Advanced Data Management (CSCI 640/490)

Reproducibility

Dr. David Koop





Provenance in Computational Science



D. Koop, CSCI 640/490, Spring 2024





2

Provenance Capture Mechanisms

- Workflow-based: Since workflow execution is controlled, keep track of all the workflow modules, parameters, etc. as they are executed
- Process-based: Each process is required to write out its own provenance information (not centralized like workflow-based)
- **OS-based**: The OS or filesystem is modified so that any activity it does it monitored and the provenance subsystem organizes it
- Tradeoffs:
 - Workflow- and process-based have better abstraction
 - OS-based requires minimal user effort once installed and can capture "hidden dependencies"









Prospective and Retrospective Provenance

- Prospective provenance is what was specified/intended
 - a workflow, script, list of steps
- Retrospective provenance is what actually happened
 - actual data, actual parameters, errors that occurred, timestamps, machine information
- Do not need prospective provenance to have retrospective provenance!
- Recipe for a cake vs. Baking a cake





PROV: Three Key Classes



An **entity** is a physical, digital, conceptual, or other kind of thing with some fixed aspects; entities may be real or imaginary.



An **activity** is something that occurs over a period of time and acts upon or with entities; it may include consuming, processing, transforming, modifying, relocating, using, or generating entities.



An **agent** is something that bears some form of responsibility for an activity taking place, for the existence of an entity, or for another agent's activity.





More provenance

- Database Provenance
- Evolution Provenance
- Provenance for Data Science









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

	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	

Q1:

SELECT a.name, a.phone

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

Result of Q_1 :

V =	
name	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

	0		
	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 Result of 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) \boldsymbol{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 :		Q_1' :	
SELECT	a.name, a.phone	SELECT	e.r
FROM	Agencies a , ExternalTours e	FROM	Ag
WHERE	a.name = e.name	WHERE	a.
	AND $e.type='boat'$		AI

name, a.phone gencies a, ExternalTours ename = e.nameND 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















D. Koop, CSCI 640/490, Spring 2024

Evolution Provenance





Data Exploration



D. Koop, CSCI 640/490, Spring 2024

[Modified from Van Wijk, Vis 2005]







Data Exploration



- Data analysis and visualization are iterative processes
- In exploratory tasks, change is the norm!

D. Koop, CSCI 640/490, Spring 2024

[Modified from Van Wijk, Vis 2005]







Exploration and Creativity Support

- Reasoning is key to the exploratory processes
- "Reflective reasoning requires the ability to store temporary results, to make inferences from stored knowledge, and to follow chains of reasoning backward and forward, sometimes backtracking when a promising line of thought proves to be unfruitful. ...the process is slow and laborious" — Donald A. Norman
- Need external aids—tools to facilitate this process
 "Creativity support tools" —Ben Shneiderman
- Need aid from people—collaboration

D. Koop, CSCI 640/490, Spring 2024





13

Change-based Provenance: Photo Editing

User Actions



Undo/Redo History







Change-based Provenance: Photo Editing

• User Actions



Undo/Redo History







Version Trees

- 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











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





Workflow Evolution Provenance







Workflow Evolution Provenance

(



D. Koop, CSCI 640/490, Spring 2024



delete module "GMapCell"

delete module "CellLocation"

delete module "ProjectTable"

delete module "SelectFromTable"

•••

. . .

add module "SelectFromTable"

add parameter "float_expr" to "SelectFromTable" with value "latitutde > 40.6"

delete parameter "float_expr" from "SelectFromTable"

add parameter "float_expr" to "SelectFromTable" with value "latitutde > 40.7"

delete parameter "float_expr" from "SelectFromTable"

add parameter "float_expr" to "SelectFromTable" with value "latitutde > 40.8"







Execution Provenance







Execution Provenance

```
<module id="12" name="vtkDataSetReader"
        start time="2010-02-19 11:01:05"
        end time="2010-02-19 11:01:07">
 <annotation key="hash"</pre>
            value="c54bea63cb7d912a43ce"/>
</module>
<module id="13" name="vtkContourFilter"
        start time="2010-02-19 11:01:07"
        end time="2010-02-19 11:01:08"/>
<module id="15" name="vtkDataSetMapper"
        start time="2010-02-19 11:01:09"
        end time="2010-02-19 11:01:12"/>
<module id="16" name="vtkActor"
        start time="2010-02-19 11:01:12"
        end time="2010-02-19 11:01:13"/>
<module id="17" name="vtkCamera"
        start time="2010-02-19 11:01:13"
        end time="2010-02-19 11:01:14"/>
<module id="18" name="vtkRenderer"
        start time="2010-02-19 11:01:14"
        end time="2010-02-19 11:01:14"/>
• • •
```







Capturing and querying fine-grained provenance of preprocessing pipelines in data science

A. Chapman, P. Missier, L. Lauro, R. Torlone





Data Provenance for Data Science

	CId	Gender	Age	Zip	ageRange
1	113	F	24	98567	young
2	241	M	28	\bot	adult
3	375	C	\bot	32768	
4	578	F	44	32768	adult



D. Koop, CSCI 640/490, Spring 2024



Northern Illinois University







Provenance Templates









<u>Assignment 5</u>

- Divvy Bikes Data
- Spatial, Graph, and Temporal Data Processing
- Use pandas, geopandas, neo4j, (modin for extra credit)









Final Exam

- Wednesday, May 8, 8:00-9:50pm, PM 252
- Similar format
- More comprehensive (questions from topics covered in Test 1 & 2)
- Will also have questions from graph/spatial/temporal data, provenance, reproducibility, machine learning









<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



D. Koop, CSCI 640/490, Spring 2024





Northern Illinois University





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 640/490, Spring 2024

OK ^{≤30} OK 130 64	> ³⁰ OK ^{Auth} 23	
	Notation	Number of papers
\mathcal{J}	HW	excluded due to replication requiring special hardware
Build	NC	excluded due to results not being backed by code
fails	EX	excluded due to overlapping author lists
9	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 ^{no}	where the author responds to an email message saying code cannot be provided
	EM ^ø	where the author does not respond to email requests within two months
	OK ^{≤30}	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 ^{Auth}	where code is available and we fail to build, and the author says the code builds with reasonable effort
$\begin{array}{ccc} \mathbf{M}^{\emptyset} & \mathbf{E}\mathbf{M}^{\mathrm{no}} \\ 30 & 146 \end{array}$	Fails	where code is available and we fail to build, and the author says the code may have problems building









28

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 640/490, Spring 2024





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	ի	lr	
	Γ)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	 email address and/or web site
Resource	 types: code, data, media, documenta availability: no access, access, NDA expense: free, non-free, free for acad distribution form: source, binary, se expiration date license comment
Support	 kinds: resolve installation issues, fix upgrade to new language and operat system versions, port to new environ improve performance, add features expense: free, non-free, free for acad expiration date

D. Koop, CSCI 640/490, Spring 2024

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













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. 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 realsity 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

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

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 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 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 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 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/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 App, users can execute Figure 1's VisTrails workflow on their own

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

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 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 simulathat 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 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 App, users can execute Figure 1's VisTrails workflow on their own computers. Of course, they can real-

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

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 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 simulathat 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 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 App, users can execute Figure 1's VisTrails workflow on their own computers. Of course, they can real-

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

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 Python Module" sidebar describes, properties of the flow that resulted rho_min is an additional parameter from our astrophysical fluid simulathat the customized Python module tion. It's not unusual for computauses; 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 workers 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 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

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 ment won't discourage readers from

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









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





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 640/490, Spring 2024

[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









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?

D. Koop, CSCI 640/490, Spring 2024







