## Advanced Data Management (CSCI 490/680)

Reproducibility

Dr. David Koop





## Provenance in Computational Science



D. Koop, CSCI 680/490, Spring 2021





2

## Database Provenance

- Motivation: Data warehouses and curated databases
  - Lots of work
  - Provenance helps check correctness
  - Adds value to data by how it was obtained
- Three Types:
  - Why (Lineage): Associate each tuple t present in the output of a query with a set of tuples present in the input
  - How: Not just existence but routes from tuples to output (multiple contrib.'s) - Where: Location where data is copied from (may have choice of different
  - tables)













# Why Provenance

### Agencies

	0		
	name	based_in	phone
$t_1$ :	BayTours	San Francisco	415-1200
$t_2$ :	HarborCruz	Santa Cruz	831-3000

### ExternalTours

name	destination	type	price
BayTours	San Francisco	cable car	\$50
BayTours	Santa Cruz	bus	\$100
BayTours	Santa Cruz	boat	\$250
BayTours	Monterey	boat	\$400
HarborCruz	Monterey	boat	\$200
HarborCruz	Carmel	train	\$90
	name BayTours BayTours BayTours BayTours HarborCruz HarborCruz	namedestinationBayToursSan FranciscoBayToursSanta CruzBayToursSanta CruzBayToursMontereyHarborCruzMontereyHarborCruzCarmel	namedestinationtypeBayToursSan Franciscocable carBayToursSanta CruzbusBayToursSanta CruzboatBayToursMontereyboatHarborCruzMontereyboatHarborCruzCarmeltrain

Q1:

SELECT a.name, a.phone

FROM Agencies a, ExternalTours e WHERE a.name = e.name AND e.type='boat'

### **Result of** $Q_1$ :

namo	nhone
IIaIIIC	phone
BayTours	415-1200
HarborCruz	831-3000

- Lineage of (HarborCruz, 831-3000): {Agencies(t2), ExternalTours(t7)}
- Lineage of (BayTours, 415-1200): {Agencies(t1), ExternalTours(t5,t6)}
- This is not really precise because we don't need both t5 and t6—only one is ok













## How Provenance

### Agencies

	name	based_in	phone
$t_1$ :	BayTours	San Francisco	415-1200
$t_2$ :	HarborCruz	Santa Cruz	831-3000

### ExternalTours

	name	destination	type	price
$t_3$ :	BayTours	San Francisco	cable car	\$50
$t_4$ :	BayTours	Santa Cruz	bus	\$100
$t_5$ :	BayTours	Santa Cruz	boat	\$250
$t_6$ :	BayTours	Monterey	boat	\$400
$t_7$ :	HarborCruz	Monterey	boat	\$200
$t_8$ :	HarborCruz	Carmel	train	\$90

### $Q_2$ :

SELECT	e.destination, $a.$ phone	<b>Result of</b> $Q_2$ :		
FROM	Agencies $a$ ,	destination	phone	
	(SELECT name,	San Francisco	415-1200	$t_1 \cdot (t_1 + t_3)$
	based_in AS destination	Santa Cruz	831-3000	$t_{2}^{2}$
	FROM Agencies $a$	Santa Cruz	415-1200	$t_1 \cdot (t_4 + t_5)$
	UNION	Monterey	415-1200	$t_1 \cdot t_6$
	SELECT name, destination	Monterey	831-3000	$t_1 \cdot t_7$
	FROM External Tours ) $e$	Carmel	831-3000	$t_1 \cdot t_8$
WHERE	a.name = e.name			

- How provenance gives more detail about how the tuples provide witnesses to the result
- Prov of (San Francisco, 415-1200):  $\{ \{ t1 \}, \{ t1, t3 \} \}$
- t1 contributes **twice**
- Uses provenance semirings (the "polynomial" shown on the right)
- $t_5)$











## Where Provenance

### Agencies

	name	based_in	phone
$t_1$ :	BayTours	San Francisco	415-1200
$t_2$ :	HarborCruz	Santa Cruz	831-3000

### ExternalTours

	name	destination	type	price
$t_3$ :	BayTours	San Francisco	cable car	\$50
$t_4$ :	BayTours	Santa Cruz	bus	\$100
$t_5$ :	BayTours	Santa Cruz	boat	\$250
$t_6$ :	BayTours	Monterey	boat	\$400
$t_7$ :	HarborCruz	Monterey	boat	\$200
$t_8$ :	HarborCruz	Carmel	train	\$90

 $Q_1$ : SELECT FROM WHERE

a.name, a.phone Agencies a, ExternalTours ea.name = e.nameAND *e*.type='boat'

 $Q'_1$ : SELECT FROM

WHERE

e.name, a.phone Agencies a, ExternalTours ea.name = e.nameAND *e*.type='boat'

### **Result of** $Q_1$ :

name	phone
BayTours	415-1200
HarborCruz	831-3000

- Where provenance traces to specific locations, not the tuple values
- Q and Q' give the same result but the name comes from different places
- Prov of HarborCruz in second output: (t2, name)
- Important in annotation-propogation

















## VisTrails

- Comprehensive provenance infrastructure for computational tasks
- Focus on exploratory tasks such as simulation, visualization, and data analysis
- Transparently tracks provenance of the discovery process from data acquisition to visualization
  - The trail followed as users generate and test hypotheses
  - Users can refer back to any point along this trail at any time
- Leverage provenance to streamline exploration
- Focus on usability—build tools for scientists





## Version Trees for Evolution Provenance

- Undo/redo stacks are linear!
- We lose history of exploration
- Old Solution: User saves files/state
- VisTrails Solution:
  - Automatically & transparently capture entire history as a tree
  - Users can tag or annotate each version
  - Users can go back to **any** version by selecting it in the tree











# Capturing Exploration: Version Tree of Workflows











## Capturing Exploration: Version Tree of Workflows











# Capturing Exploration: Version Tree of Workflows











































# Querying Provenance by Example

- Provenance is represented as graphs: hard to specify queries using text! • Querying workflows by example [Scheidegger et al., TVCG 2007; Beeri et al., VLDB 2006; Beeri et al. VLDB 2007]
- - WYSIWYQ -- What You See Is What You Query
  - Interface to create workflow is same as to query









## <u>Assignment 5</u>

- Four parts
  - Loading Data
  - Spatial Analysis
  - Graph Analysis
  - Temporal Analysis
- Due at the end of the semester (April 22, 2021)
- Questions?





## Final Exam

- Monday, April 26, 4:00-5:50pm, Online (Blackboard)
- Similar format
- More comprehensive (questions from topics covered in Test 1 & 2)
- machine learning
- Page to be posted

# • Will also have questions from temporal data, provenance, reproducibility,





## Stronger Links Between Provenance and Data



- Filenames are often the mode of identification in data exploration
- We might also use URIs or access curated data stores
  - Always expected for exploratory tasks?
  - What happens if offline?
- Solution:
  - Managed store for data associated with computations
  - Improved data identification
  - Automatic versioning





## Provenance from Data







## **Building Visualization Pipelines**











# **Building Visualization Pipelines**







## Completions

🙆 on wikingdin org /

men.wikipeula.org/		
http://en.wikipedia.org/		
http://encarta.msn.com/		
http://www.engadget.com/		
http://www.engadget.com/2008/09/09/live-from-apples-lets-rock-e		
http://en.wikipedia.org/wiki/VisTrails		
http://en.wikipedia.org/wiki/ACM_Transactions_on_Graphics		
http://en.wikipedia.org/wiki/Barack_Obama		
http://en.wikipedia.org/wiki/Columbus,_Ohio		
http://en.wikipedia.org/wiki/Joe_Biden		
http://en.wikipedia.org/wiki/John_McCain		



[Code Completion, Intellisense]

### D. Koop, CSCI 680/490, Spring 2021



## [URL Completion, Safari]

visualization	
visualizations for windows media player	1,670,000 results
visualization techniques	954,000 results
visualization tools	2,090,000 results
visualization board	3,380,000 results
visualization api	2,210,000 results
visualization toolkit	368,000 results
visualization technique	756,000 results
visualizations photography	1,830,000 results
visualization meditation	190,000 results
visualizations for media player	1,050,000 results
	<u>close</u>

[Web Search Completion, Google]





## Visualization Pipeline Completions



### D. Koop, CSCI 680/490, Spring 2021









19

## VisComplete Overview

- Mine provenance collection: Identify graph fragments that co-occur in a collection of workflows (Data-Driven)
- Predict sets of likely workflow additions to a given partial workflow











## Suggestion Interface



### D. Koop, CSCI 680/490, Spring 2021





21

## Suggestion Interface



### D. Koop, CSCI 680/490, Spring 2021





21

## VisComplete Results











## VisComplete Results











































# Generating Visualizations by Analogy





is to

is to

as

D. Koop, CSCI 680/490, Spring 2021







Northern Illinois University







# Generating Visualizations by Analogy











# Generating Visualizations by Analogy

• Compute difference  $\Delta(A,B)$  from provenance -  $D = \Delta(A,B) \circ C$  is often not a valid workflow










# Generating Visualizations by Analogy

- Compute difference  $\Delta(A,B)$  from provenance -  $D = \Delta(A,B) \circ C$  is often not a valid workflow
- Find map between A & C: map(A,C)











# Generating Visualizations by Analogy

- Compute difference  $\Delta(A,B)$  from provenance -  $D = \Delta(A,B) \circ C$  is often not a valid workflow
- Find map between A & C: map(A,C)
- Compute mapped difference  $\Delta AC(A,B) = map(A,C) \Delta(A,B)$ 
  - $D = \Delta AC(A,B) \circ C$











## VisMashup











# VisTrails for Teaching Scientific Visualization

- "Using VisTrails and Provenance for Teaching Scientific Visualization"
   [Silva et al., Eurographics Educator Program, 2010]
- Same features that scientists use for exploratory tasks can also benefit students
  Exploration: see all pipelines not just a
  - Exploration: see all pipelines not ju "final" one
  - Comparison: see different pipelines and what changes exist
  - Assessment: see how a solution was developed

D. Koop, CSCI 680/490, Spring 2021

Sheet 1 PE#0 critical\_points.vt 0









## Provenance Analysis of Projects









# Provenance Analysis of Projects



D. Koop, CSCI 680/490, Spring 2021

### Comparing Paths to Solutions for Two Students







# <u>The State of Repeatability in</u> <u>Computer Systems Research</u>

C. Collberg and T. Proebsting CACM 2016





# State of Repeatability in Computer Systems

- "Cool paper! Can you send me the system?"
- How hard is it to just re-execute published experiments
- Most people say they will share their code and data are available...
- Weak repeatability: Do authors make the source code used to create the results in their article available, and will it build?







## Experiment











## Repeatability Results



Figure 11: Study result. Blue numbers represent papers that were excluded from consideration, green numbers papers that are weakly repeatable, red numbers papers that are non-weakly repeatable, and orange numbers represent papers that were excluded (due to our restriction of sending at most one email to each author).

### D. Koop, CSCI 680/490, Spring 2021

OK <sup>≤30</sup> OK <sup>≤</sup> 64	>30 OK <sup>Auth</sup> 23	
	Notation	Number of papers
Build fails 9	HW	excluded due to replication requiring special hardware
	NC	excluded due to results not being backed by code
	EX	excluded due to overlapping author lists
	BC	where the results are backed by code
	Article	where code was found in the paper itself
	Web	where code was found through a Web search
	EM yes	where the author provides code after receiving an email message
	EM <sup>no</sup>	where the author responds to an email message saying code cannot be provided
	EMø	where the author does not respond to email requests within two months
	OK <sup>≤30</sup>	where code is available and we succeed in building the system in $\leq$ 30 minutes
	OK >30	where code is available and we succeed in building the system in >30 minutes
	OK <sup>Auth</sup>	where code is available and we fail to build, and the author says the code builds with reasonable effort
$\mathbf{M}^{\emptyset}$ $\mathbf{E}\mathbf{M}^{\mathrm{no}}$ 30 $146$	Fails	where code is available and we fail to build, and the author says the code may have problems building







33



### Excuses

- "Unfortunately the current system is not mature" "The code was never intended to be released so it is not in any shape for
- general use"
- "[Our] prototype included many moving pieces that only [student] knew how to operate... he left"
- "... the server in which my implementation was stored had a disk crash ... three disks crashed... Sorry for that"









### Excuses

- to speed than on our own research"
- "... we can't share what [we] did for this paper. ... this is not in the academic tradition, but this is a hazard in an industrial lab"
- "... based on earlier (bad) experience, we [want] to make sure that our implementation is not used in situations that it is not meant for"

• "...when we attempted to share it, we [spent] more time getting outsiders up













## Excuse Classification

- Versioning
- Available Soon
- No Intention to Share
- Personnel Issues
- Lost Code
- Academic Tradeoffs
- Industrial Lab Tradeoffs
- Obsolete HW/SW
- Controlled Usage
- Privacy/Security
- Design Issues

D. Koop, CSCI 680/490, Spring 2021





Northern Illinois University









# Some of these are (partially) people problems, not technical problems







# Examining 'Reproducibility in Computer Science'

- Repeat the experiment in reproducibility!
- Differences from original
- Shows issues with trying to classify experiments

F	າ	ır
	D	)i

All Others Purported Not 27%

- ported Not Building; 6% ••••• sputed; Not Checked
- Purported Building; Disputed; 2% •• Not Checked
  - Conflicting Checks! 0%
    - Misclassified 1% •
  - Purported Not Building But 14% ••••••••• Found Building
- Purported Building But Found 0% Not Building
  - Purported Not Building; 0% Confirmed
- Purported Building; Confirmed 0% •







## Recommendations

- Fund repeatability engineering
- Require sharing contracts

Location	<ul> <li>email address and/or web site</li> </ul>
Resource	<ul> <li>types: code, data, media, document</li> <li>availability: no access, access, NDA</li> <li>expense: free, non-free, free for aca</li> <li>distribution form: source, binary, se</li> <li>expiration date</li> <li>license</li> <li>comment</li> </ul>
Support	<ul> <li>kinds: resolve installation issues, fix upgrade to new language and operations system versions, port to new environ improve performance, add features</li> <li>expense: free, non-free, free for aca</li> <li>expiration date</li> </ul>

### D. Koop, CSCI 680/490, Spring 2021

ation access demics ervice

bugs, ting ments,

idemics





Northern Illinois University







## Reproducible Research

- Science is verified by replicating work independently
- Replication Issues:
  - Requires many resources to replicate (Sloan Digital Sky Survey) - Requires significant computing power (Climate Model Simulation) - Requires too much time or very specific circumstances (Environment

  - Epidemiology)
- Reproducibility
  - Replication of the analysis based on the collected data (not replicating the data collection itself)
  - Better if we have the actual code or available executables \_







# Reproducibility Spectrum



### D. Koop, CSCI 680/490, Spring 2021









41

# Published Papers

- "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]







# Problem: Incomplete Publications

- A paper cannot include all relevant details of the science
  - Large volumes of data
  - Complex processes
  - Code dependencies
- This makes publishing complete results more difficult!







### VISUALIZATION CORNER



Figure 2. The VisMashup window that displays when users select the "Figure 2" tab (see www.vistrails.org/index.php/User:Tohline/IVAJ/Levels2and3). The window displays an image generated by a customized VisTrails workflow using the indicated values of the three variable parameters, Omega\_frame (=  $\Delta \Omega$ ), rho\_min, and Propagation\_time. The VisMashup App generates a new image in the online article (in accordance with the workflow shown in Figure 1) if the reader selects a different set of parameters and clicks the green "Update" button. Clicking on the red "Execute on my desktop" button downloads the Figure 1 workflow to the reader's computer system for local execution.

far test particles residing in different model parameters outside those origiflow field regions will travel in a given nally discussed. sity is unity.)

on in the original printed article. By little additional time, we can bring to a modern IVAJ.) actively adjusting one or more of the the original figures to life and reap Following the local execution key model parameter values and us- additional benefits from our original of Figure 1's workflow using the ing the embedded VisMashup App to code-development efforts. generate a new figure based on those It's important to note that each in Figure 2, the App displays the values, users likely will gain a better time a user changes a parameter value rendered configuration outside the conclusions. Further, using the Wiki's performs the requested analysis on spreadsheet. (The initial download

value. As the article's "SwitchCoord satisfactorily analyze the underlying Python Module" sidebar describes, properties of the flow that resulted rho\_min is an additional parameter from our astrophysical fluid simula- computers. Of course, they can real-

appreciation of our original article's and executes the VisMashup App, it browser, in one cell of a VisTrails standard editing features, users can the original model data. That is, we've and execution can take 10 minutes or

archived the original astrophysical fluid simulation's model data to support our effort to enhance the article's content. This is a step in the right direction, as efforts to demonstrate the reproducibility of largescale numerical simulations aren't likely to succeed until the computational sciences community makes a commitment to archive simulation results. Our IVAJ-formatted article with Level 2 enhancements illustrates how such archival data can naturally enrich the content of published journal articles.

### Level 3 Enhancements

As the example at www.vistrails. as a VisTrails workflow parameter (see comment on the insights they've org/index.php/User:Tohline/IVAJ/ Figure 1) that we use to examine how gained from examining a range of Levels2and3 shows, our IVAJ article offers yet another enhancement level over traditional journal articles. By amount of time; in general, the collec- We invested considerable time in clicking the red "Execute on my tion of streamlines will shorten if we our original article, piecing together desktop" button displayed within the specify a smaller propagation\_time a visualization workflow that let us Figure 2 window of the VisMashup App, users can execute Figure 1's VisTrails workflow on their own that the customized Python module tion. It's not unusual for computa- ize this Level 3 enhancement only if uses; individual streamlines are trun- tional sciences researchers to invest they've previously installed VisTrails cated once the test particle traveling such time on postprocessing analysis (version 1.4.2 or later) as a functionalong that streamline enters a region (especially on visualization tasks). In ing application on their local system. where the gas density is less than rho\_ the original article, we captured the (VisTrails is an open source applicamin. (Densities have been normalized scientific fruits of this labor in two tion designed to run under a wide such that the model's maximum den- static images (Figures 2 and 3). Our range of operating systems, so we embedded VisMashup App executes hope this local installation require-This Level 2 enhancement lets us- exactly the same visualization work- ment won't discourage readers from ers examine more thoroughly the flow as the original article. Hence, exploring and considering the added astrophysical model that we focused with the investment of relatively value that such applications can bring

model parameters initially displayed





### VISUALIZATION CORNER Figure 3 Figure 2 amega,fame -8.845 | + ma\_min 0.001 propagation time 3.2 Update Figure 2. The VisMashup window that displays when users select the "Figure 2" tab (see www.vistrails.org/index.php/User:Tohline/IVAJ/Levels2and3). The window

displays an image generated by a customized VisTrails workflow using the indicated values of the three variable parameters, Omega\_frame (=  $\Delta \Omega$ ), rho\_min, and Propagation\_time. The VisMashup App generates a new image in the online article (in accordance with the workflow shown in Figure 1) if the reader selects a different set of parameters and clicks the green "Update" button. Clicking on the red "Execute on my desktop" button downloads the Figure 1 workflow to the reader's computer system for local execution.

flow field regions will travel in a given nally discussed. sity is unity.)

ing the embedded VisMashup App to code-development efforts. generate a new figure based on those It's important to note that each in Figure 2, the App displays the values, users likely will gain a better time a user changes a parameter value rendered configuration outside the conclusions. Further, using the Wiki's performs the requested analysis on spreadsheet. (The initial download standard editing features, users can the original model data. That is, we've and execution can take 10 minutes or

value. As the article's "SwitchCoord satisfactorily analyze the underlying Python Module" sidebar describes, properties of the flow that resulted rho\_min is an additional parameter from our astrophysical fluid simula- computers. Of course, they can realcated once the test particle traveling such time on postprocessing analysis (version 1.4.2 or later) as a functionalong that streamline enters a region (especially on visualization tasks). In ing application on their local system. embedded VisMashup App executes hope this local installation requireon in the original printed article. By little additional time, we can bring to a modern IVAJ.) actively adjusting one or more of the the original figures to life and reap Following the local execution key model parameter values and us- additional benefits from our original of Figure 1's workflow using the

appreciation of our original article's and executes the VisMashup App, it browser, in one cell of a VisTrails

archived the original astrophysical fluid simulation's model data to support our effort to enhance the article's content. This is a step in the right direction, as efforts to demonstrate the reproducibility of largescale numerical simulations aren't likely to succeed until the computational sciences community makes a commitment to archive simulation results. Our IVAJ-formatted article with Level 2 enhancements illustrates how such archival data can naturally enrich the content of published journal articles.

### Level 3 Enhancements

As the example at www.vistrails. as a VisTrails workflow parameter (see comment on the insights they've org/index.php/User:Tohline/IVAI/ Figure 1) that we use to examine how gained from examining a range of Levels2and3 shows, our IVAJ article far test particles residing in different model parameters outside those origi- offers yet another enhancement level over traditional journal articles. By amount of time; in general, the collec- We invested considerable time in clicking the red "Execute on my tion of streamlines will shorten if we our original article, piecing together desktop" button displayed within the specify a smaller propagation\_time a visualization workflow that let us Figure 2 window of the VisMashup App, users can execute Figure 1's VisTrails workflow on their own that the customized Python module tion. It's not unusual for computa- ize this Level 3 enhancement only if uses; individual streamlines are trun- tional sciences researchers to invest they've previously installed VisTrails where the gas density is less than rho\_ the original article, we captured the (VisTrails is an open source applicamin. (Densities have been normalized scientific fruits of this labor in two tion designed to run under a wide such that the model's maximum den- static images (Figures 2 and 3). Our range of operating systems, so we This Level 2 enhancement lets us- exactly the same visualization work- ment won't discourage readers from ers examine more thoroughly the flow as the original article. Hence, exploring and considering the added astrophysical model that we focused with the investment of relatively value that such applications can bring

model parameters initially displayed





### VISUALIZATION CORNER Figure 3 Figure 2 amega,fame -8.845 | + ma\_min 0.001 propagation time 3.2 Figure 2. The VisMashup window that displays when users select the "Figure 2" tab (see www.vistrails.org/index.php/User:Tohline/IVAJ/Levels2and3). The window

displays an image generated by a customized VisTrails workflow using the indicated values of the three variable parameters,  $Omega_frame (= \Delta \Omega)$ , rho\_min, and Propagation\_time. The VisMashup App generates a new image in the online article (in accordance with the workflow shown in Figure 1) if the reader selects a different set of parameters and clicks the green "Update" button. Clicking on the red "Execute on my desktop" button downloads the Figure 1 workflow to the reader's computer system for local execution.

flow field regions will travel in a given nally discussed. sity is unity.)

ing the embedded VisMashup App to code-development efforts. generate a new figure based on those It's important to note that each in Figure 2, the App displays the values, users likely will gain a better time a user changes a parameter value rendered configuration outside the conclusions. Further, using the Wiki's performs the requested analysis on spreadsheet. (The initial download standard editing features, users can the original model data. That is, we've and execution can take 10 minutes or

specify a smaller propagation\_time a visualization workflow that let us Figure 2 window of the VisMashup value. As the article's "SwitchCoord satisfactorily analyze the underlying Python Module" sidebar describes, properties of the flow that resulted rho\_min is an additional parameter from our astrophysical fluid simulacated once the test particle traveling such time on postprocessing analysis (version 1.4.2 or later) as a functionalong that streamline enters a region (especially on visualization tasks). In ing application on their local system. where the gas density is less than rho\_ the original article, we captured the (VisTrails is an open source applicaembedded VisMashup App executes hope this local installation require-This Level 2 enhancement lets us- exactly the same visualization work- ment won't discourage readers from ers examine more thoroughly the flow as the original article. Hence, exploring and considering the added on in the original printed article. By little additional time, we can bring to a modern IVAJ.) actively adjusting one or more of the the original figures to life and reap key model parameter values and us- additional benefits from our original of Figure 1's workflow using the

appreciation of our original article's and executes the VisMashup App, it browser, in one cell of a VisTrails

archived the original astrophysical fluid simulation's model data to support our effort to enhance the article's content. This is a step in the right direction, as efforts to demonstrate the reproducibility of largescale numerical simulations aren't likely to succeed until the computational sciences community makes a commitment to archive simulation results. Our IVAJ-formatted article with Level 2 enhancements illustrates how such archival data can naturally enrich the content of published journal articles.

### Level 3 Enhancements

As the example at www.vistrails. as a VisTrails workflow parameter (see comment on the insights they've org/index.php/User:Tohline/IVAJ/ Figure 1) that we use to examine how gained from examining a range of Levels2and3 shows, our IVAJ article far test particles residing in different model parameters outside those origi- offers yet another enhancement level over traditional journal articles. By amount of time; in general, the collec- We invested considerable time in clicking the red "Execute on my tion of streamlines will shorten if we our original article, piecing together desktop" button displayed within the App, users can execute Figure 1's VisTrails workflow on their own computers. Of course, they can realthat the customized Python module tion. It's not unusual for computa- ize this Level 3 enhancement only if uses; individual streamlines are trun- tional sciences researchers to invest they've previously installed VisTrails min. (Densities have been normalized scientific fruits of this labor in two tion designed to run under a wide such that the model's maximum den- static images (Figures 2 and 3). Our range of operating systems, so we astrophysical model that we focused with the investment of relatively value that such applications can bring

Following the local execution model parameters initially displayed







### VISUALIZATION CORNER Figure 3 Figure 2 amaga,fama -8.845 (\* ma\_min 0.001 propagation time 3.2 Figure 2. The VisMashup window that displays when users select the "Figure 2" tab (see www.vistrails.org/index.php/User:Tohline/IVAJ/Levels2and3). The window

displays an image generated by a customized VisTrails workflow using the indicated values of the three variable parameters,  $Omega_frame (= \Delta \Omega)$ , rho\_min, and Propagation\_time. The VisMashup App generates a new image in the online article (in accordance with the workflow shown in Figure 1) if the reader selects a different set of parameters and clicks the green "Update" button. Clicking on the red "Execute on my desktop" button downloads the Figure 1 workflow to the reader's computer system for local execution.

flow field regions will travel in a given nally discussed. sity is unity.)

ing the embedded VisMashup App to code-development efforts. generate a new figure based on those It's important to note that each in Figure 2, the App displays the values, users likely will gain a better time a user changes a parameter value rendered configuration outside the conclusions. Further, using the Wiki's performs the requested analysis on spreadsheet. (The initial download standard editing features, users can the original model data. That is, we've and execution can take 10 minutes or

specify a smaller propagation\_time a visualization workflow that let us Figure 2 window of the VisMashup value. As the article's "SwitchCoord satisfactorily analyze the underlying Python Module" sidebar describes, properties of the flow that resulted rho\_min is an additional parameter from our astrophysical fluid simulacated once the test particle traveling such time on postprocessing analysis (version 1.4.2 or later) as a functionalong that streamline enters a region (especially on visualization tasks). In ing application on their local system. embedded VisMashup App executes hope this local installation require-This Level 2 enhancement lets us- exactly the same visualization work- ment won't discourage readers from on in the original printed article. By little additional time, we can bring to a modern IVAJ.) actively adjusting one or more of the the original figures to life and reap Following the local execution key model parameter values and us- additional benefits from our original of Figure 1's workflow using the

appreciation of our original article's and executes the VisMashup App, it browser, in one cell of a VisTrails

archived the original astrophysical fluid simulation's model data to support our effort to enhance the article's content. This is a step in the right direction, as efforts to demonstrate the reproducibility of largescale numerical simulations aren't likely to succeed until the computational sciences community makes a commitment to archive simulation results. Our IVAJ-formatted article with Level 2 enhancements illustrates how such archival data can naturally enrich the content of published journal articles.

### Level 3 Enhancements

As the example at www.vistrails. as a VisTrails workflow parameter (see comment on the insights they've org/index.php/User:Tohline/IVAJ/ Figure 1) that we use to examine how gained from examining a range of Levels2and3 shows, our IVAJ article far test particles residing in different model parameters outside those origi- offers yet another enhancement level over traditional journal articles. By amount of time; in general, the collec- We invested considerable time in clicking the red "Execute on my tion of streamlines will shorten if we our original article, piecing together desktop" button displayed within the App, users can execute Figure 1's VisTrails workflow on their own computers. Of course, they can realthat the customized Python module tion. It's not unusual for computa- ize this Level 3 enhancement only if uses; individual streamlines are trun- tional sciences researchers to invest they've previously installed VisTrails where the gas density is less than rho\_ the original article, we captured the (VisTrails is an open source applicamin. (Densities have been normalized scientific fruits of this labor in two tion designed to run under a wide such that the model's maximum den- static images (Figures 2 and 3). Our range of operating systems, so we ers examine more thoroughly the flow as the original article. Hence, exploring and considering the added astrophysical model that we focused with the investment of relatively value that such applications can bring

model parameters initially displayed







### VISUALIZATION CORNER Figure 3 Figure 2 amega,fame -8.845 | + ma\_min 0.001 propagation time 3.2 Figure 2. The VisMashup window that displays when users select the "Figure 2" tab (see www.vistrails.org/index.php/User:Tohline/IVAJ/Levels2and3). The window

displays an image generated by a customized VisTrails workflow using the indicated values of the three variable parameters,  $Omega_frame (= \Delta \Omega)$ , rho\_min, and Propagation\_time. The VisMashup App generates a new image in the online article (in accordance with the workflow shown in Figure 1) if the reader selects a different set of parameters and clicks the green "Update" button. Clicking on the red "Execute on my desktop" button downloads the Figure 1 workflow to the reader's computer system for local execution.

flow field regions will travel in a given nally discussed. Python Module" sidebar describes, properties of the flow that resulted sity is unity.)

ing the embedded VisMashup App to code-development efforts. generate a new figure based on those It's important to note that each in Figure 2, the App displays the values, users likely will gain a better time a user changes a parameter value rendered configuration outside the conclusions. Further, using the Wiki's performs the requested analysis on spreadsheet. (The initial download

amount of time; in general, the collec- We invested considerable time in tion of streamlines will shorten if we our original article, piecing together desktop" button displayed within the specify a smaller propagation\_time a visualization workflow that let us value. As the article's "SwitchCoord satisfactorily analyze the underlying rho\_min is an additional parameter from our astrophysical fluid simulathat the customized Python module tion. It's not unusual for computacated once the test particle traveling such time on postprocessing analysis (version 1.4.2 or later) as a functionalong that streamline enters a region (especially on visualization tasks). In ing application on their local system. embedded VisMashup App executes hope this local installation require-This Level 2 enhancement lets us- exactly the same visualization workon in the original printed article. By little additional time, we can bring to a modern IVAJ.) actively adjusting one or more of the the original figures to life and reap key model parameter values and us- additional benefits from our original

appreciation of our original article's and executes the VisMashup App, it browser, in one cell of a VisTrails standard editing features, users can the original model data. That is, we've and execution can take 10 minutes or

archived the original astrophysical fluid simulation's model data to support our effort to enhance the article's content. This is a step in the right direction, as efforts to demonstrate the reproducibility of largescale numerical simulations aren't likely to succeed until the computational sciences community makes a commitment to archive simulation results. Our IVAJ-formatted article with Level 2 enhancements illustrates how such archival data can naturally enrich the content of published journal articles.

### Level 3 Enhancements

As the example at www.vistrails. as a VisTrails workflow parameter (see comment on the insights they've org/index.php/User:Tohline/IVAJ/ Figure 1) that we use to examine how gained from examining a range of Levels2and3 shows, our IVAJ article far test particles residing in different model parameters outside those origi- offers yet another enhancement level over traditional journal articles. By clicking the red "Execute on my Figure 2 window of the VisMashup App, users can execute Figure 1's VisTrails workflow on their own computers. Of course, they can realize this Level 3 enhancement only if uses; individual streamlines are trun- tional sciences researchers to invest they've previously installed VisTrails where the gas density is less than rho\_ the original article, we captured the (VisTrails is an open source applicamin. (Densities have been normalized scientific fruits of this labor in two tion designed to run under a wide such that the model's maximum den- static images (Figures 2 and 3). Our range of operating systems, so we ment won't discourage readers from ers examine more thoroughly the flow as the original article. Hence, exploring and considering the added astrophysical model that we focused with the investment of relatively value that such applications can bring

Following the local execution of Figure 1's workflow using the model parameters initially displayed







44

# Challenges

- Re-using results
- Adding results to publications
- Obtaining results, computations, and input from publications
- Publishing interactive experiments
- Searching executable paper collections
- Reviewers: execution environments, checking different parameters
- Longevity/maintenance
- Resource constraints:
  - analyses run on supercomputers
  - large datasets
  - privacy or intellectual property concerns





45

# General Strategies for Reproducibility

- Preserving the Mess:
  - Just save a virtual machine
  - Trace dependencies
- Encouraging Cleanliness:
  - Use a system (e.g. Umbrella, VisTrails)
  - Use literate programming environments
  - Use code and data repositories
  - Use packaging system (ReproZip)

### D. Koop, CSCI 680/490, Spring 2021

[Categories from H. Meng et al., 2016]



Northern Illinois University







## Literate Programming

- Knuth's WEB system
- Mathematica
- Code this is well-documented using comments
- Jupyter Notebooks





## Data and Code Availability

- Code Repositories:
  - GitHub
  - GitLab

- ...

- Data Repositories:
  - figshare, freebase, dryad, DataONE
  - Also many domain-specific repositories
  - http://oad.simmons.edu/oadwiki/Data\_repositories





# 10 Rules for Reproducible Computational Research

- Rule 1: For Every Result, Keep Track of How It Was Produced
- Rule 2: Avoid Manual Data Manipulation Steps
- Rule 3: Archive the Exact Versions of All External Programs Used
- Rule 4: Version Control All Custom Scripts
- Rule 5: Record All Intermediate Results, When Possible in Standardized Formats







49

# 10 Rules for Reproducible Computational Research

- Rule 6: For Analyses That Include Randomness, Note Underlying Random Seeds
- Rule 7: Always Store Raw Data behind Plots
- Rule 8: Generate Hierarchical Analysis Output, Allowing Layers of Increasing Detail to Be Inspected
- Rule 9: Connect Textual Statements to Underlying Results • Rule 10: Provide Public Access to Scripts, Runs, and Results









## Rules or Benefits?

- Laws to make sure people don't cheat or lie or steal Is that a good incentive? You won't be mislabeled as a criminal?
- Benefits of Reproducibility
  - Reproducible programs can be compared
  - Reproducible software and results are documented
  - Reproducible software is portable
  - Reproducible experiments are cited









# Reproducible Experiments Classification

- Depth: how much is available?
  - figures
  - scripts
  - raw data
  - experiments
  - software system
- Portability: what machine specs are necessary?
  - same machine
  - similar machine
  - different OS
- Coverage: how much can be reproduced?











# (Database) Research Topics

- Design and Management of Experiment Repositories
- Querying and Searching Experiments
- Mining Experiments











## A Large-scale Study about Quality and Reproducibility of Jupyter Notebooks

J. F. Pimentel, L. Murta, V. Braganholo, and J. Freire





## Notebooks and Hidden State



D. Koop, CSCI 680/490, Spring 2021

In [1]: co = 0 In [1]: co = 0 In [1]: co = 0In [3]: co += 1 In [2]: co += 2 In [4]: In [3]: co In [3]: co CO Out[4]: 2 Out[3]: 1 Out[3]: 1

[Pimentel et al., 2019]












## Notebook Composition



### D. Koop, CSCI 680/490, Spring 2021



[Pimentel et al., 2019]











# Notebook Reproducibility

- Use notebooks from Github (~1 million) - Unambiguous cell order? 81.99%
- Study notebook dependencies
  - Dependencies Available? 13.72%
  - Dependencies Install? 5.03%
- Study notebook executability
  - Execute: 24.11% of unambiguous cell order
  - Matched results: 4.03%

## D. Koop, CSCI 680/490, Spring 2021











## Best Practices

- Use short titles with a restrict charset (A-Z a-z 0-9 . -) for notebook files and markdown headings for more detailed ones in the body
- Pay attention to the bottom of the notebook. Check whether it can benefit from descriptive markdown cells or can have code cells executed or removed
- Abstract code into functions, classes, and modules and test them
- Declare the dependencies in requirement files & pin versions of all packages
- Use a clean environment to test if dependencies are properly declared
- Put imports at the beginning of notebooks
- Use relative paths for accessing data in the repository
- Re-run notebooks top to bottom before committing











# Problem: What is df at any point in time?

In [5]:	<pre>import pandas as pd df = pd.read_csv('guardian-top100-female-2019.csv')</pre>												
Out[5]:		Name	Rank	Position	Age	on 1 Dec 201	9	Nationality					
	0	Sam Kerr	1	Forward		2	26	Australia					
	···· ··· ···												
	99	Ludmila	100	Forward		2	25	Brazil					
	100	rows × 5 co	lumns										
In [6]:	df	= df.r	enam	e(colur	nns=	{'Age on	1	Dec 201	9': 'Age'})				
Out[6]:		Name	Rank	Position	Age	Nationality							
	0	Sam Kerr	1	Forward	26	Australia							
	99	Ludmila	100	Forward	25	Brazil							
	100 rows × 5 columns												

D. Koop, CSCI 680/490, Spring 2021

In [3]:	df	= df[df.Ag	ſe >=	31]		
Out[3]:		Name	Rank	Position	Age	Nationality
	2	Megan Rapinoe	3	Midfielder	34	USA
	96	Cláudia Neto	97	Midfielder	31	Portugal
	19 rc	ows × 5 columns				
In [7]:	df	= df[df.Ag	ſe <=	24]		
Out[7]:		Name	Rank	Position	Age	Nationality
	3	Ada Hegerberg	4	Forward	24	Norway
	•••					
	98	Lena Oberdorf	99	Midfielder	17	Germany
	25 rc	ows × 5 columns				





59



# Dataflow Notebooks: Resolve Notebook Ambiguities

In [d51f8eab]:	<pre>import pandas as pd df = pd.read_csv('guardian-top100-female-2019.csv')</pre>					In [over30]:	df = df <b>\$full</b> [df <b>\$full</b> .Age >= 31]						
df:	Name	Rank Position	Age on 1 Dec 2019	df:	Name Rank Position Age Nationality								
	0 Sam Kerr	1 Forward	26	Australia			2 Me	egan Rapinoe	3	Midfielder	34	USA	
	99 Ludmila	100 Forward	25	Brazil			96 (	Cláudia Neto	97 I	Midfielder	31	Portugal	
	100 rows × 5 co	lumns					19 rows :	× 5 columns					
In [full]:	df = df.1	cename(colur	nns={'Age on 1	Dec 2019	<pre>9': 'Age'})</pre>	In [under25]:	df =	df <mark>\$full</mark> [d	df <mark>\$fu</mark>	ll.Age	<= 2	24]	
df.	Name	Rank Position	df:		Name I	Rank	Position	Age N	Nationality				
ur.	0 Sam Kerr	1 Forward	26 Australia				3 Ada	a Hegerberg	4	Forward	24	Norway	
	99 Ludmila	100 Forward	25 Brazil				98 Le	ena Oberdorf	99 N	/lidfielder	17	Germany	
	100 rows × 5 co	lumns					25 rows :	× 5 columns					

## D. Koop, CSCI 680/490, Spring 2021







# Dataflow Notebooks: Dependency Graph



### D. Koop, CSCI 680/490, Spring 2021

- Shows connections between cells
- Can see which cells would be affected by a change
- Same colors indicate which parts of the graph are stale
  - Linked to the notebook
  - Hover to show a cell's code
  - Can also execute in the graph







