



# Best Practices: Data and Metadata



<https://learning.nceas.ucsb.edu>  
<https://dataone.org>



# Computational Reproducibility

- Preservation enables:
  - Understanding
  - Evaluation
  - Reuse
- Future You!




Metadata



Software




# Computational Workflows





# Data Packages == Research Objects



[Home](#) / [Search](#) / [Metadata](#)

Benjamin Halpern, Melanie Frazier, John Potapenko, Kenneth Casey, Kellee Koenig, et al. 2015. Cumulative human impacts: pressure and cumulative impacts data (2013, all pressures). Knowledge Network for Biocomplexity.  
doi:10.5063/F15718ZN.

[Citations](#)

1

[Downloads](#)

3.6K

[Views](#)

1.8K

[Copy Citation](#)[Quality report](#)[↑ Parent dataset: Cumulative human impacts: Supplementary data](#)

Files in this dataset Package: urn:uuid:975e3a96-c912-4e41-a888-7cccab216bf6

| Name   | File type                 | Size     | Download All |
|--|---------------------------|----------|--------------|
| Metadata: Cumulative human impacts: pressure and cumulative impacts data (2013, all pressures) | EML v2.1.1                | 20 KB    | 1768 views   |
| cumulative_impact_one_2013_global_cumul_impact_2013_mol.zip                                    | <a href="#">More info</a> | ZIP file | 2 GB         |
| pressure_one_2013_artisanal_fishing_mol.zip  | <a href="#">More info</a> | ZIP file | 17 MB        |
| pressure_one_2013_demersal_destructive_fishing_mol.zip   | <a href="#">More info</a> | ZIP file | 218 MB       |

[▶ Show 17 more items in this data set](#)



# Practical Reproducibility



Preserve the data

Preserve the software workflow

Document what you did

Describe how to interpret it all






# **Data and Metadata Guidelines**




# A Data Life Cycle





# A Data Life Cycle





# Guidelines

<https://arcticdata.io/submit/>

- Organizing Data
- File Formats
- Large Data Packages
- Metadata
- Data Identifiers
- Provenance





# Organizing Data

- Understand basics of “tidy” data models
- Design and create effective data tables
  
- **Benefits of tidy data systems**
- Powerful search and filtering
- Handle large, complex data sets
- Enforce data integrity
- Decrease errors from redundant updates





# Not Tidy: Multiple Tables

|         |             | main trunks | reiterated trunks | limbs   | branches | leaves |
|---------|-------------|-------------|-------------------|---------|----------|--------|
| species | tree        | kg          | kg                | kg      | kg       | kg     |
| SESE    | Atlas       | 255144.9    | 48020.6           | 5477.7  | 13433.2  | 1101.2 |
| SESE    | Ballantine  | 221966.4    | 7851.6            | 5922.9  | 11210.0  | 1084.8 |
| SESE    | Bell        | 253246.4    | 5454.3            | 5792.6  | 48500.7  | 1043.4 |
| SESE    | Broken Top  | 130928.9    | 4805.2            | 1608.1  | 5137.4   | 729.9  |
| SESE    | Buena Vista | 128833.0    | 3486.5            | 0.0     | 8552.1   | 518.4  |
| SESE    | Demeter     | 155896.0    | 11085.6           | 3204.3  | 10054.1  | 768.7  |
| SESE    | Epimetheus  | 226987.0    | 12915.7           | 1797.2  | 13585.2  | 1029.4 |
| SESE    | Iluvatar    | 349586.6    | 65003.9           | 12315.6 | 13987.0  | 1481.8 |
| SESE    | Kronos      | 134154.1    | 12204.4           | 7232.7  | 5036.1   | 597.3  |
| SESE    | Pleiades I  | 182385.2    | 3735.0            | 1935.2  | 10846.6  | 762.2  |
| SESE    | Pleiades II | 235838.8    | 11183.4           | 4306.0  | 11306.5  | 877.7  |
| SESE    | Prometheus  | 230414.0    | 25228.9           | 1612.6  | 12458.2  | 1086.0 |
| SESE    | Rhea        | 147114.1    | 487.6             | 730.1   | 5524.2   | 691.2  |
| SESE    | Zeus        | 241367.1    | 2885.5            | 1620.4  | 19104.7  | 954.3  |
| SESE    | 3           | 76.1        | 0.0               | 0.0     | 87.6     | 41.4   |
| SESE    | 4           | 6312.0      | 356.0             | 73.5    | 214.1    | 43.8   |
| SESE    | 5           | 206.0       | 0.0               | 0.0     | 8.7      | 2.5    |
| SESE    | 6E          | 18697.4     | 0.0               | 0.0     | 1055.2   | 66.3   |
| SESE    | 6W          | 14651.5     | 7.7               | 0.0     | 626.3    | 49.6   |
| SESE    | 11          | 614.4       | 0.0               | 0.0     | 28.1     | 17.0   |
| SESE    | 12          | 232.1       | 0.0               | 0.0     | 11.2     | 10.3   |
| SESE    | 18          | 15632.0     | 0.0               | 0.0     | 946.3    | 106.8  |
| SESE    | 19          | 11805.5     | 0.0               | 0.0     | 770.1    | 80.3   |
| SESE    | 20          | 309.5       | 0.0               | 0.0     | 12.5     | 5.9    |
| SESE    | 22          | 25618.3     | 0.0               | 0.0     | 1504.0   | 120.2  |
| SESE    | 23          | 483.7       | 0.0               | 0.0     | 18.9     | 4.5    |
| SESE    | 25          | 87.7        | 0.0               | 0.0     | 4.1      | 1.3    |
| SESE    | 30          | 512.1       | 1.8               | 0.0     | 18.7     | 8.7    |

Table 1

| type  | species | main trunk | reiteration | dry masses (kg) |        |       | TOTAL   | % total |
|-------|---------|------------|-------------|-----------------|--------|-------|---------|---------|
|       |         |            |             | limb            | branch | leaf  |         |         |
| tree  | SESE    | 3569312    | 213247      | 53714           | 230945 | 17192 | 4084409 | 95.3491 |
| tree  | PSME    | 135815     | 0           | 0               | 8338   | 961   | 145114  | 3.3876  |
| tree  | THSE    | 31799      | 0           | 0               | 6343   | 864   | 39006   | 0.9105  |
| tree  | ACMA    | 4444       | 0           | 0               | 925    | 264   | 5634    | 0.1315  |
| tree  | UMCA    | 2921       | 0           | 0               | 937    | 273   | 4131    | 0.0964  |
| shrub | RUSP    | 0          | 0           | 0               | 1974   | 686   | 2660    | 0.0620  |
| fern  | POMU    | 0          | 0           | 0               | 0      | 1271  | 1271    | 0.0296  |
| shrub | VAOV    | 0          | 0           | 0               | 56     | 26    | 552     | 0.0129  |
| shrub | COCO    | 0          | 0           | 0               | 84     | 6     | 289     | 0.0067  |
| fern  | POSC    | 0          | 0           | 0               | 107    | 89    | 196     | 0.0045  |
| tree  | RHPU    | 100        | 0           | 0               | 44     | 18    | 162     | 0.0037  |
| herb  | OXOR    | 0          | 0           | 0               | 0      | 112   | 112     | 0.0026  |
| shrub | VAPA    | 0          | 0           | 0               | 94     | 4     | 99      | 0.0023  |
| tree  | PISI    | 0          | 0           | 0               | 1      | 0     | 1       | 0.0000  |
| tree  | CHLA    | 0          | 0           | 0               | 1      | 0     | 1       | 0.0000  |
| shrub | GASH    | 0          | 0           | 0               | 0      | 0     | 0       | 0.0000  |
| shrub | SACA    | 0          | 0           | 0               | 0      | 0     | 0       | 0.0000  |
|       |         | 3744390    | 213247      | 53714           | 250519 | 21767 | 4283636 |         |

Table 2

|          | main trunk | reiteration | limb  | branch | leaf  | total   | proportion |            |
|----------|------------|-------------|-------|--------|-------|---------|------------|------------|
|          |            |             |       |        |       |         | geophytic  | ecological |
| SESE geo | 3569312    | 213247      | 53714 | 230945 | 17192 | 4084409 | 1.00       | 0.00       |
| SESE epi | 0          | 0           | 0     | 0      | 0     | 0       | 0.00       | 1.00       |
| PSME geo | 135815     | 0           | 0     | 8338   | 961   | 145114  | 1.00       | 0.00       |
| PSME epi | 0          | 0           | 0     | 0      | 0     | 0       | 0.00       | 1.00       |
| TSHE geo | 31740      | 0           | 0     | 6332   | 860   | 38932   | 0.99       | 0.01       |
| TSHE epi | 59         | 0           | 0     | 12     | 4     | 74      | 0.00       | 1.00       |
| ACMA geo | 4444       | 0           | 0     | 925    | 264   | 5634    | 1.00       | 0.00       |
| ACMA epi | 0          | 0           | 0     | 0      | 0     | 0       | 0.00       | 1.00       |

Table 3



# Not Tidy: Inconsistent observations

| AtlasGroveCOMPLETE.xls |             |             |                   |         |          |        |  |  |  |  |  |  |  |  |  |  |  |
|------------------------|-------------|-------------|-------------------|---------|----------|--------|--|--|--|--|--|--|--|--|--|--|--|
|                        |             |             |                   |         |          |        |  |  |  |  |  |  |  |  |  |  |  |
| species                | tree        | main trunks | reiterated trunks | limbs   | branches | leaves |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Atlas       | 255144.9    | 48020.6           | 5477.7  | 13433.2  | 1101.2 |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Ballantine  | 221966.4    | 7651.6            | 5922.9  | 11210.0  | 1084.8 |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Bell        | 253248.4    | 5454.3            | 5792.6  | 48500.7  | 1043.4 |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Broken Top  | 130928.9    | 4805.2            | 1608.1  | 5137.4   | 729.9  |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Buena Vista | 128833.0    | 3486.5            | 0.0     | 8552.1   | 518.4  |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Demeter     | 155896.0    | 1104.3            | 3204.3  | 10054.1  | 768.7  |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Epimetheus  | 226987.0    | 12915.7           | 1797.2  | 13585.2  |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Iluvatar    | 349586.6    | 65003.9           | 11715.6 | 13987.0  |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Kronos      | 134154.1    | 12204.4           | 7237.7  | 5036.1   |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Pleiades I  | 182385.2    | 3735.0            | 1935.2  | 10846.6  |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Pleiades II | 235838.8    | 11183.4           | 4306.0  | 1306.5   |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Prometheus  | 239414.0    | 25228.9           | 1612.6  | 1293.2   |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Rhea        | 143710.4    | 487.8             | 730.1   | 5524.2   |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | Zeus        | 243385.7    | 2885.5            | 1620.4  | 19104.7  |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 3           | 1761.3      | 0.0               | 0.0     | 87.6     |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 4           | 6312.0      | 356.0             | 73.5    | 214.1    |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 5           | 206.0       | 0.0               | 0.0     | 8.7      |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 6E          | 18697.4     | 0.0               | 0.0     | 1055.2   |        |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 6W          | 14651.5     | 7.7               | 0.0     | 626.3    | 49.6   |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 11          | 614.4       | 0.0               | 0.0     | 28.1     | 17.0   |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 12          | 232.1       | 0.0               | 0.0     | 11.2     | 10.3   |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 18          | 15632.0     | 0.0               | 0.0     | 946.3    | 106.8  |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 19          | 11805.5     | 0.0               | 0.0     | 770.1    | 80.3   |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 20          | 309.5       | 0.0               | 0.0     | 12.5     | 5.9    |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 22          | 25618.3     | 0.0               | 0.0     | 1504.0   | 120.2  |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 23          | 483.7       | 0.0               | 0.0     | 18.9     | 4.5    |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 25          | 87.7        | 0.0               | 0.0     | 4.1      | 1.3    |  |  |  |  |  |  |  |  |  |  |  |
| SESE                   | 30          | 512.1       | 1.8               | 0.0     | 18.7     | 8.7    |  |  |  |  |  |  |  |  |  |  |  |

All the same observation?  
No.



# Not Tidy: Inconsistent variables

| AtlasGroveCOMPLETE.xls |             |             |                   |         |          |        |  |  |  |  |  |  |  |          |          |       |         |  |
|------------------------|-------------|-------------|-------------------|---------|----------|--------|--|--|--|--|--|--|--|----------|----------|-------|---------|--|
| species                | tree        | main trunks | reiterated trunks | limbs   | branches | leaves |  |  |  |  |  |  |  | dry mass | ses (kg) |       |         |  |
|                        |             | kg          | kg                | kg      | kg       | kg     |  |  |  |  |  |  |  | branch   | leaf     | TOTAL | % total |  |
| SESE                   | Atlas       | 255144.9    | 48020.6           | 5477.7  | 13433.2  | 1101.2 |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Ballantine  | 221966.4    | 7651.6            | 5922.9  | 11210.0  | 1084.8 |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Bell        | 253246.4    | 5454.3            | 5792.6  | 48500.7  | 1043.4 |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Broken Top  | 130928.9    | 4805.2            | 1608.1  | 5137.4   | 729.9  |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Buena Vista | 128833.0    | 3486.5            | 0.0     | 8552.1   | 518.4  |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Demeter     | 155896.0    | 11085.6           | 3204.3  | 10054.1  | 768.7  |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Epimetheus  | 226987.0    | 12915.7           | 1797.2  | 13585.2  | 1029.4 |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Iluvatar    | 349586.6    | 65003.9           | 12315.6 | 13987.0  | 1481.8 |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Kronos      | 134154.1    | 12204.4           | 7232.7  | 5036     |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Pleiades I  | 182385.2    | 3735.0            | 1935.2  | 10846    |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Pleiades II | 235838.8    | 11183.4           | 4306.0  | 11306    |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Prometheus  | 239414.0    | 25228.9           | 1612.6  | 12456    |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Rhea        | 143710.4    | 487.8             | 730.1   | 5524     |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | Zeus        | 243385.7    | 2885.5            | 1620.4  | 19104    |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 3           | 1761.3      | 0.0               | 0.0     | 87       |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 4           | 6312.0      | 356.0             | 73.5    | 214      |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 5           | 206.0       | 0.0               | 0.0     | 8        |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 6E          | 18897.4     | 0.0               | 0.0     | 1055     |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 6W          | 14651.5     | 7.7               | 0.0     | 626      |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 11          | 614.4       | 0.0               | 0.0     | 28       |        |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 12          | 232.1       | 0.0               | 0.0     | 11.2     | 10.3   |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 18          | 15632.0     | 0.0               | 0.0     | 946.3    | 106.8  |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 19          | 11805.5     | 0.0               | 0.0     | 770.1    | 80.3   |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 20          | 309.5       | 0.0               | 0.0     | 12.5     | 5.9    |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 22          | 25618.3     | 0.0               | 0.0     | 1504.0   | 120.2  |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 23          | 483.7       | 0.0               | 0.0     | 18.9     | 4.5    |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 25          | 87.7        | 0.0               | 0.0     | 4.1      | 1.3    |  |  |  |  |  |  |  |          |          |       |         |  |
| SESE                   | 30          | 512.1       | 1.8               | 0.0     | 18.7     | 8.7    |  |  |  |  |  |  |  |          |          |       |         |  |

All the same variable?  
No.





# Not Tidy: Marginal info

AtlasGroveCOMPLETE.xls

|         |             | C           | D                 | E       | F        | G      | H | I | J | K      | L       | M       | N      | O      | P      | Q       |            |         |
|---------|-------------|-------------|-------------------|---------|----------|--------|---|---|---|--------|---------|---------|--------|--------|--------|---------|------------|---------|
| species | tree        | main trunks | reiterated trunks | limbs   | branches | leaves |   |   |   |        |         |         |        |        |        | % total |            |         |
| SESE    | Atlas       | 255144.9    | 48020.6           | 5477.7  | 13433.2  | 1101.2 |   |   |   | tree   | SESE    | 3569312 | 213247 | 53714  | 230945 | 17192   | 4084409    | 95.3491 |
| SESE    | Ballantine  | 221966.4    | 7651.6            | 5922.9  | 11210.0  | 1084.8 |   |   |   | tree   | PSME    | 135815  | 0      | 0      | 8338   | 961     | 145114     | 3.3876  |
| SESE    | Bell        | 253246.4    | 5454.3            | 5792.6  | 48500.7  | 1043.4 |   |   |   | tree   | THSE    | 31799   | 0      | 0      | 6343   | 864     | 39006      | 0.9105  |
| SESE    | Broken Top  | 130928.9    | 4805.2            | 1608.1  | 5137.4   | 729.9  |   |   |   | tree   | ACMA    | 4444    | 0      | 0      | 925    | 264     | 5634       | 0.1315  |
| SESE    | Buena Vista | 128833.0    | 3486.5            | 0.0     | 8552.1   | 518.4  |   |   |   | tree   | UMCA    | 2921    | 0      | 0      | 937    | 273     | 4131       | 0.0964  |
| SESE    | Demeter     | 155896.0    | 11085.6           | 3204.3  | 10054.1  | 768.7  |   |   |   | shrub  | RUSP    | 0       | 0      | 0      | 1974   | 686     | 2660       | 0.0620  |
| SESE    | Epimetheus  | 226987.0    | 12915.7           | 1797.2  | 13585.2  | 1029.4 |   |   |   | fern   | POMU    | 0       | 0      | 0      | 0      | 1271    | 1271       | 0.0296  |
| SESE    | Iluvatar    | 349586.6    | 65003.9           | 12315.6 | 13987.0  | 1481.8 |   |   |   | shrub  | VAOV    | 0       | 0      | 0      | 526    | 26      | 552        | 0.0129  |
| SESE    | Kronos      | 134154.1    | 12204.4           | 7232.7  | 5036.1   | 597.3  |   |   |   | shrub  | COCO    | 0       | 0      | 0      | 284    | 6       | 289        | 0.0067  |
| SESE    | Pleiades I  | 182385.2    | 3735.0            | 1935.2  | 10846.6  | 762.2  |   |   |   | fern   | POSC    | 0       | 0      | 0      | 107    | 89      | 196        | 0.0045  |
| SESE    | Pleiades II | 235838.8    | 11183.4           | 4306.0  | 11306.5  | 877.7  |   |   |   | tree   | RHPU    | 100     | 0      | 0      | 44     | 18      | 162        | 0.0037  |
| SESE    | Prometheus  | 239414.0    | 25228.9           | 1612.6  | 12458.2  | 1086.0 |   |   |   | herb   | OXOR    | 0       | 0      | 0      | 0      | 112     | 112        | 0.0026  |
| SESE    | Rhea        | 143710.4    | 487.8             | 730.1   | 5524.2   | 691.2  |   |   |   | shrub  | VAPA    | 0       | 0      | 0      | 94     | 4       | 99         | 0.0023  |
| SESE    | Zeus        | 243385.7    | 2885.5            | 1620.4  | 19104.7  | 954.3  |   |   |   | tree   | PISI    | 0       | 0      | 0      | 1      | 0       | 1          | 0.0000  |
| SESE    | 3           | 1761.3      | 0.0               | 0.0     | 87.6     | 41.4   |   |   |   | tree   | CHLA    | 0       | 0      | 0      | 1      | 0       | 1          | 0.0000  |
| SESE    | 4           | 6312.0      | 356.0             | 73.5    | 214.1    | 43.8   |   |   |   | shrub  | GASH    | 0       | 0      | 0      | 0      | 0       | 0          | 0.0000  |
| SESE    | 5           | 206.0       | 0.0               | 0.0     | 8.7      | 2.5    |   |   |   | shrub  | SACA    | 0       | 0      | 0      | 0      | 0       | 0          | 0.0000  |
| SESE    | 6E          | 18897.4     | 0.0               | 0.0     | 1055.2   | 66.3   |   |   |   |        | 3744390 | 213247  | 53714  | 250519 | 21767  | 4283636 |            |         |
| SESE    | 6W          | 14651.5     | 7.7               | 0.0     | 626.3    | 49.6   |   |   |   |        |         |         |        |        |        |         | proportion |         |
| SESE    | 11          | 614.4       | 0.0               | 0.0     | 28.1     | 17.0   |   |   |   |        |         |         |        |        |        |         | geophytic  |         |
| SESE    | 12          | 232.1       | 0.0               | 0.0     | 11.2     | 10.3   |   |   |   |        |         |         |        |        |        |         | 1.00       |         |
| SESE    | 18          | 15632.0     |                   |         |          |        |   |   |   | SESE   | SE epi  | 3569312 | 213247 | 53714  | 230945 | 17192   | 4084409    |         |
| SESE    | 19          | 11805.5     |                   |         |          |        |   |   |   | SE epi | 0       | 0       | 0      | 0      | 0      | 0       |            |         |
| SESE    | 20          | 309.5       |                   |         |          |        |   |   |   | ME geo | 135815  | 0       | 0      | 8338   | 961    | 145114  |            |         |
| SESE    | 22          | 25618.3     |                   |         |          |        |   |   |   | ME epi | 0       | 0       | 0      | 0      | 0      | 0       |            |         |
| SESE    | 23          | 483.7       |                   |         |          |        |   |   |   | HE geo | 31740   | 0       | 0      | 6332   | 860    | 38932   |            |         |
| SESE    | 25          | 87.7        |                   |         |          |        |   |   |   | HE epi | 59      | 0       | 0      | 12     | 4      | 74      |            |         |
| SESE    | 30          | 512.1       |                   |         |          |        |   |   |   | MA geo | 4444    | 0       | 0      | 925    | 264    | 5634    |            |         |
|         |             |             |                   |         |          |        |   |   |   | MA epi | 0       | 0       | 0      | 0      | 0      | 0       |            |         |

Marginal  
sums and  
totals

main trunk

reiteration

limb

branch

leaf

total

main trunk

reiteration

limb



# Data Modeling 101

| <b>id</b> | <b>date</b> | <b>site</b> | <b>elev</b> | <b>sp1code</b> | <b>sp1height</b> | <b>sp2code</b> | <b>sp2height</b> |
|-----------|-------------|-------------|-------------|----------------|------------------|----------------|------------------|
| 1         | 2017-10-10  | 1           | 3.7         | DAPU           | 4.6              | DAMA           | 4.5              |
| 2         | 2017-09-05  | 2           | 3.2         | DAMA           | 3.5              | DAPU           | 3.9              |

- Denormalized data (aka, not Tidy)
- Observations about different entities combined



# Tidy Data (observe one entity per table)

- Species observations

| <b>id</b> | <b>date</b> | <b>site</b> | <b>spcode</b> | <b>height</b> |
|-----------|-------------|-------------|---------------|---------------|
| 1         | 2017-10-10  | 1           | DAPU          | 4.6           |
| 2         | 2017-09-05  | 2           | DAMA          | 3.5           |
| 3         | 2017-10-10  | 1           | DAMA          | 4.5           |
| 4         | 2017-09-05  | 2           | DAPU          | 3.9           |

- Site observations

| <b>site</b> | <b>name</b> | <b>elev</b> | <b>temp</b> |
|-------------|-------------|-------------|-------------|
| 1           | Taku        | 3.7         | 21.2        |
| 2           | Lituya      | 3.2         | 23.1        |



# Tidy Data (Relational)

Join Key

- Species observations

| <b>id</b> | <b>date</b> | <b>site</b> | <b>spcode</b> | <b>height</b> |
|-----------|-------------|-------------|---------------|---------------|
| 1         | 2017-10-10  | 1           | DAPU          | 4.6           |
| 2         | 2017-09-05  | 2           | DAMA          | 3.5           |
| 3         | 2017-10-10  | 1           | DAMA          | 4.5           |
| 4         | 2017-09-05  | 2           | DAPU          | 3.9           |

- Site observations

| <b>site</b> | <b>name</b> | <b>elev</b> | <b>temp</b> |
|-------------|-------------|-------------|-------------|
| 1           | Taku        | 3.7         | 21.2        |
| 2           | Lituya      | 3.2         | 23.1        |



# Organizing Data: Best Practices

- **Some Simple Guidelines for Effective Data Management.**
  - Borer et al. 2009. Bulletin of the Ecological Society of America. <https://doi.org/10.1890/0012-9623-90.2.205>
- **Nine simple ways to make it easier to (re)use your data.**
  - White et al. 2013. Ideas in Ecology and Evolution 6. <https://doi.org/10.4033/iee.2013.6b.6.f>



# Organizing Data: Best Practices

- **Scripts** for all data manipulation
  - Uncorrected raw data file
  - Document processing in scripts
- **Design to add rows, not columns**
  - Each column one variable
  - Each row one observation
- **Nonproprietary file formats**
  - Descriptive names, no spaces
  - Header line



# File Formats

<https://arcticdata.io/submit/#file-format-guidelines>

- **Open Formats**
  - **Text** - support long term access and preservation
  - **Open binary formats** (NetCDF, HDF5)
- Any (meta)data is better than none
  - Microsoft Excel: common but proprietary
  - Export GIS data to ESRI shapefiles
  - Export MATLAB, IDL, etc. to NetCDF

Always bet  
on text!





# Large Data Packages (~ Terabytes)

- Talk to the data center early
- Tile data structures by subset
  - Spatial regions
  - Temporal windows
  - Measured variables
- Use efficient tools (NetCDF, HDF)
  - Compact data format
  - Parallel read/write libraries



# **Metadata Guidelines**



# Metadata: the Goal

- Target a typical researcher (maybe you!)
- 30+ years from now
  
- Goal
  - Understand
  - Interpret
  - Re-use





# Metadata: the Goal

- **What** was measured?
- **Who** did it?
- **When** and **where**?
- **How?** (data structure & methods)
- **Why?** (science context)
- **Attribution & Licensing**





# Metadata: Bibliographic Details

- **Global Identifier** (e.g., DOI)
- **Descriptive title**
  - topic, geographic location, dates, and, if applicable, the scale of the data
- **Descriptive abstract**
  - brief overview of the specific contents and purpose of the data package.
- **Funding** information (award number and sponsor).
- **People and organizations**
  - **Creators** – who should be cited for the data set
  - Contacts
  - Contributors
  - Sponsors, and more



**Metadata**



# Metadata: Discovery Details

- **Geospatial coverage**
  - Field and laboratory sampling locations
  - including place names and precise coordinates
- **Temporal Coverage**
  - When measurements were made
  - To what time period do measurements apply
  - Might be calendar times, or geologic times
- **Taxonomic Coverage**
  - What species were measured
  - Taxonomy standards and procedures
- Other contextual information



**Metadata**



# Metadata: Interpretation Details

- Field and laboratory data **collection methods**
- Full **experimental and project design**, and relationship to data
- Full field and laboratory sample **processing methods**
- **Sampling quality control** procedures
  
- Analysis and modeling methods
  - **Provenance** information
  - **Hardware** and **software** used
    - including make, model, and version
  - **Computing quality control** procedures
    - testing, code review, etc.



**Metadata**



# Metadata: Data Structure and Contents

- **Data model description**
- **Data object descriptions (granules)**
  - Tables
  - Images
  - Matrices
  - Spatial layers, etc.
- **Variable information** (attributes/parameters)
  - Definitions / link to methods
  - Standardized measurement types
  - Units
  - Coded values
  - Missing value codes



**Metadata**



# Metadata: Rights and Attribution

- **Scientific rights and expectations**
  - **Citation format**
  - **Attribution expectations**
  - **Reuse rights**
    - Who may reuse data, and for what purposes
  - **Redistribution rights**
    - Who may copy and redistribute data and metadata
- **Legal terms and conditions**
  - **Licensing terms**





# Metadata Standards

- Ecological Metadata Language (EML)
- Geospatial Metadata Standards
  - (ISO 19115\*, ISO 19139)
- Biological Data Profile (BDP)
- Dublin Core
- Darwin Core
- PREMIS and METS
- ... and the list goes on



Metadata

Research and Analysis Section. 2017. Resident vs Nonresident Workers Wages in the Alaskan Seafood and Fishing Processing Industry. KNB Test Node. urn:uuid:d52fa737-fdc1-4192-9c60-b2ad145aa7f9.

| Files  | Size  | Type | Status  |   |
|--|-------|------|---|---|
|  Resident vs Nonresident Workers Wages in the Alaskan Seafood and Fishing Processing Industry | 26 KB |      |   |  |
|  AISFPOver.pdf   | 6 KB  | Data |  |  |
|  processingWorkersWages4.csv   | 6 KB  | Data |  |  |
|  ANSFPOver.pdf   | 6 KB  | Data |  |  |

## Overview \*

### Overview

#### Title \*

A title for this dataset. Include the topic, geographic location, dates, and if applicable, the scale of the data. Write out all abbreviations.

Resident vs Nonresident Workers Wages in the Alaskan Seafood and Fishing Processing Industry

#### Abstract \*

Provide a brief overview that summarizes the specific contents and purpose of this dataset.

These data were taken from Alaska's Department of Labor and Workforce Development website (<http://live.laborstats.alaska.gov/seafood/>), Research and Analysis Section. The csv data file is extracted from the pdfs included in the data package. The data file contains the average wages of resident and nonresident workers in the Alaskan seafood and fishing processing industry from 2001-2015. The data are organized into 8 regions, and 1 'Statewide' region encompassing all 8 regions. For the Northern region data, the large jump in workers in 2013 was due to an employer previously in a different industry being recoded into the seafood processing industry.



## Data Identifiers

Nina J. Karnovsky and Ann M. A. Harding. 2016. At-sea density of foraging little auks (*Alle alle*) near Hornsund Fjord. Arctic Data Center. doi:10.5065/D6MK6B17.

- DOI == Digital Object Identifier
- We assign a DOI to each published data set
- Researchers should cite data they use

 A newer version of this dataset exists. [View it now.](#)

[Home](#) / [Search](#) / [Metadata](#)

Julie McKnight. 2015. **Thule, Greenland CO<sub>2</sub> flux, soil moisture and temperature - 2015**. Arctic Data Center. [doi:10.18739/A2ZK3V](https://doi.org/10.18739/A2ZK3V).



- Each update has a unique identifier
- Cite the exact version used
- Newer versions are clearly indicated



# Data Usage Metrics

[← Back to search](#) | [Home](#) / [Search](#) / [Metadata](#)

Hajo Eicken. 2009. **The State of the Arctic Sea Ice Cover: Sustaining the integrated seasonal ice zone observing network.** Arctic Data Center. urn:uuid:3fb067ab-a8c6-4297-863f-511f1d39233b.



Citations

5

Downloads

101.7K

Views

4.3K

[Copy Citation](#)

[Quality report](#)



5 Citations

x

I.J. Smith, H. Eicken, A.R. Mahoney, R. Van Hale, A.J. Gough, et al. 2016. Surface water mass composition changes captured by cores of Arctic land-fast sea ice. *Continental Shelf Research*. Vol. 118. pp. 154-164. <https://doi.org/10.1016/j.csr.2016.02.008>.

Daisuke Hirano, Yasushi Fukamachi, Eiji Watanabe, Kay I. Ohshima, Katsushi Iwamoto, et al. 2016. A wind-driven, hybrid latent and sensible heat coastal polynya off Barrow, Alaska. *Journal of Geophysical Research: Oceans*. Vol. 121. pp. 980-997. <https://doi.org/10.1002/2015JC011318>.

Megan O&apos;Sadnick, Malcolm Ingham, Hajo Eicken, and Erin Pettit. 2016. In situ field measurements of the temporal evolution of low-frequency sea-ice dielectric properties in relation to temperature, salinity, and microstructure. *The Cryosphere*. Vol. 10. pp. 2923-2940. <https://doi.org/10.5194/tc-10-2923-2016>.

Megan O&apos;Sadnick, Malcolm Ingham, Hajo Eicken, and Erin Pettit. 2016. In situ field measurements of the temporal evolution of low-frequency sea-ice dielectric properties in relation to temperature, salinity, and microstructure. *The Cryosphere*. Vol. 10. pp. 2923-2940. <https://doi.org/10.5194/tc-10-2923-2016>.

P. J. Griewank and D. Notz. 2015. A 1-D modelling study of Arctic sea-ice salinity. *The Cryosphere*. Vol. 9. pp. 305-329. <https://doi.org/10.5194/tc-9-305-2015>.

 101.7K Downloads


For all versions of this data set, the number of times that all or part of this data set was downloaded over time.

These download counts are COUNTER compliant, meaning that downloads from some Internet robots and repeat downloads within a certain time window are excluded.

Drag the slider to visualize a specific time window for the download events.

### 42030 Downloads from Mar 2015 to Sep 2019

Zoom to  year  month  all




 Citations

Views 




# Provenance Metadata

- Simplified view of complex workflows



## Data Table, Image, and Other Data Details


4 sources



### Data Table


|                          |  |                        |   |                  |      |                       |        |                    |  |                 |   |
|--------------------------|--|------------------------|---|------------------|------|-----------------------|--------|--------------------|--|-----------------|---|
| Entity Name              | Total_Aromatic_Alkanes_PWS.csv   |                        |   |                  |      |                       |        |                    |  |                 |   |
|                          | <a href="#">Download</a>   |                        |   |                  |      |                       |        |                    |  |                 |   |
| 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 | <a href="https://cn.dataone.org/cn/v2/resolve/urn:uuid:44108e76-405d-4d58-b1b3-fb4b55e3fff9">https://cn.dataone.org/cn/v2/resolve/urn:uuid:44108e76-405d-4d58-b1b3-fb4b55e3fff9</a>  |                        |   |                  |      |                       |        |                    |  |                 |   |
| 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><b>Simple Text</b></td><td></td></tr><tr><td>Field Delimiter</td><td>,</td></tr></table> | Number of Header Lines | 1 | Record Delimiter | #x0A | Attribute Orientation | column | <b>Simple Text</b> |  | Field Delimiter | , |
| Number of Header Lines   | 1  |                        |   |                  |      |                       |        |                    |  |                 |   |
| Record Delimiter         | #x0A   |                        |   |                  |      |                       |        |                    |  |                 |   |
| Attribute Orientation    | column   |                        |   |                  |      |                       |        |                    |  |                 |   |
| <b>Simple Text</b>       |  |                        |   |                  |      |                       |        |                    |  |                 |   |
| 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\_SI\_Ns.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




# Data package with Provenance





# Rmarkdown as Provenance

```
01-brood-table-integration.Rmd < 31
32 ## Datasets
33
34 As part of the SASAP project, brood tables for 48 Sockeye salmon stocks 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 ````````{r, echo = FALSE}
41 datatable(stocks[, c('Stock.ID','Stock' , 'Region', 'Sub.Region')], rownames = FALSE,
42 caption = "Stock information")
43
44 These stocks range geographically from Washington to Alaska. Although temporal coverage
45 varies by stock, many of the brood tables were updated in 2016, and some have
46 reconstructions dating back to 1922.
47
48 ````{r, echo = FALSE}
49 salmon <- makeIcon("images/salmon_tiny.png",
50                     "images/salmon_big.png",
51                     26, 14)
52
53 m <- leaflet(stocks) %>%
54   setView(~median(stocks$Lon), median(stocks$Lat), zoom = 4) %>%
55   addTiles() %>%
56   addMarkers(~Lon, ~Lat, icon = salmon)
57
58 m
59
60
61
62 Figure 2.1 indicates the approximate location of the salmon stocks in Table 2.1.
63
64 ````{r, echo = FALSE}
65 salmon <- makeIcon("images/salmon_tiny.png",
66                     "images/salmon_big.png",
67                     26, 14)
68
69 m <- leaflet(stocks) %>%
70   setView(~median(stocks$Lon), median(stocks$Lat), zoom = 4) %>%
71   addTiles() %>%
72   addMarkers(~Lon, ~Lat, icon = salmon)
73
74 m
75
76
77
78 Figure 2.1: Location of stocks used in this data integration. Salmonid icon by Servien
79 (vectorized by T. Michael Keesey)
80 [CC-BY-SA](https://creativecommons.org/licenses/by-sa/3.0/), available at
81 [PhyloIcon](http://phyloicon.org/)
```



## 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.

| 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.



# Citing multi-generational workflows

Transitive Credit  
Via  
Provenance





# Guidelines

<https://arcticdata.io/submit/>

- Organizing Data
- File Formats
- Large Data Packages
- Metadata
- Data Identifiers
- Provenance





**<https://learning.nceas.ucsb.edu>**



**<https://dataone.org>**