### Advanced Data Management (CSCI 490/680)

### Data Citation

Dr. David Koop





### What is Data?



Less than 0 Change for U.S.: 32,712,033

D. Koop,

- here

to matin I chlorent gamin blin nent and delition & MIT it S'an . It fills to phile in the addition of the the method

new house is o live in ey want him plastering. He para gitti.

Pisa Griffin

{much money went} Has a tractor.

### Date: July 1980 Place:Sakaltutan Zafor:

Household now Zafor and wife; Nazif Unal and wife and youngest son, still a boy. They run two dolmuß; one with a driver from Süleymanti. Goes in and out once a day. He gets 8,000 a month. Zafor then said, keskin deoil. { not sharp - i.e.? not profitable} I said he did very well on 8,000 TL with only two journeys a day. Nazif Unal has "bought" a Durak (dolmuß stop} from Belediye and works all day in Kayseri.

http://onlineqda.hud.ac.uk/Intro\_QDA/Examples\_of\_Qualitative\_Data.php











### What is data?

- "Data are representations of observations, objects, or other entities used as evidence of phenomena for the purposes of research or scholarship." [C. L. Borgman]
- Data can be digital but can also be physical (e.g. sculptures)
- Semantics are important (e.g. temperature to engineer and biologist)
- Grey Data: surveys, student records—think about privacy









## Sharing Data

- Required/encouraged by universities, funding agencies, publishers
- used to support the arguments." [C. L. Borgman]
- Questions:
  - How is data maintained? Who is responsible?
  - What is the process for curating data?
  - How long should data be kept?
  - How should data collection and curation be acknowledged?

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

# "Publications are arguments made by authors, and data are the evidence





### Data Curation Lifecycle

### The DCC Curation Lifecycle Model













## Sequential Actions in Data Curation

- Create or Receive: Create/receive data and make sure metadata exists
- preservation
- Preservation Action: Data cleaning, validation (ensure that data remains) authentic, reliable and usable)
- Store: Store the data in a secure manner adhering to relevant standards Access, Use and Reuse: Make sure is accessible to users and reusers

• Conceptualize: Plan creation of data—capture method and storage options. Appraise and Select: Evaluate data and select for long-term curation and

Ingest: Transfer data to an archive, repository, data centre or other custodian

Transform: Create new data from the original (migrate formats, subsets, etc.)





6

## FAIR Principles

- computers
- Accessible: Users need to know how data can be accessed, possibly including authentication and authorization
- Interoperable: Can be integrated with other data, and can interoperate with applications or workflows for analysis, storage, and processing
- Reusable: Optimize the reuse of data. Metadata and data should be welldescribed so they can be replicated and/or combined in different settings

### • Findable: Metadata and data should be easy to find for both humans and









### Findable: DataCite Workflow















## Accessible: DOI to Landing Page with Metadata



### Document citing the data

D. Koop, CSCI 680/490, Spring 2021

Repository housing the data

Data store











### Interoperable: Standard vocabularies

View as Table	View as Grid				« 1	2	3	4 5	6	7
ort by										
Name		\$	Registry	Name	Abbreviati	ion		Туре	Subjec	t
ecommended Records			க்	ABA Adult Mouse Brain	ABA			Standard	🕜 Neu	iroscieno
Recomm	ended		E	Access to Biological	ABCD			Standard	<ul> <li>Biod</li> </ul>	diversity
Associated Publication?				Collection Data					Life Science	
No Publication	Has Publication									
Claimed?		_								
No Maintainer	Has Maintainer									
Uncertain Deprecated	n development Re	eady		Access to Biological Collection Databases Extended for Geosciences	ABCDEFG			Standard	<ul><li>Eart</li><li>Pale</li></ul>	th Scien eontolog
Standard Type				Access to Biological	ABCDDNA			Standard	Biod	diversity
Terminology Artifact		771		Collection Data						
Model/Format		405								
Reporting Guideline		163								
Metric		30		.ACE format	.ACE forma	at		Standard		Science
Identifier Schema		15								
	Shov	v More	கீ	AdaLab-meta ontology	ADALAB-N	IETA		Standard	None	
Domains			Å	AdaLab ontology	ADALAB			Standard	None	
Report		141		Adverse Drug	EU-ADR M	IL		Standard	None	
		134		Language						

						nc5	~	/ \C/ V	ancea	J										
				Sho	owingr	records	1 - 50	of <b>138</b> 4	ŀ.											
8 9	10	) 11	12	13	14	15	16	17	18 1	.9 20	21	22	23	24	25	26	27	28	»	
		Domain				Taxonoi	my	Relate	d Database		Related	Standard	d I	Related F	Policy	ln Co	llection/	/Recom	mendatio	n Status
		<ul><li>Brain</li><li>Brain In</li></ul>	Gene I	Expression		🖌 Mus	musculus	Neuro	lorpho.Org		None			None		No	ne			R
Biology		None				All		GBIF ALA IP Reposi GBIF S Reposi Canade SiB Co Colomb Plus 1	F - GBIF Aus ory pain IPT - G ory ensys IPT - C ensys Repos ombia IPT - oia Repositor more	stralia BIF Spain GBIF sitory GBIF ry	ABCDDN	IA G		None			DWG Biodivi	versity Inform	nation Standards	
<ul><li>Geology</li><li>Soil Scier</li></ul>	nce	None				II AII		GeoCA	Se Data Por	rtal	XML ABCD		I	None		No	ne			R
Biology		<ul> <li>DNA S</li> <li>Experiit</li> <li>Sequel</li> <li>Deoxyr</li> <li>Polyme</li> <li>Plus 1</li> </ul>	equence Data ment Metadat nce ribonucleic Ac arase Chain R more	a id Reaction		I All		GenBa	nk		MOD-CC ABCD	)	1	None			DWG Biodive	versity Inform	nation Standards	
		DNA S     Deoxyr	equence Data ribonucleic Ac	a 🛷 Co id 🗣 G	enome	🖌 All		None			None		I	None		No	ne			R
		None				🛷 All		None			None		I	None		No	ne			R
		None				🗣 All		None			None		I	None		No	ne			R
		<ul><li>Advers</li><li>Electro</li></ul>	e Reaction	cord			o sapiens	None			XML		I	None		No	ne		[	fairs





### Reusable: Licensing

- Citation of a dataset is expected as a scholarly norm, not by law
- CC0:
  - "I hereby waive all copyright and related or neighboring rights together with all associated claims and causes of action with respect to this work to the extent possible under the law"
- CC BY: license, not a waiver as CC0
  - "You must give appropriate credit, provide a link to the license, and indicate if changes were made."
- Data Use Agreements (DUA): Used when data are restricted due to proprietary or privacy concerns.





### Reusable: Data Citation & Metrics



D. Koop, CSCI 680/490, Spring 2021





Northern Illinois University



12

### <u>Assignment 4</u>

- World Education Data
- Collected/collated by UNESCO, World Bank, and OECD
- Transform World Bank Data
- Impute missing year data
- Integrate teacher and student numbers
- Fuse three datasets





## Studying Data Availability

- Who mandates data sharing, and what is the impact?
  - Government
  - Funding agencies
  - Institutions
  - Journals
- How does the age of a publication/data item affect availability?
  - If not curated, how to locate?
  - What factors influence this?





### Data Availability by Journal Policy









## Data Availability by Year

Year	No Working E-Mail	No Response to E-Mail	Response Did Not Give Status of Data	Data Lost	Data Exist, Unwilling to Share	Data Received	Data Extant (Unwilling to Share + Received)	Number of Papers
1991	9 (35%)	9 (35%)	2 (8%)	4 (15%)	1 (4%)	1 (4%)	2 (8%)	26
1993	14 (39%)	11 (31%)	3 (8%)	7 (19%)	0 (0%)	1 (3%)	1 (3%)	36
1995	11 (31%)	9 (26%)	0 (0%)	7 (20%)	2 (6%)	6 (17%)	8 (23%)	35
1997	11 (37%)	9 (30%)	1 (3%)	2 (7%)	3 (10%)	4 (13%)	7 (23%)	30
1999	19 (48%)	13 (32%)	1 (2%)	1 (2%)	0 (0%)	6 (15%)	6 (15%)	40
2001	13 (30%)	15 (35%)	3 (7%)	4 (9%)	0 (0%)	8 (19%)	8 (19%)	43
2003	9 (20%)	20 (43%)	4 (9%)	2 (4%)	0 (0%)	11 (24%)	11 (24%)	46
2005	11 (24%)	14 (31%)	6 (13%)	1 (2%)	0 (0%)	13 (29%)	13 (29%)	45
2007	12 (18%)	31 (47%)	2 (3%)	4 (6%)	1 (2%)	16 (24%)	17 (26%)	66
2009	9 (13%)	34 (49%)	3 (4%)	5 (7%)	6 (9%)	12 (17%)	18 (26%)	69
2011	13 (16%)	29 (36%)	8 (10%)	0 (0%)	7 (9%)	23 (29%)	30 (38%)	80
Totals	131 (25%)	194 (38%)	33 (6%)	37 (7%)	20 (4%)	101 (19%)	121 (23%)	516

Data are displayed as n (%); the percentages are calculated by rows.







## Working Email







### Received Response









### Status of Data







## Data Extant (Shared or Exists)









### Lots of Data is Shared...







### Genome Sequence and Structure Data









### ...but how much isn't shared?

- What isn't shared?
- Who isn't sharing?
- Why not?
- How much does it matter?
- What can be done about it?















## Why Share Data? Increased Citations





Articles with Data Articles with Data Shared (n=41)





Northern Illinois University





## What Factors Impact Sharing?













### Factors

- Has journal policy
- Count of R01 & other NIH grants
- Authors prev GEOAE sharing & OA & microarray creation
  - NO K funding or P funding
    - Journal impact
  - Journal policy consequences & long halflife
    - Institution high citations & collaboration
      - NOT animals or mice
  - Institution is government & NOT higher ed
  - Last author num prev pubs & first year pub
    - Large NIH grant
    - Humans & cancer
  - NO geo reuse + YES high institution output
  - First author num prev pubs & first year pub













## Why not data sharing? (self-reported)

sharing is too much effort want student or jr faculty to publish more they themselves want to publish more

- commercial value of results

D. Koop, CSCI 680/490, Spring 2021

cost industrial sponsor confidentiality



0%

### Nature data availability and data citations

- Policy as of July 2016
- data-citations.pdf

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

### <u>http://www.nature.com/authors/policies/data/data-availability-statements-</u>









# The Evolution of Data Citation: From Principles to Implementation

M. Altman and M. Crosas





## Data Sharing Policies

- Science:
  - "all data necessary to understand, assess, and extend the conclusions of the manuscript must be available to any reader of Science"
  - "citations to unpublished data and personal communications cannot be used to support claims in a published paper"
- Often this is only used as reason to retract work when issues arise
- Need:
  - Recognition of data authorship
  - Robust citation practices and infrastructure







## Chronology of Data Citation

**Exemplar Systems** 



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

### **Core Principles**

### Key Work

<ul> <li>Facilitate description</li> <li>&amp; information retrieval</li> <li>Describe data in archives</li> <li>Describe as works not media</li> <li>Provide author, title, version.</li> </ul>		[Avram 1975] [Dodd 1979] [ISBD 1990] [ISO 1997]
--	--	--

- Facilitate access
& persistence
- Cite research data in all
publications that use it.
<ul> <li>Provide actionable URI's</li> </ul>
- Provide persistent identifiers
- Use persistent institutions

[Altman, et al. 2001] [Ryssevik & Musgrave 2001]

- Facilitate verification & reproducibility - Provide bit- or semantic- fixity - Provide granularity

[Brase 2004] [Buneman 2006] [Altman & King 2007]

- Facilitate integration - Include data citations in standard locations in text - Index data citations in existing catalogs - Integrate data citation with

[Uhlir (ed.) 2012] [CODATA 2013] [Data Synthesis Group 2014]







## Phases of Data Citation (1977-2009)

- citation? (Libraries)
- 2. Support data access and persistence: if citations to data in publications, need methods to discover information about data
- 3. Support verification and **reproducibility**: allow verification of claims based on the data (wider integration into publishing

1. Support description and information retrieval: what should be included in a









## Joint Declaration of Data Citation Principles

- 1. Importance. Data should be considered legitimate, citable products of research. Data citations should be accorded the same importance in the scholarly record as citations of other research objects, such as publications.
- 2. Credit and Attribution. Data citations should facilitate giving scholarly credit and normative and legal attribution to all contributors to the data, recognizing that a single style or mechanism of attribution may not be applicable to all data.
- 3. Evidence. In scholarly literature, whenever and wherever a claim relies upon data, the corresponding data should be cited.
- 4. Unique Identification. A data citation should include a persistent method for identification that is machine actionable, globally unique, and widely used by a community.









## Joint Declaration of Data Citation Principles

- referenced data.
- 6. **Persistence**. Unique identifiers, and metadata describing the data they describe.

5. Access. Data citations should facilitate access to the data themselves and to such associated metadata, documentation, code, and other materials, as are necessary for both humans and machines to make informed use of the

the data, and its disposition, should persist -- even beyond the lifespan of









## Joint Declaration of Data Citation Principles

- access to, and verification of the specific data that support a claim. originally cited.
- citation practices across communities.

7. Specificity and Verifiability. Data citations should facilitate identification of, Citations or citation metadata should include information about provenance and fixity sufficient to facilitate verifying that the specific timeslice, version and/or granular portion of data retrieved subsequently is the same as was

8. Interoperability and flexibility. Data citation methods should be sufficiently flexible to accommodate the variant practices among communities, but should not differ so much that they compromise interoperability of data













### Generic Data Citation

- Archive, version or subset
- Authors, repository  $\rightarrow$  Principle 2
- Global Persistent Identifier: Principle 4 and 6

### Author(s), Year, Dataset Title, Global Persistent Identifier, Data Repository or

• Year and title  $\rightarrow$  not related to principle but consistent with other citations







## More Information

- Provide via the web
  - Metadata
  - Fixity and provenance information
- Community Indices:
  - CrossRef
  - DataCite

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

### • Structured Identifiers (ORCID, ISNI) preferred over unstructured metadata





### Example Repositories with Citations

- Dryad, Dataverse, Figshare
- Dataverse:
  - Draft citation **automatically** generated
  - Includes versioning information







## Remaining Challenges

- Provenance: chain of ownership
- Identity: equivalence and derivation relationships
  - Equivalence: if not bitwise equal, can data still be interchangeable?
  - Versioning: if data is updated, how to find updated version?
  - Granularity: How to describe subsets of data (deep citation)
- Attribution: ensure that the correct people and institutions receive credit







### DataCite

### www.datacite.org





### Why Data Citation is a Computational Problem

P. Buneman, S. Davidson, and J. Frew





## Computational Data Citation

- Given a database D and a query Q, generate an appropriate citation. Automatic Citation requires the answers to two questions:
- - Does the citation depend on both Q and D or just on the data Q(D)extracted by Q from D?
  - If we have appropriate citations for some queries, can we use them to construct citations for other queries?
- If the data is an image or numbers, cannot expect the citation to live in that data
- If the query returns an empty dataset, we still may wish to cite that People know how to cite certain parts of a dataset but not all...

[Buneman et al., 2016]







## Computational Data Citation (GtoPdb)



Search Database
IACOLOGY
d rat (Rn).
all sections Collapse all sections
page
itation:
r, Daniel J. Drucker, Dominique Bataille, Susan Chan, Philippe Delagrange, Burkhard Göke, Kelly E. Mayo, Bernard Thorens, Rebecca Hills. or family. Accessed on 16/06/2016. IUPHAR/BPS Guide to PHARMACOLOGY, http://www.guidetopharmacology.org splayForward?familyId=29.
PHARMACOLOGY citation:
Davenport AP, Kelly E, Marrion N, Peters JA, Benson HE, Faccenda E, Pawson AJ, Sharman JL, Southan C, Davies JA and CGTP 115) <b>The Concise Guide to PHARMACOLOGY 2015/16: G protein-coupled receptors.</b> Br J Pharmacol. <b>172</b> : 5744-5869.
ily introduction, please use the following:
er, Daniel J. Drucker, Dominique Bataille, Susan Chan, Philippe Delagrange, Burkhard Göke, Kelly E. Mayo, Bernard Thorens, Rebecca Hills. For family, introduction. Last modified on 10/08/2015. Accessed on 16/06/2016. IUPHAR/BPS Guide to PHARMACOLOGY, etopharmacology.org/GRAC/FamilyIntroductionForward?familyId=29.







## Computational Data Citation (MODIS)









## Views and Citable Units

- Views describe "areas of responsibility" for parts of a database
- Use views to create "citable units"
- Determine which view V answers a particular query Q and generate a citation for the view
- What happens if two different views can answer the same query?











### Citable Views and Partial Citations











### Hierarchies of Views

- In GtoPdb, three classes of views
- Family view:
  - /Root/Family[FamilyName=\$\$f]
- Introduction view:
  - /Root/Family[FamilyName=\$\$f]/ Introduction
- Target view:
  - /Root/Family[FamilyName=\$\$f]/ Target[TargetName=\$\$t]





## Citation Rule and Partial Result (GtoPdb)

• Rule:

←

- { Title: "IUPHAR/BPS Guide to Pharmacology", Version: \$v, Family: \$\$f, Contributors: \$a, URI: "www.iuphar.org" }
  - /Root[VersionNumber: \$v]/Family[FamilyName: \$\$f]/Introduction[Contributorlist: \$a]
- Citation:
  - { Title: "IUPHAR/BPS Guide to Pharmacology", Version: 26, Family: "Calcitonin", Contributors: ["Debbie Hay", "David R. Poyner"], URI: "www.iuphar.org" }







## Citation Rule and Sample Result (MODIS)

- { author: m\_auth(\$p,\$\$v), m\_year:(\$p,\$\$v), title: m\_title(\$p), version: \$v, bounding-box : [\$\$minlong, \$\$minlat, \$\$maxlong, \$\$maxlat], interval: [\$\$mint, \$\$maxt], organization: m\_org(\$p), url: m\_url(\$p), accessed: DATE(), doi =  $m_{doi}(p, v)$ 
  - /root/product[ProdName=\$p]/version[vnum=\$\$v] /file[Lat  $\geq$  \$\$minlat and Lat  $\leq$  \$\$maxlat and Lon  $\geq$  \$\$minlon and Lon  $\leq$  \$\$maxlon and Time  $\geq$  \$\$mint and Time  $\leq$  \$\$maxt]
- { author: "E. Vermote", title: "MOD09A1 ... SIN Grid", version: 6, bounding-box: [-125, 32, -119, 42], interval: [2008-01-25, 2008-01-25], accessed: "2015-09-01", doi: "10.5067/MODIS/MOD09A1.006" }

-

organization: "NASA EOSDIS ... South Dakota", URL: "https://lpdaac.usgs.gov",







