

ACCELERATING SYNTHESIS SCIENCE THROUGH REPRODUCIBLE SCIENCE PRACTICES

Matthew B. Jones

*National Center for Ecological Analysis and Synthesis
University of California Santa Barbara*



@metamattj

jones@nceas.ucsb.edu

<https://orcid.org/0000-0003-0077-4738>



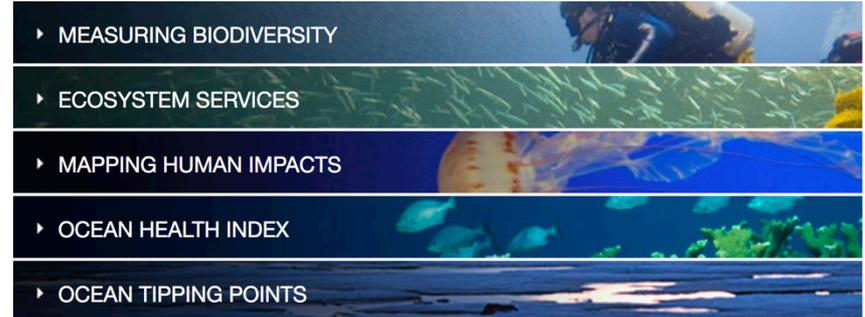
Marine Systems



Threats and Population Declines



Understanding Ocean Health



Climate and Ecosystems





Provenance



Citation



Synthesis

Reproducible Science



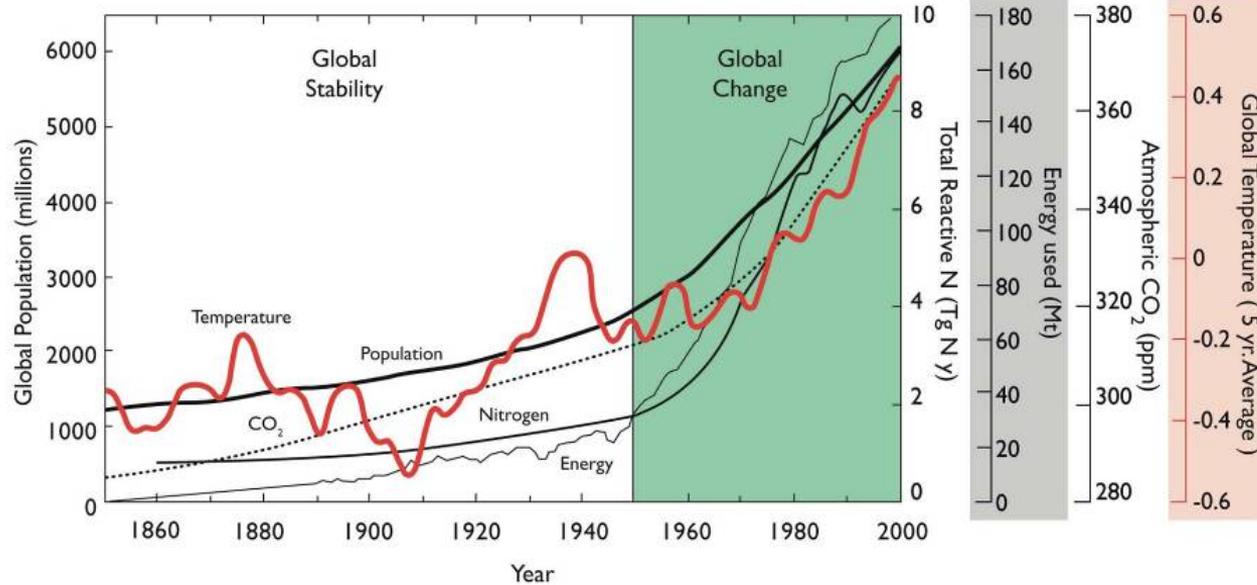
Climate Change
Fisheries
Sustainability
Subsistence



Science
Governance
Regulation
Policy



Trust in Science



What **data**?
What **methods**?
What **parameter settings**?

Can we **trust** these data and methods?

Smith et al. (2009) Ecology doi:10.1890/08-1815.1

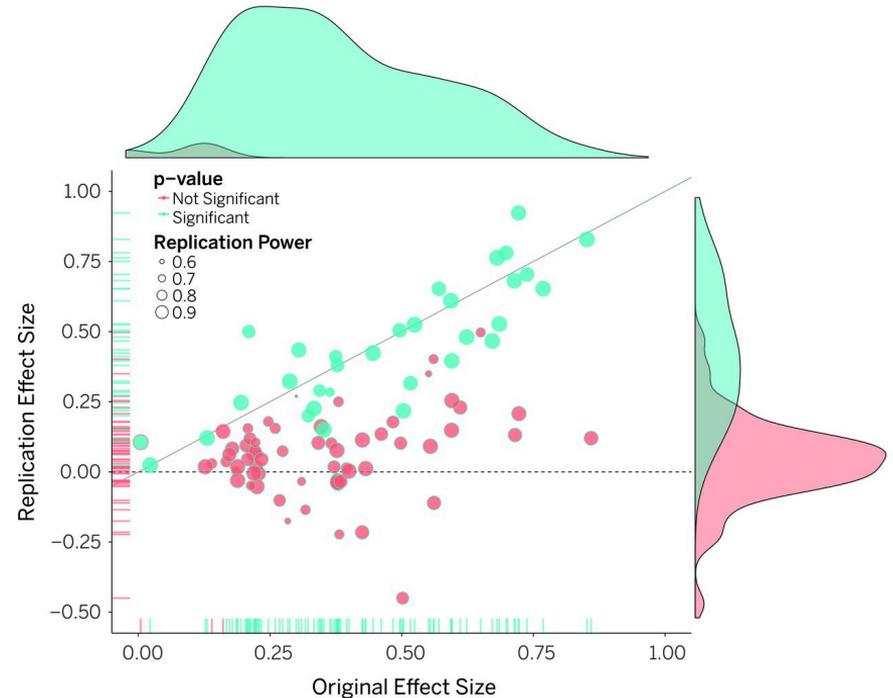
Reproducibility Crisis

“Most research findings are false for most research designs and for most fields”

Ioannidis, 2005

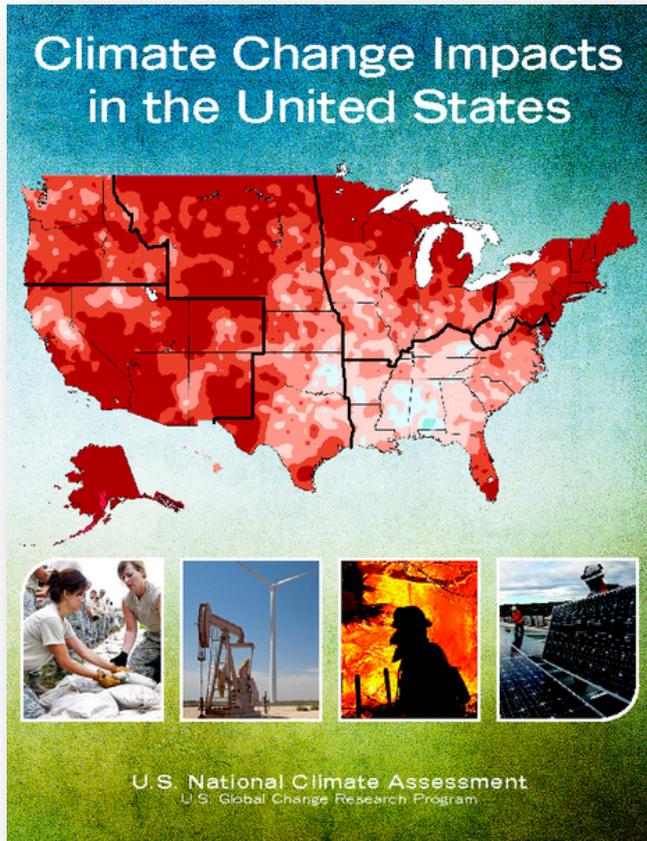
“Most replication effects were smaller than original results”

Open Science Collaboration, 2015



doi:10.1126/science.aac4716

National Climate Assessment



“This report is the result of a **three-year** analytical effort by a team of **over 300 experts**, overseen by a broadly constituted Federal Advisory Committee of **60 members**. It was developed from information and analyses gathered in over 70 workshops and listening sessions held across the country.”

Computational Reproducibility

Facilitate transparency by **capturing** and **communicating** scientific workflows

Increase **trust in science**



Stand on the shoulders of giants
(build on work that came before)

Give credit for that **secondary** usage enabling **easy attribution**

Practical Reproducibility

Preserve the data

Preserve the software workflow

Document what you did

Describe how to interpret it all



✕ Clear all filters

Search ?

Search phrase



My Search

sasap ✕

Filter by:

▶ Data attribute

▶ Data files

▶ Creator

▶ Year

▶ Identifier

▶ Taxon

▶ Location

DATASETS 1 TO 25 OF 44

1 2 Next

Sort by Most recent ▼

Jeanette Clark and Rich Brenner. 2017. **Sockeye salmon brood tables, northeastern Pacific, 1922-2016**. Knowledge Network for Biocomplexity. urn:uuid:c11dff42-b988-437a-afee-58fc62dcd1dc.

5

Commercial Fisheries Entry Commission. 2018. **Commercial Fisheries Entry Commission Basic Information Table, 1975-2016**. Knowledge Network for Biocomplexity. urn:uuid:8f351735-baf9-451a-b821-c1117ebf5a5e

12

Andrew Munro and Eric Volk. 2018. **Summary of Pacific Salmon Escapement Goals in Alaska with a Review of Escapements from 2001 to 2009**. Knowledge Network for Biocomplexity. urn:uuid:d62539fd-3025-48d0-a1c3-5a903de1f269.

10

Alaska Department of Labor and Workforce Development, Research and Analysis Section. 2018. **Alaskan fishing industry employee counts by month, grouped by region and fish species from 2000-2016**. Knowledge Network for Biocomplexity. urn:uuid:32958097-0ad3-428a-aba9-c37e804be0ef.

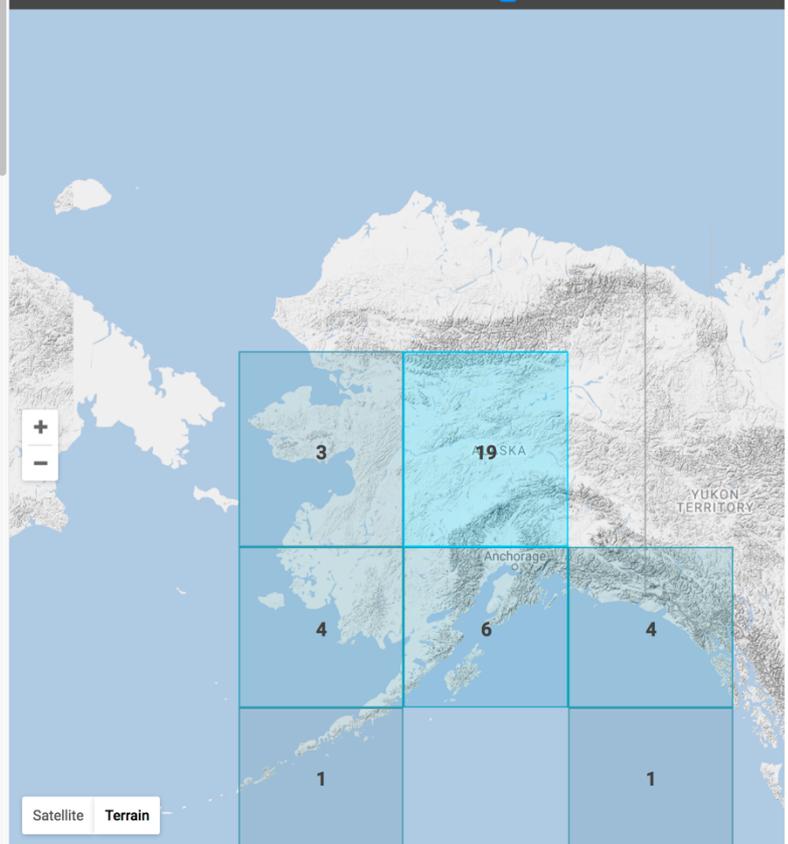
9

Alaska Department of Labor and Workforce Development Research & Analysis Section. 2018. **Alaskan fishing industry employee counts by month, subsetted by region and fish species**. Knowledge Network for Biocomplexity. urn:uuid:4bbc9577-e81f-40f4-b4ca-9c740092baba.

5

Commercial Fisheries Entry Commission. 2018. **Commercial Fisheries Entry Commission Permit Earnings, 1975-2016**. Knowledge Network for Biocomplexity.

Hide Map >

 Limit my search to the map area+
-

Satellite Terrain

Google

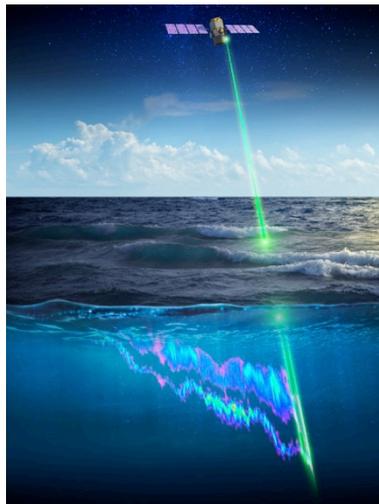
Map data ©2018 Google, INEGI, SK telecom, ZENRIN 500 km Terms of Use



Global
Data Coverage



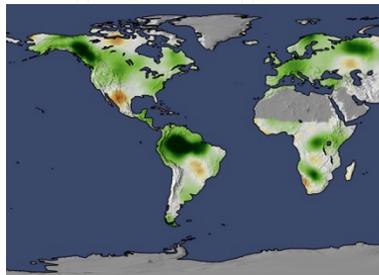
800K
Data Packages



40
Repositories



143K
Contributors





Reproducible
Science



Provenance



Citation



Synthesis

Computational Provenance

Origin, processing history of data

- Input data
- Workflow/scripts
- Output data
- Figures
- Understand methods, dataflow, and dependencies

Texas Summer 2011: Record Heat and Drought
Cooperative Institute for Climate and Satellites - NC
Laura Stevens

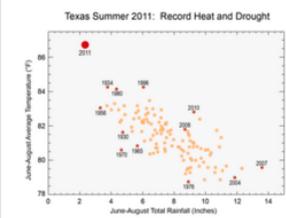
The **time range** for this image is January 01, **1895** (00:00 AM) to December 31, **2012** (00:00 AM).
This image was created on July 03, 2013.

The **spatial range** for this image is 25.83° to 36.50° latitude, and -106.65° to -93.52° longitude.

Attributes : Temperature, precipitation, observed, Texas.

This image **was derived from** dataset `nca3-cddv2-r` **using** the activity `02c53cf7-nca3-cddv2-r1-process`.

This image is part of this figure :



data and "code" / method linked

alt formats

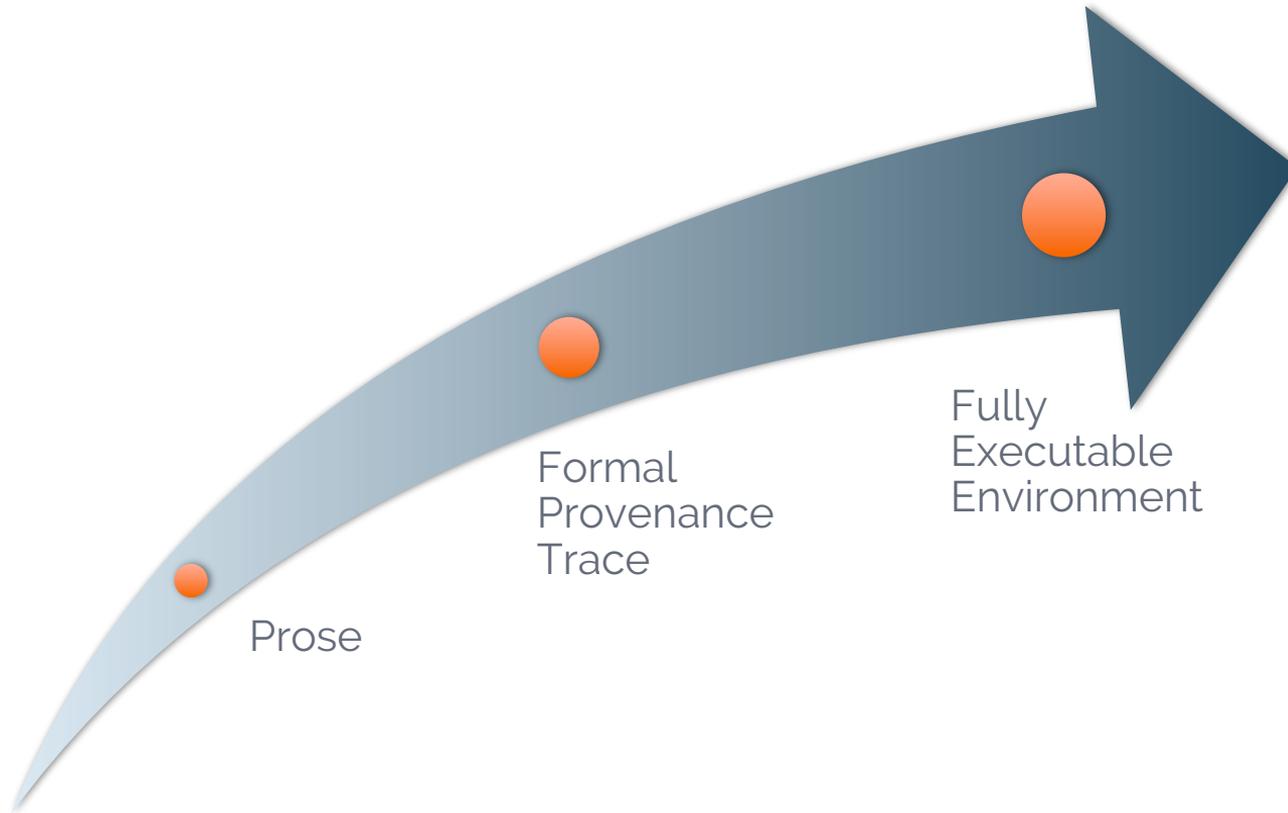
You are viewing `/image/02c53cf7-75f8-4243-a925-f59a0025f04e` in [HTML](#)

Alternatives : [JSON](#) [YAML](#) [Turtle](#) [N-Triples](#) [JSON Triples](#) [RDF/XML](#) [RDF+JSON](#) [Graphviz](#) [SVG](#)

GlobalChange.gov

Provenance

Origin and processing history of artifacts



Provenance in DataONE

Phase II Goal: Facilitate reproducible science

- Track **data derivation** history
- Track data **inputs** and **outputs** of analyses
- Track analysis and model **executions**
- Preserve and document software **workflows**
- Link all of these to **publications**



ProvONE

Extended PROV model for workflow provenance.

Prov Index

DataONE support for indexing, searching, and displaying provenance.

R and Matlab

Libraries in R, MATLAB, Java for generating and manipulating provenance records.

Web Provenance

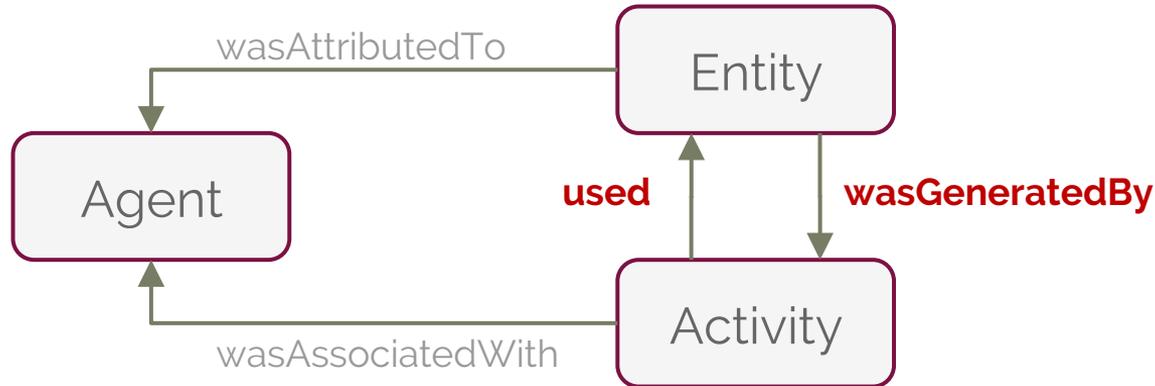
Web-based user interface for displaying and editing provenance.

Modeling Provenance



W3C PROV

See w3.org/TR/prov-o/

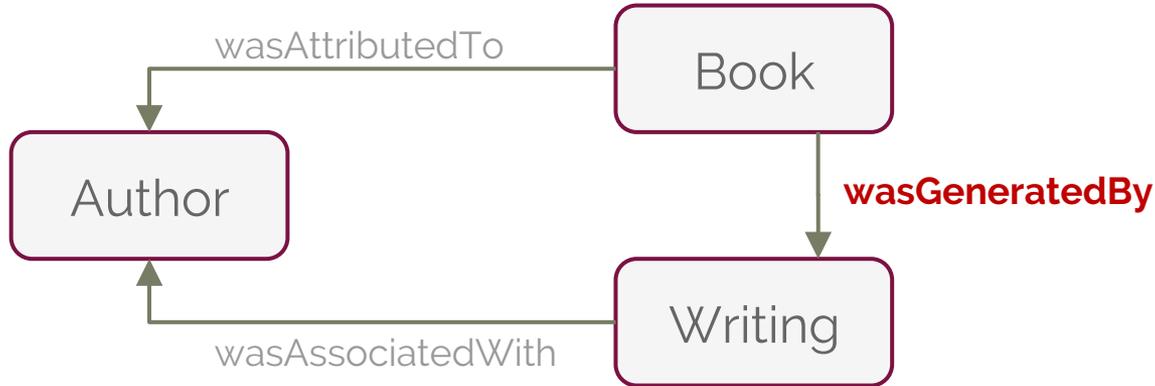


Modeling Provenance

ProvONE

W3C PROV

See w3.org/TR/prov-o/



Provenance for Science Workflows



ProvONE – an extension of W3C PROV

See purl.dataone.org/provone-v1-dev

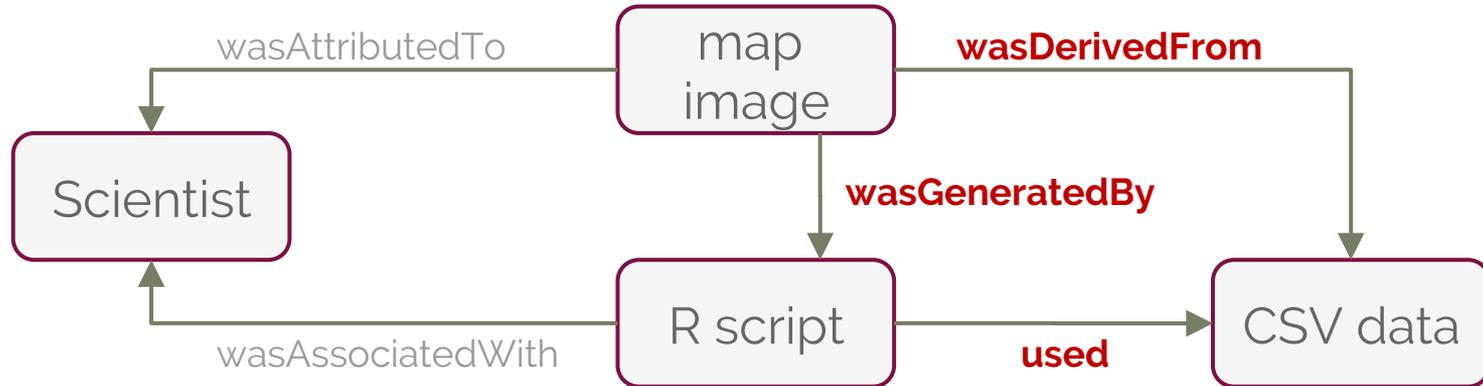


Provenance for Science Workflows

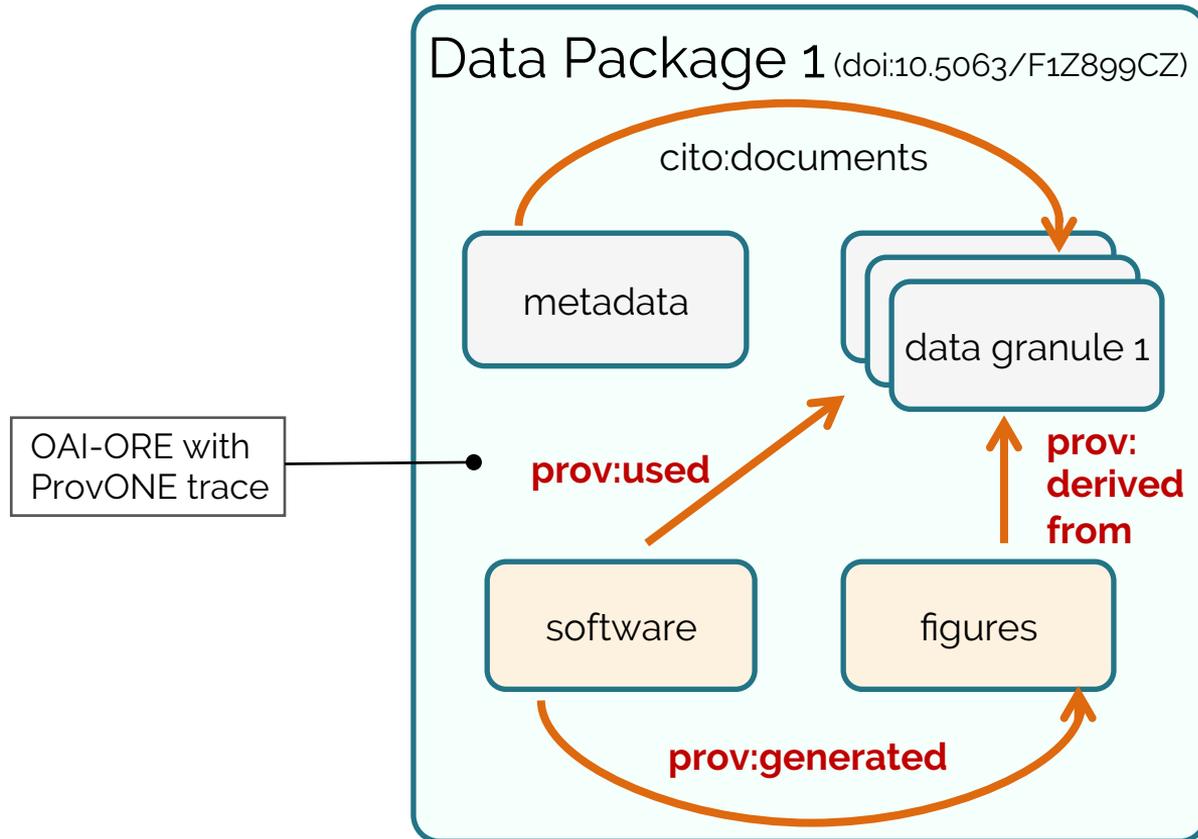


ProvONE – an extension of W3C PROV

See purl.dataone.org/provone-v1-dev

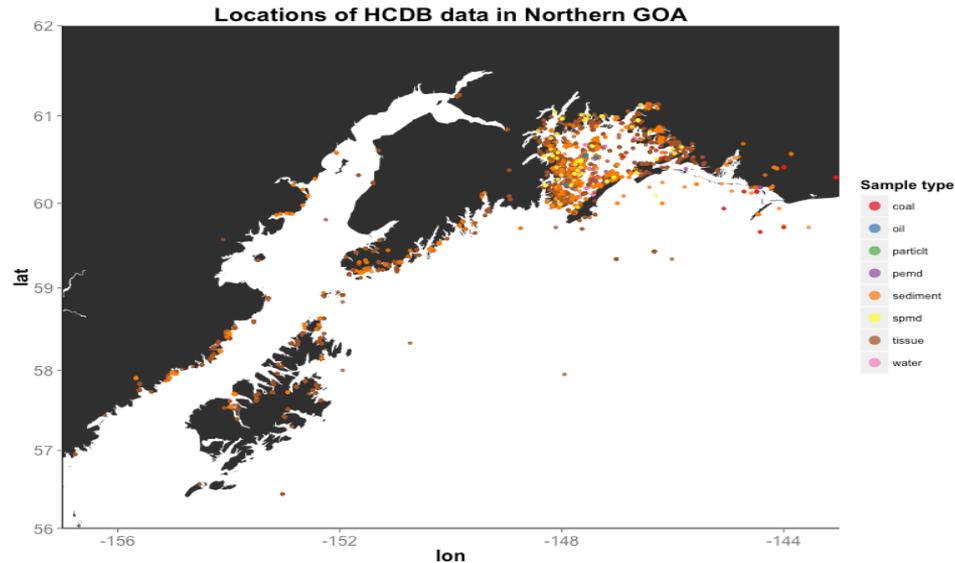


Data Package with Provenance



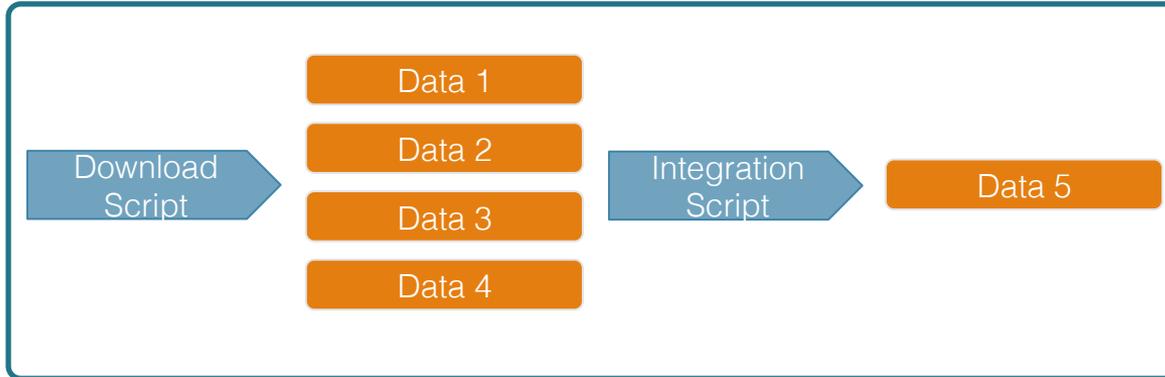
Hydrocarbon Data Example

Mark Carls. 2017. Analysis of hydrocarbons following the Exxon Valdez oil spill, Gulf of Alaska, 1989 - 2014. Arctic Data Center.

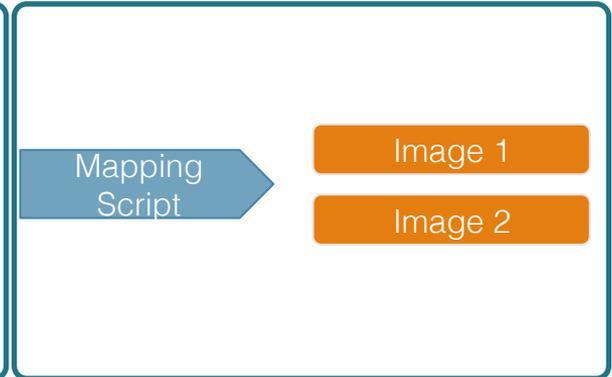


Publishing Data Workflows

Dataset C



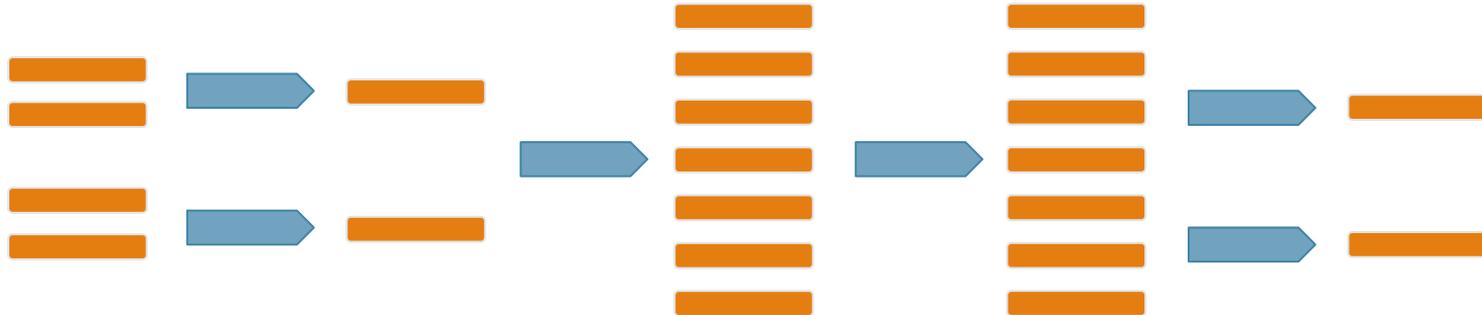
Dataset D



Hydrocarbon Data Example

Complex Workflows

Simplified view of complex workflows



Provenance Display

DataONE Search

About News Participate Resources Education Data

DATAONE SEARCH: Search Summary Jump to: DOI or ID Go

Sign in or Sign up

< Back to search | Search / Metadata

Mark Carls. 2017. Analysis of hydrocarbons following the Exxon Valdez oil spill, Gulf of Alaska, 1989 - 2014. Gulf of Alaska Data Portal. urn:uuid:3249ada0-afe3-4dd6-875e-0f7928a4c171.



[Copy Citation](#)

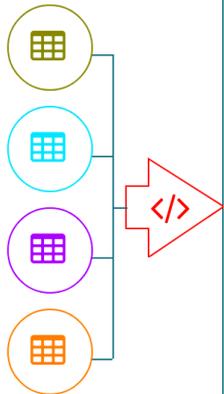
Files in this dataset Package: urn:uuid:1d23e155-3ef5-47c6-9612-027c80855e8d

Name	File type	Size	Download all
Metadata: metadata.xml	EML v2.1.1	140 KB	112 views
Total_Aromatic_Alkanes_PWS.csv	More info text/csv	3 MB	3 downloads
CollectionMethods.csv	More info text/csv	793 B	2 downloads
Non-EVOS_SINs.csv	More info text/csv	3 KB	

[Show 8 more items in this data set](#)

Data Table, Image, and Other Data Details

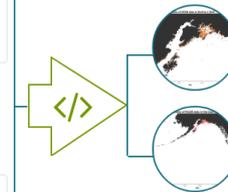
4 sources



Data Table

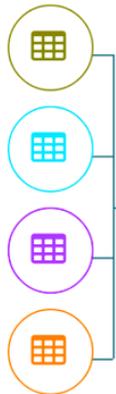
Entity Name	Total_Aromatic_Alkanes_PWS.csv										
	Download										
Description	Combined dataset from PAH, Alkane and Sample tables documenting samples collected after the Exxon Valdez oil spill in Prince William Sound, AK										
Object Name	Total_Aromatic_Alkanes_PWS.csv										
Online Distribution Info	https://cn.dataone.org/cn/v2/resolve/urn:uuid:44108e76-405d-4d58-b1b3-fb4b55e3fff9										
Size	2801033 byte										
Text Format	<table><tr><td>Number of Header Lines</td><td>1</td></tr><tr><td>Record Delimiter</td><td>#x0A</td></tr><tr><td>Attribute Orientation</td><td>column</td></tr><tr><td colspan="2">Simple Text</td></tr><tr><td>Field Delimiter</td><td>,</td></tr></table>	Number of Header Lines	1	Record Delimiter	#x0A	Attribute Orientation	column	Simple Text		Field Delimiter	,
Number of Header Lines	1										
Record Delimiter	#x0A										
Attribute Orientation	column										
Simple Text											
Field Delimiter	,										
Number Of Records	12142										

2 derivations



Data Table, Image, and Other Data Details

4 sources



Source Program

Total_PAH_and_Alkanes_GoA_Hydrocarbons_Clean.R

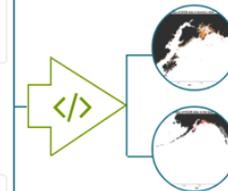
Citation

[View »](#)

This program generated the data you are currently viewing, **Total_Aromatic_Alkanes_PWS.csv**.

This program used **PAH.csv**, **Sample.csv**, **Non-EVOS_SINs.csv** and **(and 1 more)**.

2 derivations



Text Format

Number of Header Lines	1
Record Delimiter	#x0A
Attribute Orientation	column
Simple Text	
Field Delimiter	,

Number Of Records

12142

Web Provenance Editor

Deployed by Arctic Data Center

The screenshot displays the NSF Arctic Data Center Web Provenance Editor interface. At the top, the header includes the NSF Arctic Data Center logo, navigation links for Data, Support, and About, a green Submit Data button, and a user profile for Christopher Jones. The main content area is titled "Data Table, Image, and Other Data Details" and features a "Data Table" section. This section displays the following information:

- Entity Name:** Total_Aromatic_Alkanes_PWS.csv
- Description:** Combined dataset from PAH, Alkane and Sample tables documenting samples collected after the Exxon Valdez oil spill in Prince William Sound, AK
- Object Name:** Total_Aromatic_Alkanes_PWS.csv
- Online Distribution Info:** <https://cn-stage.test.dataone.org/cn/v2/resolve/urn:uuid:df984766-dd89-4e57-b97e-350506d7007e>
- Size:** 2801033 byte
- Text Format:**

Number of Header Lines	1
Record Delimiter	#x0A

Additional interface elements include "0 sources" and "0 derivations" counts, each with a circular "Add" button and a right-pointing arrow. A "Download" button with a cloud icon is located below the Entity Name field.



Add source data to Total_Aromatic_Alkanes_PWS.csv

Choose files in this dataset:

CollectionMethods.csv
hcdbSamplesGOA.png
hcdbSampleLocs.png
PAH.csv
Alkane.csv
Non-EVOS_SINs.csv
Sample.csv

Done

Online Distribution Info

<https://cn-stage.test.dataone.org/cn/v2/resolve/urn:uuid:df984766-dd89-4e57-b97e-350506d7007e>

Size 2801033 byte

Text Format

Number of Header Lines	1
Record Delimiter	#x0A
Attribute Orientation	column
Simple Text	
Field Delimiter	,

Number Of Records 12142



Data Table, Image, and Other Data Details

4 sources



Data Table

Entity Name **Total_Aromatic_Alkanes_PWS.csv**
[Download](#)

 Description
 Combined dataset from PAH, Alkane and Sample tables documenting samples collected after the Exxon Valdez oil spill in Prince William Sound, AK
Object Name **Total_Aromatic_Alkanes_PWS.csv**
 Online Distribution Info
<https://cn-stage.test.dataone.org/cn/v2/resolve/urn:uuid:df984766-dd89-4e57-b97e-350506d7007e>
Size **2801033 byte**

Text Format	
Number of Header Lines	1
Record Delimiter	#x0A
Attribute Orientation	column
Simple Text	
Field Delimiter	,

0 derivations


[Save](#)

Provenance Editing



Matlab DataONE Toolbox



Recordr R Library



YesWorkflow Tool

MetacatUI Web Provenance Editor

Data Table, Image, and Other Data Details

0 sources

0 derivations

Data Table	
Entity Name	Total_Aromatic_Alkanes_PWS.csv
	Download
Description	Combined dataset from PAH, Alkane and Sample tables documenting samples collected after the Exxon Valdez oil spill in Prince William Sound, AK



Reproducible
Science



Provenance



Citation



Synthesis

Credit where credit is due

Indexing and exposing data citations in international data repository networks



ALFRED P. SLOAN
FOUNDATION



University of California
CDL
California Digital Library



Data  **ONE**

Force11 Data Citation Principles

1. Importance of data citation
2. **Credit and Attribution**
3. **Evidence**
4. Unique Identification
5. Access
6. **Persistence**
7. **Specificity** and Verifiability
8. Interoperability and Flexibility

Transitive Credit

When a user cites a pub, we know:

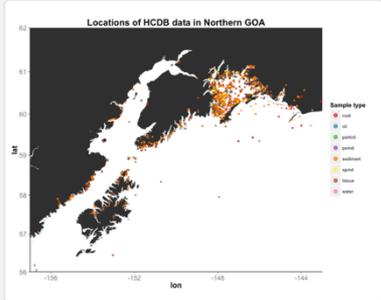
- **Which data** produced it
- **What software** produced it
- What was **derived** from it
- **Who to credit** down the attribution stack

See: Katz & Smith. 2014. **Implementing Transitive Credit with JSON-LD**. arXiv:1407.51

Derived image

Map of sampling locations in the Northern Gulf of Alaska

Citation
Mark Carls. 2015. **Hydrocarbon database, Gulf of Alaska**. MN Demo 2. urn:uuid:bf71c38b-22b2-469e-8983-734ec0ab19cb.



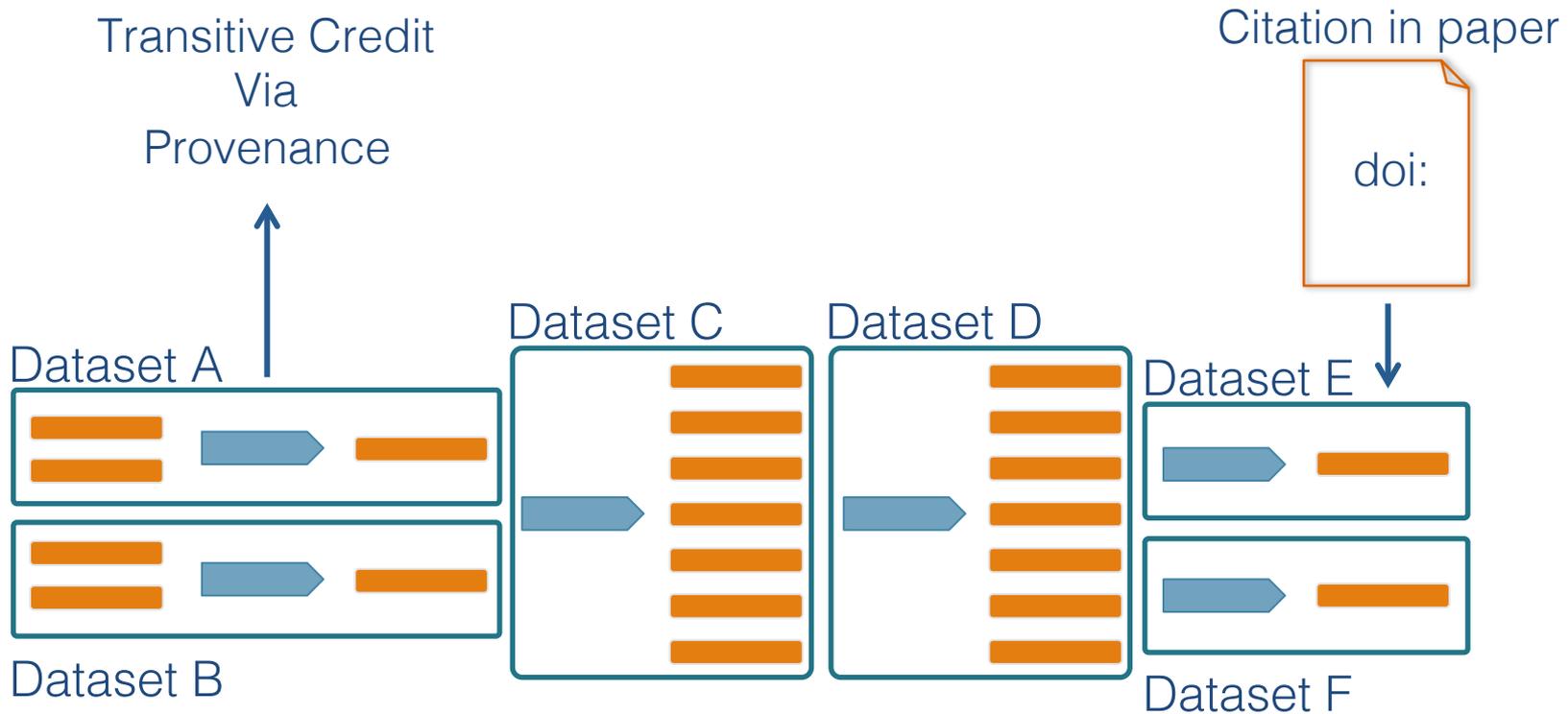
Locations of HCDB data in Northern GOA

View »

This image was generated by the program you are currently viewing, `</> Locations map R script`.

This image was derived from `Total_Aromatic_Alkanes_PWS.csv`.

Citing multi-generational workflows



Evolution of the Living Paper

Scholarly Publications



1 st Gen	Prose		
2 nd Gen	Prose	+ Data	
3 rd Gen	Prose	+ Data	+ Code

Prose + Data + Code + **Provenance**

Prose + Data + Code + **Provenance** +
Execution Environment

DataONE

WHOLETALE



Reproducible
Science



Provenance



Citation



Synthesis



NCEAS

National Center for Ecological Analysis and Synthesis

State of Alaska's Salmon and People

8 SASAP working groups

1: Bio-physical State of Knowledge of Salmon Distribution & Habitat

Leads: Peter Westley and Dan Rinella

2: Sociocultural and Economic Dimensions of Salmon Systems

Leads: Courtney Carothers, Jessica Black, Tobias Schworer

3: Governance and Subsistence

Leads: Steve Langdon, Taylor Brelsford, James Fall

4: Consistency, Causes, and Consequences of Declining Size and Age of Alaskan Salmon

Leads: Eric P. Palkovacs, Peter Westley, Bert Lewis

5: Well-Being and Alaska Salmon Systems

Leads: Rachel Donkersloot, Jessica C. Black, Courtney Carothers

6: Interacting Effects of Ocean Climate and At-Sea Competition on Alaskan Salmon

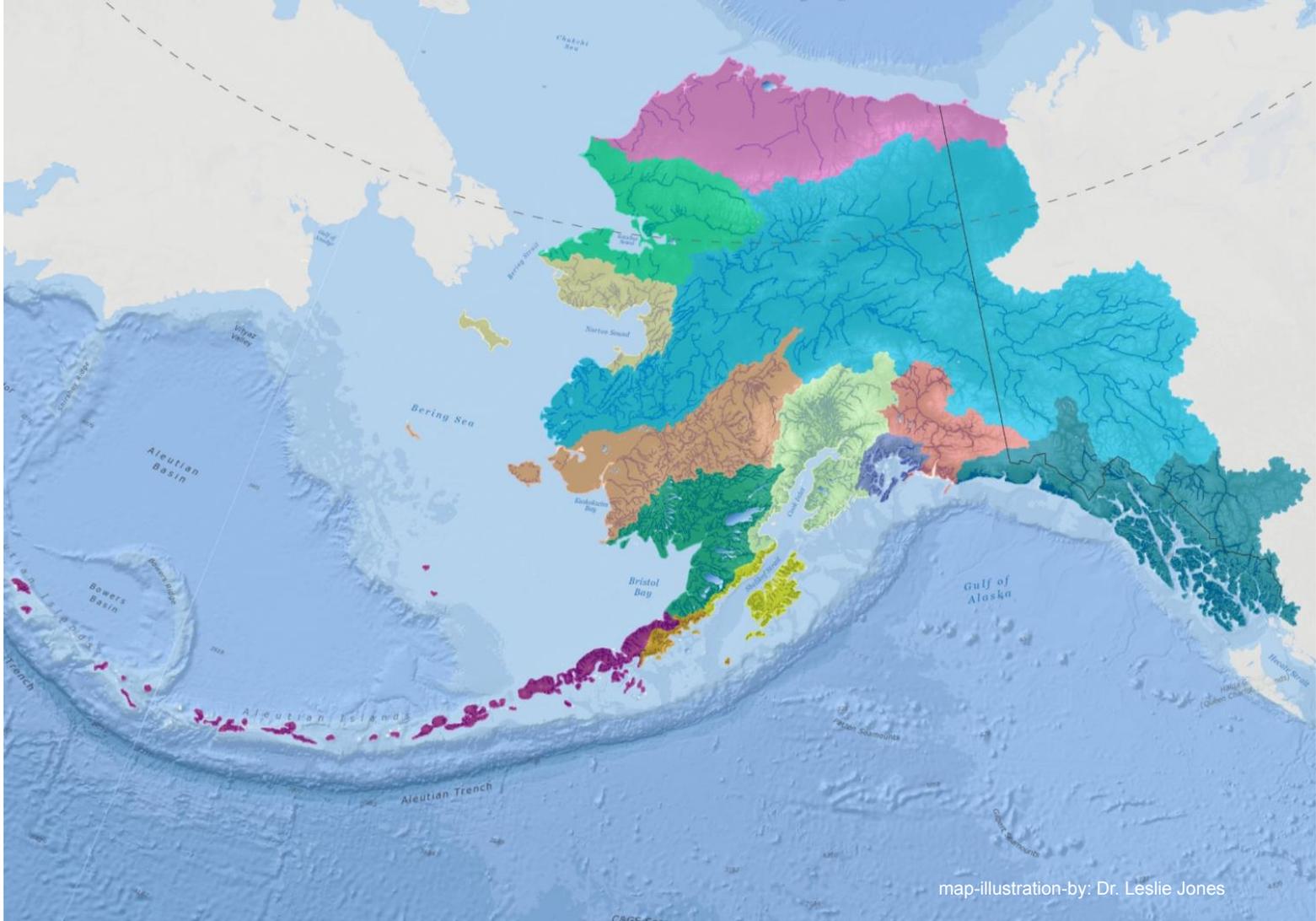
Leads: Peter S. Rand, Robert W. Campbell, Kristen B. Gorman

7: Using Participatory Modeling to Empower Community Engagement in Salmon Science

Leads: Michael L. Jones

8: Kenai Lowlands Salmon Research Synthesis and Design Tools for Integrated Watershed Management

Leads: Coowe Walker, Mark Rains, Ryan King, Charles Simenstad, Dennis Whigham



map-illustration-by: Dr. Leslie Jones

Jeanette Clark and Rich Brenner. 2017. Sockeye salmon brood tables, northeastern Pacific, 1922-2016. Knowledge Network for Biocomplexity. urn:uuid:c11dff42-b988-437a-afee-58fc62dcd1dc.

[Copy Citation](#)[Quality report](#)

Files in this dataset Package: resource_map_urn:uuid:c11dff42-b988-437a-afee-58fc62dcd1dc

Name	File type	Size	Downloads	Download All
Metadata: broodTable_metadata.xml	EML v2.1.1	37 KB	5 views	Download
BroodTables.csv	More info text/csv	449 KB	61 downloads	Download
StockInfo.csv	More info text/csv	19 KB	2 downloads	Download
SourceInfo.csv	More info text/csv	723 B	2 downloads	Download
broodTableProcessing.Rmd	More info application/R	19 KB	3 downloads	Download
broodTableProcessing.html	More info HTML	1 MB	9 downloads	Download

[Show less](#)

30 inputs



view more

Other Entity

Entity Name

Download

Data Object Type:
Other

Physical Structure Description:

Object Name

Source Data

urn:uuid:514f65fa-7f6b-4276-b502-4f46834d309b

Citation

View »

This data prov_hasDerivations
[BroodTables.csv](#).

This data was used by the program you
are currently viewing, `</>`
broodTableProcessing.Rmd.

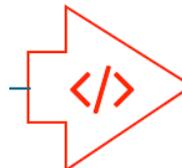
This data was used as an input to
create [BroodTables.csv](#).

1 outputs



Rmarkdown as Provenance

```
01-brood-table-integration.Rmd x
31
32 ## Datasets
33
34 As part of the SASAP project, brood tables for 48 Sockeye salmon were collected.
35 Table 2.1 shows a list of these stocks, along with other regional and location
36 information.
37
38 ```{r, echo = FALSE}
39 stocks <- read.csv('data/original/StockInfo.csv', stringsAsFactors = F)
40
41 ```{r, echo = FALSE}
42 datatable(stocks[, c('Stock.ID', 'Stock', 'Region', 'Sub.Region')], rownames = FALSE,
43 caption = "Stock information")
44
45
46 These stocks range geographically from Washington to Alaska. Although temporal coverage
47 varies by stock, many of the brood tables were updated in 2016, and some have
48 reconstructions dating back to 1922.
49
50 Figure 2.1 indicates the approximate location of the salmon stocks in Table 2.1.
51
52 ```{r, echo = FALSE}
53 salmon = makeIcon('images/salmon_tiny.png',
54 'images/salmon_big.png',
55 26, 14)
56
57 m <- leaflet(stocks) %>%
58   setView(-median(stocks$Lon), median(stocks$Lat), zoom = 4) %>%
59   addTiles() %>%
60   addMarkers(~Lon, ~Lat, icon = salmon)
61
62 m
63
64
65 Figure 2.1: Location of stocks used in this data integration. Salmonid icon by Servien
66 (vectorized by T. Michael Keesey)
67 [CC-BY-SA](https://creativecommons.org/licenses/by-sa/3.0/), available at
68 Phylonic(https://phylonic.org/)
37:72 Chunk 2
```



2.2 Datasets

As part of the SASAP project, brood tables for 48 Sockeye salmon stocks were collected. Table 2.1 shows a list of these stocks, along with other regional and location information.

Show 10 entries Search:

Stock.ID	Stock	Region	Sub.Region
101	Washington	WA	WA
102	E.Stuart	Fraser River	Fraser Early Stuart
103	Bowron	Fraser River	Fraser Early Summer
104	Fennell	Fraser River	Fraser Early Summer
105	Gates	Fraser River	Fraser Early Summer
106	Nadina	Fraser River	Fraser Early Summer
107	Pitt	Fraser River	Fraser Early Summer
108	Raft	Fraser River	Fraser Early Summer
109	Scotch	Fraser River	Fraser Early Summer
110	Seymour	Fraser River	Fraser Early Summer

Showing 1 to 10 of 54 entries Previous 1 2 3 4 5 6 Next

These stocks range geographically from Washington to Alaska. Although temporal coverage varies by stock, many of the brood tables were updated in 2016, and some have reconstructions dating back to 1922.

Figure 2.1 indicates the approximate location of the salmon stocks in Table 2.1.



Figure 2.1: Location of stocks used in this data integration. Salmonid icon by Servien (vectorized by T. Michael Keesey)

SASAP

Group

Group Id: SASAP

4 years, 6 months

Contributor since August 4, 2013

2 contributions

4,862 downloads

24 members

Krista B Oke

<http://orcid.org/0000-0002-5415-3534>

Josh Baron

<http://orcid.org/0000-0002-4286-6959>

Rich Brenner

<http://orcid.org/0000-0001-7209-9757>

★ Jeanette Clark

<http://orcid.org/0000-0003-4703-1974>
[First](#)
[1](#)
[2](#)
[3](#)
[4](#)
[5](#)
[6](#)
[Last](#)

DATASETS 1 TO 5 OF 60

[1](#)
[2](#)
[3](#)
[...](#)
[12](#)
[Next](#)
Sort by [Most recent](#)

Alaska Department of Fish and Game, Division of Commercial Fisheries, Central Region. 2018. **Chinook age, sex, and length data from East Side Cook Inlet, Alaska, 1970-2012**. Knowledge Network for Biocomplexity. urn:uuid:16763faf-9ad6-4a95-bcfc-97d60957e499.



6



Jeanette Clark and Rich Brenner. 2017. **Sockeye salmon brood tables, northeastern Pacific, 1922-2016**. Knowledge Network for Biocomplexity. urn:uuid:c11dff42-b988-437a-afee-58fc62dcd1dc.



6



Alaska Department of Fish and Game. 2018. **Salmon age, sex, and length data from Lower Cook Inlet, Alaska, 1961-2014**. Knowledge Network for Biocomplexity. urn:uuid:99e94ab7-b822-458e-88b3-df0ed1964378.



7



Jared Kibebe and Leslie Jones. 2018. **Glaciers in Alaska with subsetting by watershed and SASAP region**. Knowledge Network for Biocomplexity. urn:uuid:874e1ba2-48d2-4d31-b3fb-682aaf7e984b.



7



Jeanette Clark, Rich Brenner, and Bert Lewis. 2018. **Compiled age, sex, and length data for Alaskan salmon**. Knowledge Network for Biocomplexity. urn:uuid:63a9c8df-3543-44fe-a5d0-746469318f18.



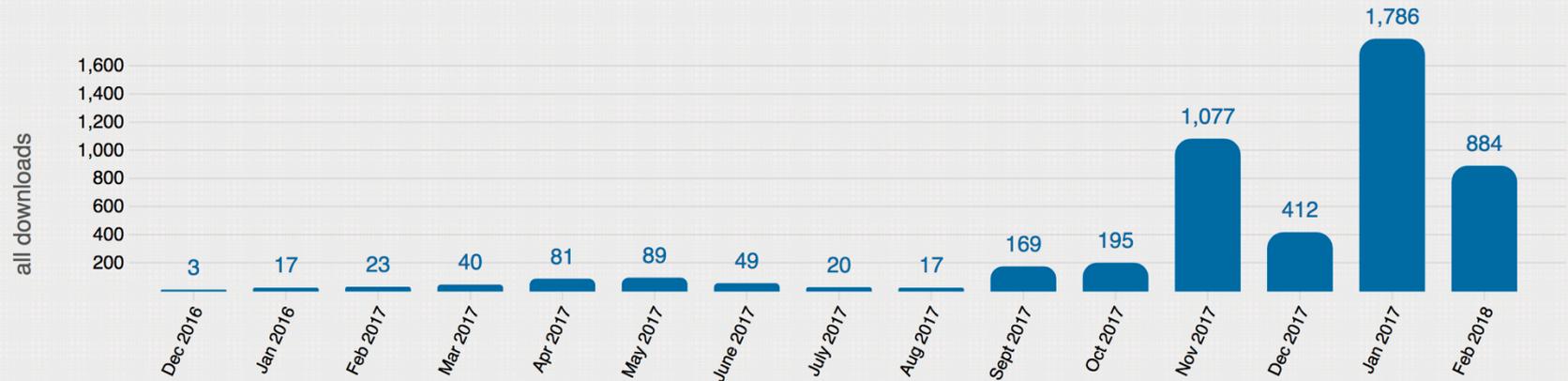
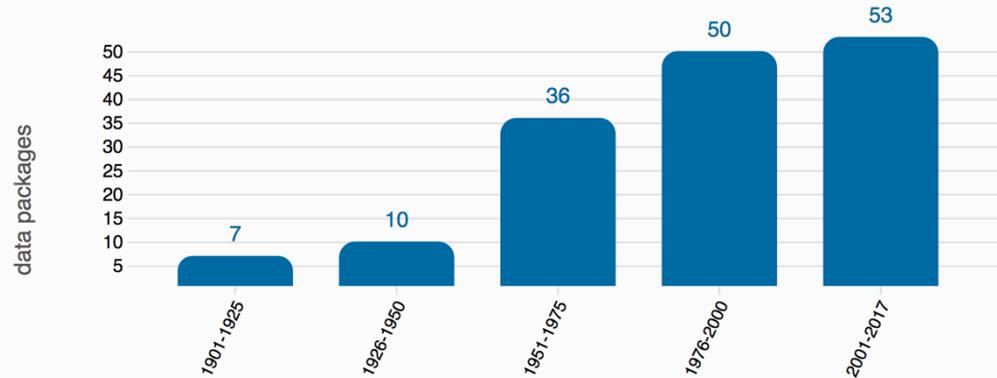
10



Time period of data

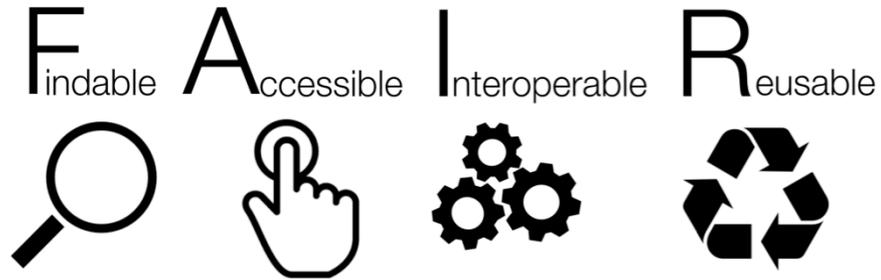
1901 - 2017

The years in which data was collected, regardless of upload date. Only the most recent version of the data package is counted.



Foundational Infrastructure

Providing *findable, accessible* data with *interoperable* infrastructure enabling long term data *reuse* for synthesis



<https://www.force11.org/fairprinciples>