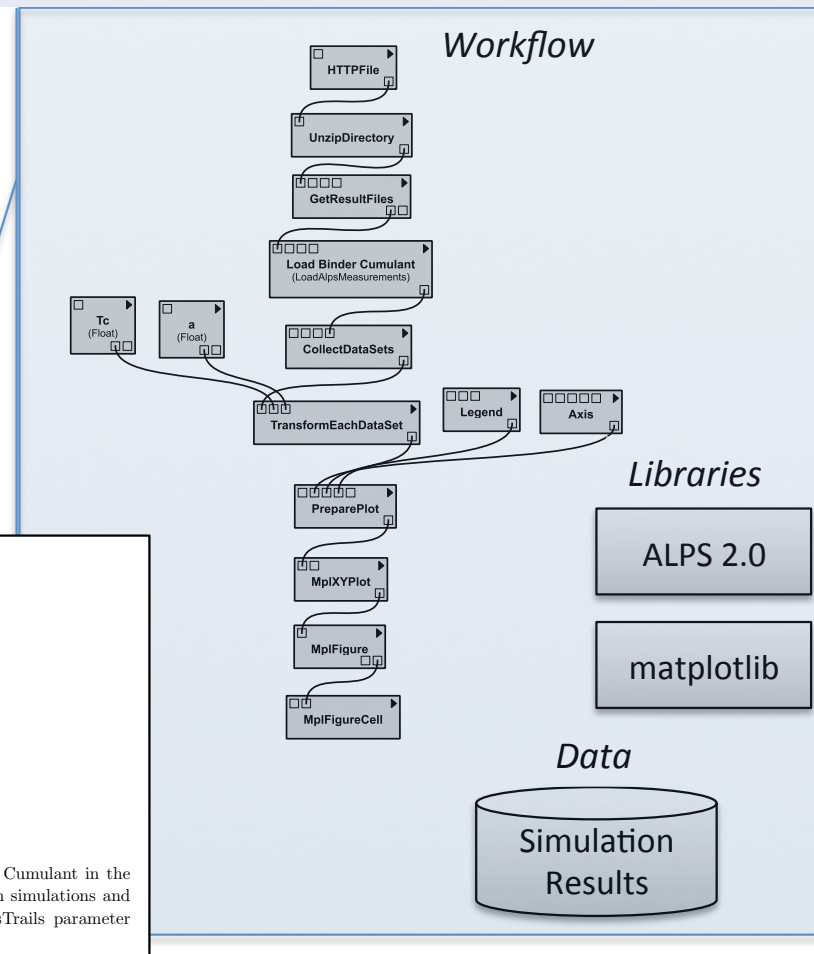[7]Columbia University, New York, NY 10027, USA
[8]Bethe Center for Theoretical Physics, Universität Bonn, Nussallee 12, 53115 Bonn, Germany

**Figure 3.** In this example we show a data collapse of the Binder Cumulant in the classical Ising model. The data has been produced by remotely run simulations and the critical exponent has been obtained with the help of the VisTrails parameter exploration functionality.

[Bauer et al., JSTAT 2011]

http://adsabs.harvard.edu/abs/2011arXiv1101.2646B

# A *Reproducible* Paper: ALPS2.0



[Bauer et al., JSTAT 2011]

http://adsabs.harvard.edu/abs/2011arXiv1101.2646B

# Some Videos

Editing an executable paper written using LaTeX and VisTrails
http://www.vistrails.org/download/download.php?type=MEDIA&id=executable_paper_latex.mov

Exploring a Web-hosted paper using server-based computation
http://www.vistrails.org/download/download.php?type=MEDIA&id=executable_paper_server.mov

An interactive paper on a Wiki*
http://www.vistrails.org/index.php/User:Tohline/CPM/Levels2and3

# Reproducible Papers

An interactive paper on a Wiki*
http://www.vistrails.org/index.php/User:Tohline/CPM/Levels2and3

The ALPS 2.0 paper
http://adsabs.harvard.edu/abs/2011arXiv1101.2646B

# Writing & Development

*An author benefits from working in an environment that simplifies the creation of an executable paper*

◆ First prototype: Leverage VisTrails' infrastructure

[Koop et al., ICCS 2011]

# The VisTrails System

- Workflow-based system for data analysis and visualization
- Comprehensive *provenance infrastructure*
- *Transparently* tracks provenance of the discovery process---from data acquisition to visualization
  - The *trail* followed as users generate and test hypotheses
- *Leverage provenance to streamline exploration*
  - Support for reflective reasoning and collaboration
  - Query and mine provenance
- Focus on usability—build tools for scientists
- The system is *open source*: http://www.vistrails.org
  - Multi-platform: Linux, Mac, Windows
  - Written in Python + Qt

# The VisTrails System

- ◆ Workflow-based system for data analysis and visualization

- ◆ Comprehensive *provenance infrastructure*

- ◆ *Transparently* tracks provenance of the discovery process---from data acquisition to visualization
  - – The *trail* followed as users generate and test hypotheses

- ◆ *Leverage provenance to streamline exploration*
  - – Support for reflective reasoning and collaboration
  - – Query and mine provenance

- •Visualizing environmental simulations (CMOP STC)
- •Simulation for solid, fluid and structural mechanics (Galileo Network, UFRJ Brazil)
- •Quantum physics simulations (ALPS, ETH Switzerland)
- •Climate analysis (CDAT)
- •Habitat modeling (USGS)
- •Open Wildland Fire Modeling (U. Colorado, NCAR)
- •High-energy physics (LEPP, Cornell)
- •Cosmology simulations (LANL)

- •Study on the use of tms for improving memory (Pyschiatry, U. Utah)
- •eBird (Cornell, NSF DataONE)
- •Astrophysical Systems (Tohline, LSU)
- •NIH NBCR (UCSD)
- •Pervasive Technology Labs (Heiland, Indiana University)
- •Linköping University (Sweden)
- •University of North Carolina, Chapel Hill
- •UTEP

# Writing & Development

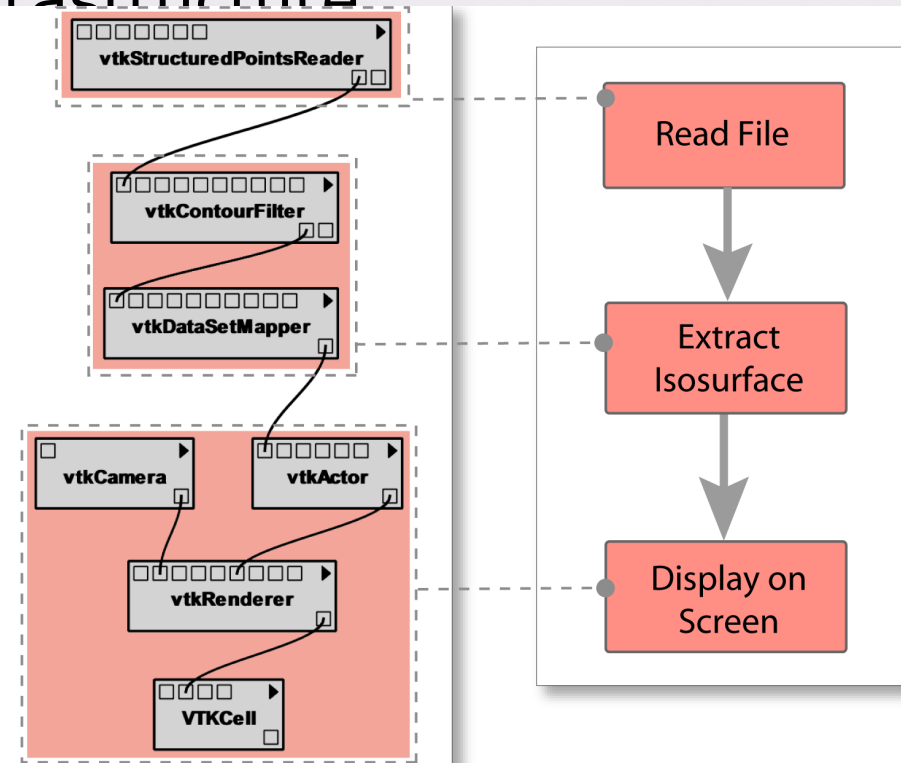*An author benefits from working in an environment that simplifies the creation of an executable paper*

- Leverage VisTrails' infrastructure
- Computations specified as workflows
  - Ability to combine tools
  - Support different levels of granularity facilitates the understanding of the computations and results

[Koop et al., ICCS 2011]

# Writing & Development

A ... working in an environment that s ... f an executable paper

... rastructure ...

```
1    import vtk
2
3    data = vtk.vtkStructuredPointsReader()
4    data.SetFileName("../examples/data/head.120.vtk")
5
6    contour = vtk.vtkContourFilter()
7    contour.SetInput(0,data.GetOutput())
8    contour.SetValue(0, 67)
9
10   mapper = vtk.vtkPolyDataMapper()
11   mapper.SetInput(contour.GetOutput())
12   mapper.ScalarVisibilityOff()
13
14   actor = vtk.vtkActor()
15   actor.SetMapper(mapper)
16
17   cam = vtk.vtkCamera()
18   cam.SetViewUp(0,0,-1)
19   cam.SetPosition(745,-453,369)
20   cam.SetFocalPoint(135,135,150)
21   cam.ComputeViewPlaneNormal()
22
23   ren = vtk.vtkRenderer()
24   ren.AddActor(actor)
25   ren.SetActiveCamera(cam)
26   ren.ResetCamera()
27
28   renwin = vtk.vtkRenderWindow()
29   renwin.AddRenderer(ren)
30
31   style = vtk.vtkInteractorStyleTrackballCamera()
32   iren = vtk.vtkRenderWindowInteractor()
33   iren.SetRenderWindow(renwin)
34   iren.SetInteractorStyle(style)
35   iren.Initialize()
36   iren.Start()
```

[Koop et al., ICCS 2011]

# Writing & Development

*An author benefits from working in an environment that simplifies the writing of an executable paper*

◆ Provenance of data and computations: workflow provenance is not sufficient

[Koop et al., ICCS 2011]

# Sharing an Experiment

◆ Juliana creates an experiment

/Users/juliana/head.vtk

Read File

VTK 1.2

Extract
Isosurface

Display on
Screen

◆ Ian tries to run Juliana's experiment

/Users/juliana/head.vtk

Read File

**File not found!**

VTK 1.2

Extract
Isosurface

**Cannot execute**

Display on
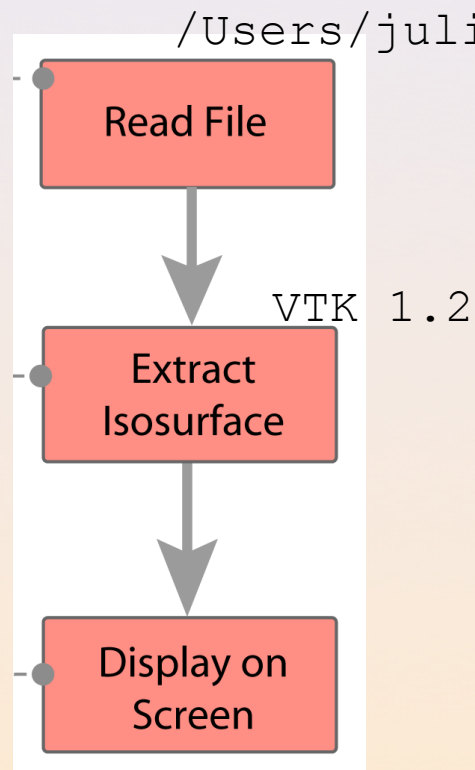Screen

# Writing & Development

*An author benefits from working in an environment that simplifies the writing of an executable paper*

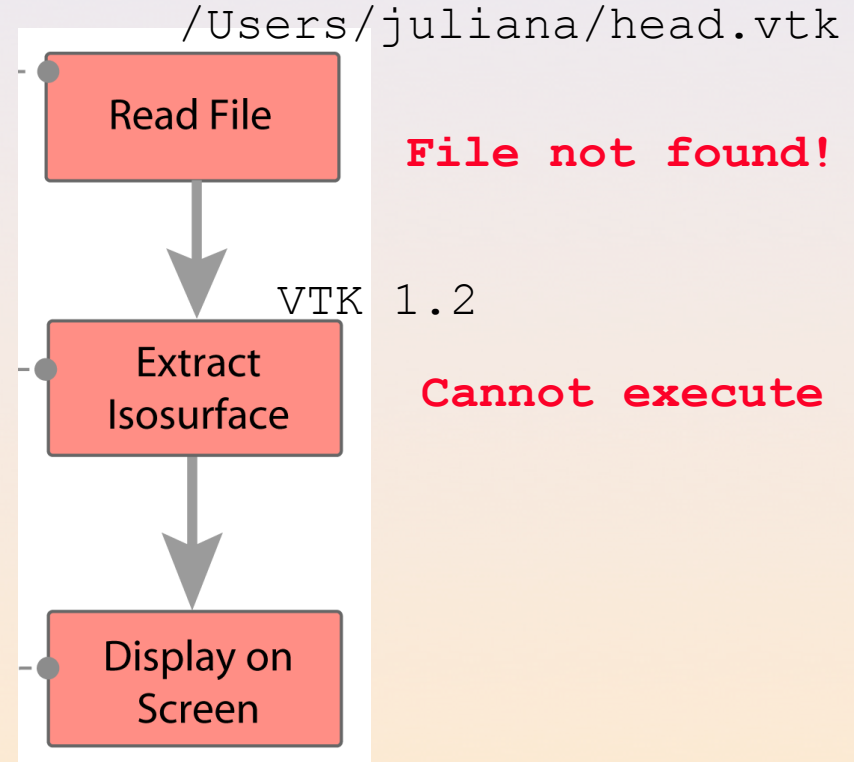- ◆ Provenance of data and computations: workflow is not sufficient
- ◆ Need 'more' information: computational environment (OS, library versions, etc.)
  - – Also use virtual machines, CDEPack
- ◆ Need better file management
  - – Designed support for strong links between data and their provenance [Koop@SSDBM2010]
  - – Use versioning servers (e.g., GIT, SVN, Oracle DBFS)
- ◆ Connect results to their provenance
  - – Support LateX, Word, Powerpoint, HTML, wiki

[Koop et al., ICCS 2011]

# Review & Validation

*Improve the quality of reviews: reviewers have the ability to explore and validate conclusions*

- ◆ Execution environment
  - – Use provenance, virtual machines, CDEPack to deal with software dependencies
  - – Support local, remote, and mixed execution: alternatives to handle proprietary code and data, special hardware
- ◆ Testing and validating computations and their results
  - – Reproduce
  - – Workability: explore parameters and configurations the authors might not have described in the paper
  - – VisTrails' data exploration infrastructure comes in handy here
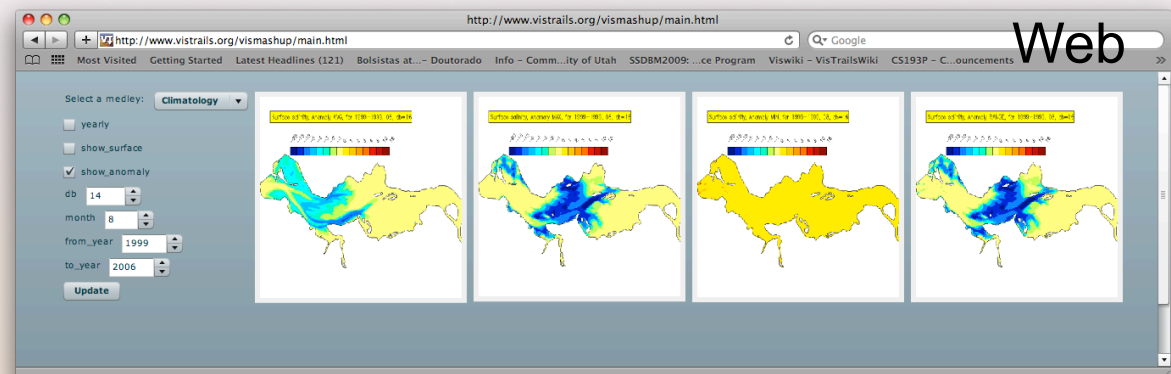
[Koop et al., ICCS 2011]

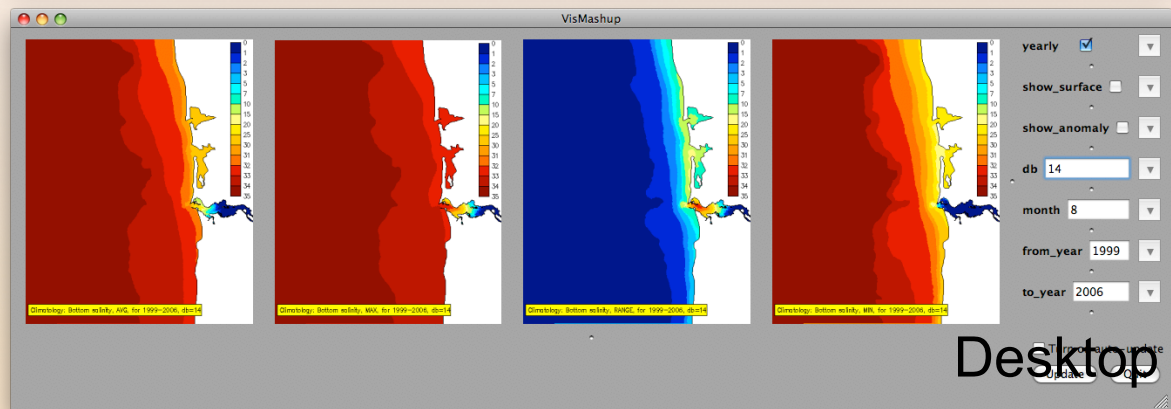# Publishing, Maintenance, & Re-Use

◆ Simplify interaction: the VisMashup system
[Santos@TVCG2009]

◆ Publish using different media, not just *documents*



Web

Portable
Devices

Desktop

# Publishing, Maintenance, & Re-Use

◆ Simplify interaction: the VisMashup system
  [Santos@TVCG2009]

◆ Publish using different media, not just *documents*

◆ Maintenance and longevity

  – Software evolves: need to *upgrade* experiments
    [Koop@IPAW2010]

◆ Querying and re-using published experiments [Freire
  et al., VLDB 2011]

  – Opportunities for knowledge discovery and re-use

  – A search/query engine for experiments: text + structure
    [Scheidegger@TVCG2007]

  – Can we discover better approaches to a given problem? Or
    discover relationships among workflows and problems?

  – Can we combine multiple results?

# Current Uses and Experiences

◆ ALPS community: ETH group has published a number of reproducible papers!

◆ Simulations of computational fluid dynamics

◆ Database research:
  – experiments using distributed database systems, querying Wikipedia
  – http://www.vistrails.org/index.php/RepeatabilityCentral
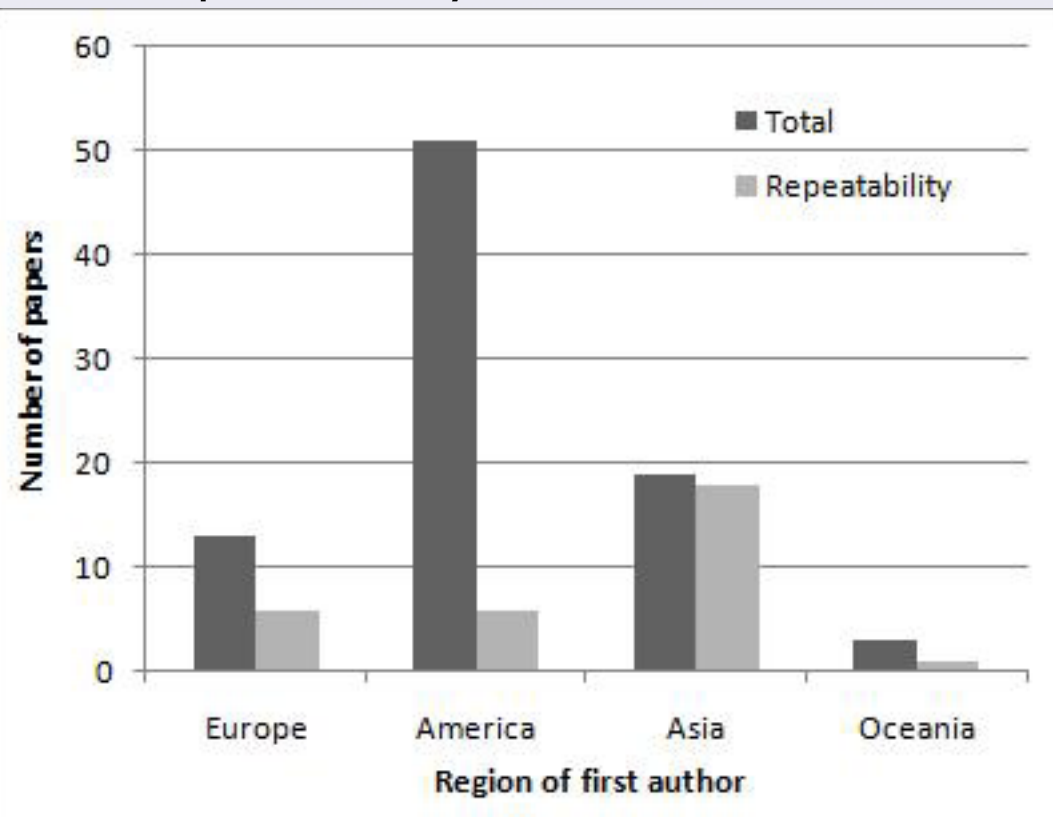
# Current Uses and Experiences

- ◆ ACM SIGMOD repeatability effort [Bonnet et al., SIGMOD Record 2011 to appear]
  - Since 2008 verifies the experiments published in accepted papers
  - Papers submitted for reproducibility evaluation: 2010-20 submissions; 2011-31 submissions
  - In 2011, lay out a set of guidelines to simplify and expedite the reviewing process; provided tutorials
  - Review was still challenging
    - » Common problem: setup failed due to implicit dependencies
    - » Easy to solve with a virtual machine…
  - Reasons for not submitting:
    - » Intellectual property rights on software
    - » Sensitive data
    - » Specific hardware requirements

http://www.sigmod2011.org/calls_papers_sigmod_ research_repeatability.shtml

# Current Uses and Experiences

◆ ACM SIGMOD repeatability effort [Bonnet et al., SIGMOD Record 2011 to appear]

- Since _____ in accepted paper

- Paper _____: 2010-20 subm___

- In 20_____ and expedite the re___

- Revie_____

  » Co_____ependencies

  » Ea_____

- Reaso_____

  » Int_____

  » Sensitive data

  » Specific hardware requirements



http://www.sigmod2011.org/calls_papers_sigmod__research_repeatability.shtml

# Going Forward

- ◆ **Need more and better incentives:**
  - – seal of quality, higher quality software/experiments, easier for newcomers in a project, citations, **recognition**

- ◆ **Need a whip(?): Some disciplines require data for publications, should we require computational experiments too?**          **ETH does!**

- ◆ **Need better tools**
  - – There is no one-size-fits-all solution
  - – Many groups building tools---we should join forces and build a *Reproducibility Toolkit*

- ◆ **Need standardS and guidelines for authors and tool developers**

- ◆ **Need provenance support in applications**
  - – Integrate provenance from different sources, connect the results

# Provenance Everywhere



Bird_table version 123476 in Oracle
- STEM v 1.2 + predictions
- R Script+results
- BirdVis vis spec

eBird

STEM

*predictions*          *results*

*predictions*

*Bird_table version 123476*

- *matplotlib script+results*

*Provenance-rich presentation*

# A Little History and a Challenge

◆ A long time ago, when I was a PhD student, generating the reference list for papers was **very** time consuming

– Find proceedings on the shelf (or walk to library), flip pages to obtain page numbers, type (title, authors, proceedings name, etc.)

◆ Today

– Google/Bing author or part of paper title, DBLP, ACM DL, IEEE Explore

– Copy bib entry in one of many formats (bibtex, EndNote, plain text), paste in paper, v*oilà!*

◆ *Can we do the same for scientific experiments?*

# Conclusions and Future Work

◆ Provenance is crucial for science and an enabler for *executable* papers

◆ Provenance must be at the center of the scientific process!

◆ Built an end-to-end solution based on VisTrails--- currently working on integrating infrastructure with other systems

– Provenance-enabling other tools

◆ Many challenges and several open research questions
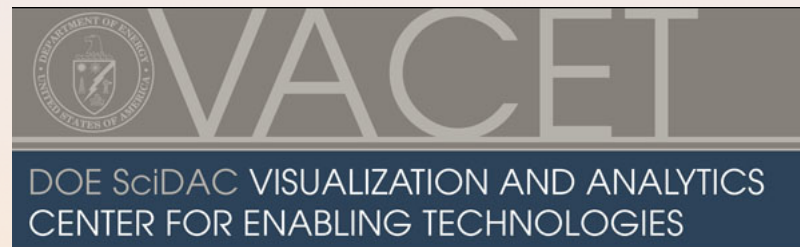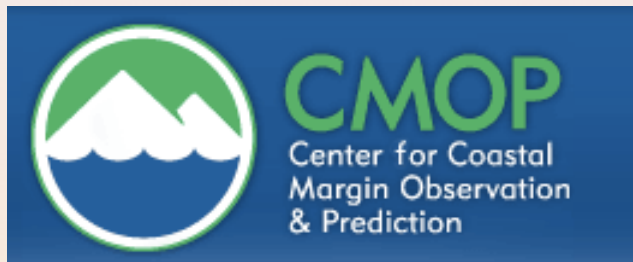
◆ Great opportunity to have impact in science

# Additional Information

◆ The VisTrails System http://www.vistrails.org

◆ An infrastructure to support the creation, review and re-use of reproducible papers
http://www.vistrails.org/index.php/ExecutablePapers

# Acknowledgments

◆ Thanks to: Philippe Bonnet, Philip Mates, Matthias Troyer, Dennis Shasha, Emanuele Santos, Claudio Silva, Joel Tohline, Huy T. Vo, and the VisTrails team

◆ This work is partially supported by the National Science Foundation, the Department of Energy, and IBM Faculty Awards.

Merci
*Ευχαριστω*
Thank you
Obrigada