Advanced Data Management (CSCI 490/680)

Reproducibility

Dr. David Koop





<u>Assignment 5</u>



D. Koop, CSCI 490/680, Spring 2020

- Work with time series & spatial data
- Shorter assignment
- Cleaning, spatial rollup, rolling average
- Due April 30
- Questions?





2

Final Exam and Review

- Final Exam
 - Tuesday, May 5 from 4-5:50pm
 - Online
 - Similar format to Test 2
 - Comprehensive but with more focus on last few weeks of class
- Review
 - Thursday, April 30
 - Submit questions via email or discussion

D. Koop, CSCI 490/680, Spring 2020





3

Reading Quiz on Tuesday

- Final reading quiz next Tuesday
- <u>SageDB paper</u>





Provenance in Computational 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

	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 :

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

	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) 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
<i>t</i> 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.na
FROM	Agencies a , ExternalTours e	FROM	Age
WHERE	a.name = e.name	WHERE	a.na
	AND $e.type='boat'$		AN

ame, a.phone encies a, ExternalTours eame = e.nameD 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

































D. Koop, CSCI 490/680, Spring 2020







13

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









Visualization Pipeline Completions











Visualization by Analogy







Visualization by Analogy







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

Sheet 1 PE#0 critical_points.vt 0









<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 490/680, Spring 2020

$OK^{\leq 30}$ $OK^{>}_{64}$	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
	EMø	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
${f M}^{\emptyset} {f E} {f M}^{ m no} {f 146}$	Fails	where code is available and we fail to build, and the author says the code may have problems building







21



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 490/680, Spring 2020





Northern Illinois University







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







Recommendations

- Fund repeatability engineering
- Require sharing contracts

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

D. Koop, CSCI 490/680, Spring 2020

ation access Idemics ervice

bugs, ting iments,

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 standard editing features, users can the original model data. That is, we've and execution can take 10 minutes or

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 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









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

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

model parameters initially displayed









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

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 amount of time; in general, the collec- We invested considerable time in clicking the red "Execute on my 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











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

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 amount of time; in general, the collec- We invested considerable time in clicking the red "Execute on my App, users can execute Figure 1's VisTrails workflow on their own computers. Of course, they can real-

model parameters initially displayed











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 490/680, Spring 2020

[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?

(Database) Research Topics

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

D. Koop, CSCI 490/680, Spring 2020

40

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 490/680, Spring 2020

In [1]: co = 0In [1]: co = 0 In [1]: co = 0In [2]: co += 2 In [3]: co += 1 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 490/680, Spring 2020

[Pimentel et al., 2019]

43

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%

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

In [a0a358]:	raw_df
raw_df:	
	0 1 2
In [aaa3c6]:	column_
column_names:	{0: 'Na:
In [a249ea]:	named_d
named_df:	
	0 1 2
In [aab079]:	named_d
Out[aab079]:	Country USA Canada Brazil

David Koop

<pre>= pd.read_csv</pre>	/("fifal	7-top2	<pre>20-women.txt",sep="-",header=None)</pre>
0	1	2	
Caroline Seger	Sweden	85	
Wendie Renard	France	85	
Steph Houghton	England	85	
_names = {0: "	'Name",	1: "Co	<pre>ountry", 2: "Rating"}</pre>
ame', 1: 'Coun	try', 2	: 'Rat	ing'}
df = raw_df.re	ename(co	lumns=	column_names)
Name	Country	Rating	
Caroline Seger	Sweden	85	
Wendie Renard	France	85	

df.groupby("Country").size().sort_values(ascending=False)

85

...

...

Steph Houghton England

...

У			
	6		
	3		
	3		
•	• •		[D. Koop
		\sim	

David Koop

<pre>= pd.read_csv</pre>	/("fifal	7-top2	<pre>20-women.txt",sep="-",header=None)</pre>
0	1	2	
Caroline Seger	Sweden	85	
Wendie Renard	France	85	
Steph Houghton	England	85	
_names = {0: "	'Name",	1: "Co	<pre>ountry", 2: "Rating"}</pre>
ame', 1: 'Coun	try', 2	: 'Rat	ing'}
df = raw_df.re	ename(co	lumns=	column_names)
Name	Country	Rating	
Caroline Seger	Sweden	85	
Wendie Renard	France	85	

named_df.groupby("Country").size().sort_values(ascending=False)

85

...

Steph Houghton

...

England

...

У			
	6		
	3		
	3		
•	• •		[D. Koop
		\sim	

- Persistent Identifiers
- Named Outputs

Steph Houghton

...

England

...

David Koop

<pre>= pd.read_csv</pre>	/("fifal	7-top2	<pre>20-women.txt",sep="-",header=None)</pre>
0	1	2	
Caroline Seger	Sweden	85	
Wendie Renard	France	85	
Steph Houghton	England	85	
_names = {0: "	'Name",	1: "Co	<pre>ountry", 2: "Rating"}</pre>
ame', 1: 'Coun	try', 2	: 'Rat	ing'}
df = raw_df.re	ename(co	lumns=	column_names)
Name	Country	Rating	
Caroline Seger	Sweden	85	
Wendie Renard	France	85	

In [aab079]: named_df.groupby("Country").size().sort_values(ascending=False)

85

...

У			
	6		
	3		
	3		
•	• •		[D. Koop
		\sim	

- Persistent Identifiers
- Named Outputs
- Unnamed Outputs

In [a0a358]:	raw_df
raw_df:	
	0
	1
	2
In [aaa3c6]:	column_
column_names:	{0: 'Na
In [a249ea]:	named_c
named_df:	
	0
	1
	2
In [aab079]:	named_c
Out[aab079]:	Country USA Canada
	•••

Steph Houghton

...

England

...

<pre>= pd.read_csv</pre>	/("fifal	7-top2	<pre>20-women.txt",sep="-",header=None)</pre>
0	1	2	
Caroline Seger	Sweden	85	
Wendie Renard	France	85	
Steph Houghton	England	85	
_names = {0: "	'Name",	1: "Co	<pre>ountry", 2: "Rating"}</pre>
ame', 1: 'Coun	try', 2	: 'Rat	ing'}
df = raw_df.re	ename(co	lumns=	column_names)
Name	Country	Rating	
Caroline Seger	Sweden	85	
Wendie Renard	France	85	

df.groupby("Country").size().sort_values(ascending=False)

85

...

У			
	6		
	3		
	3		
•	• •		[D. Koop
		\sim	

- Persistent Identifiers
- Named Outputs
- Unnamed Outputs
- Connection by Variable Reference

In [a0a358]:	raw_df
raw_df:	
	0 1 2
In [aaa3c6]:	column_
column_names:	{0: 'Na
In [a249ea]:	named_c
named_df:	
	0 1 2
In [aab079]:	named_c
Out[aab079]:	Country USA Canada Brazil

• • •

= pd.read_csv	("fifal	7-top
0	1	2
Caroline Seger	Sweden	85
Wendie Renard	France	85
Steph Houghton	England	85
Oteph Houghton	Ligidia	00
•••	•••	
names = {0: "	Name",	1: "Co
ame', 1: 'Coun	try', 2	: 'Rat
df = raw_df.re	name(co	lumns
Name	Country	Rating
Caroline Seger	Sweden	85
Wendie Renard	France	85
Steph Houghton	England	85
	0	
df.groupby("Co	untry")	.size
	<u> </u>	
У		
6		
3		
3		

[D. Koop et al.]

Dataflow Notebooks: Dependency Graph

David Koop

- 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

47