

Research Data Management

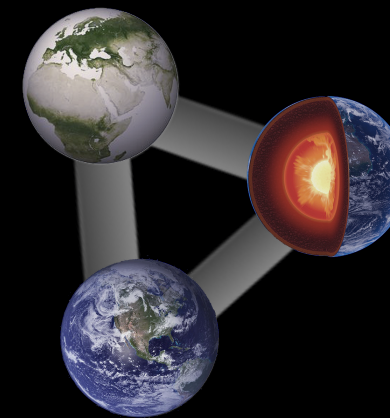
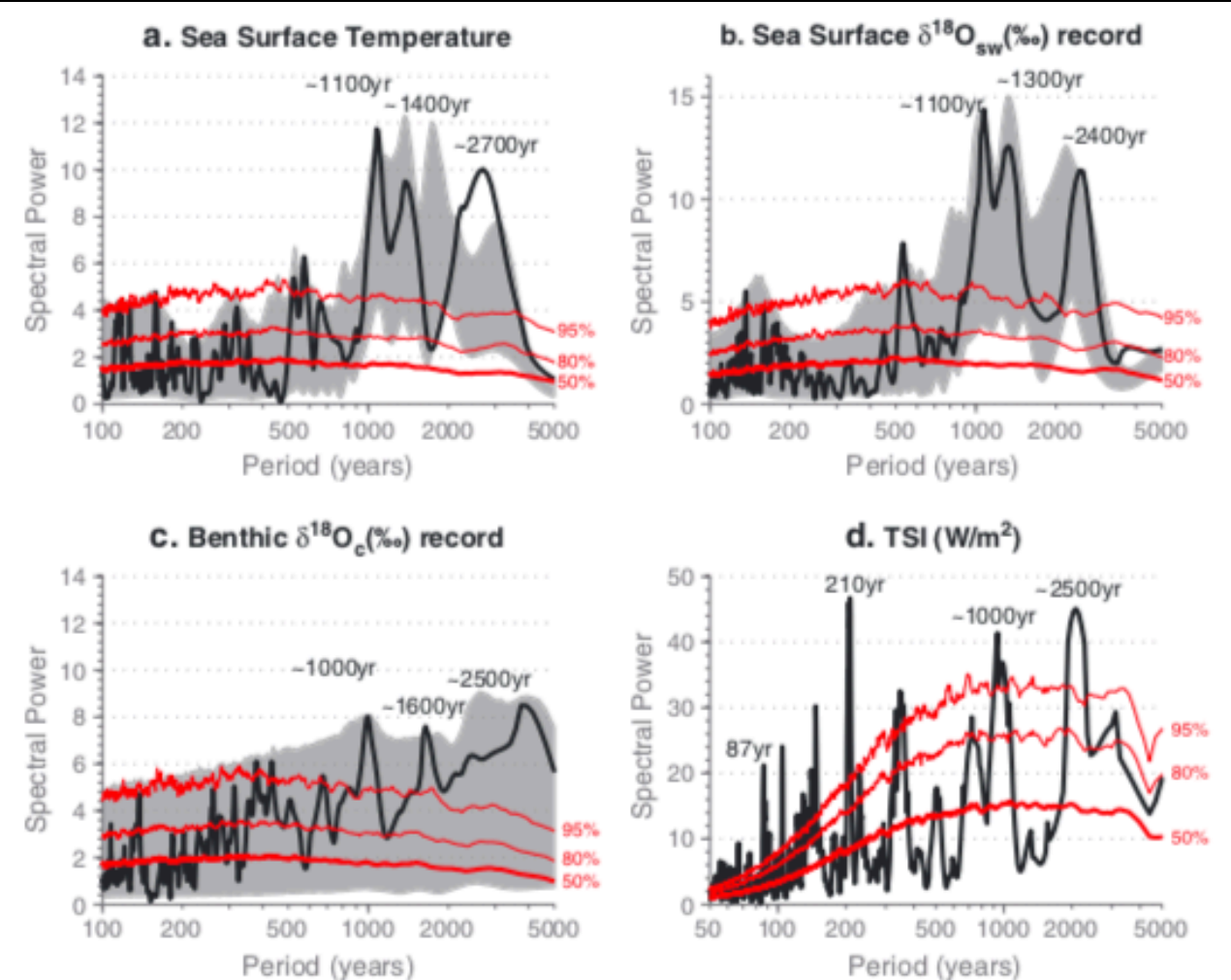
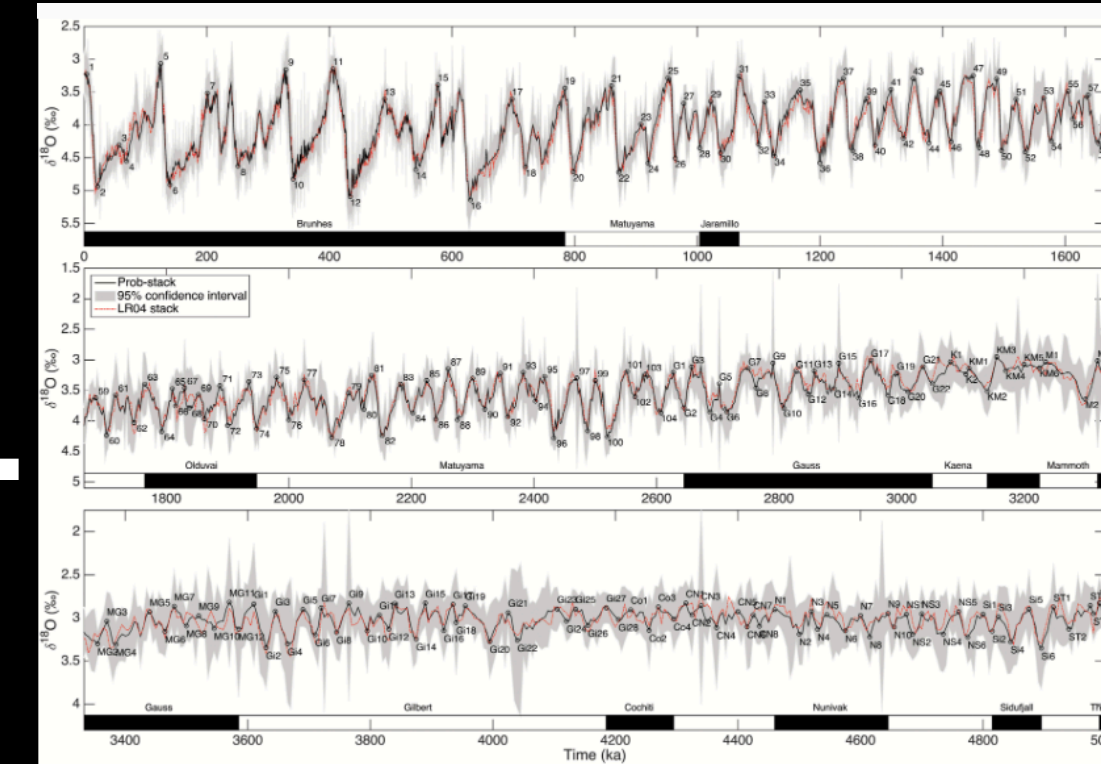
Deborah Khider



USC Viterbi
School of Engineering
Information Sciences Institute

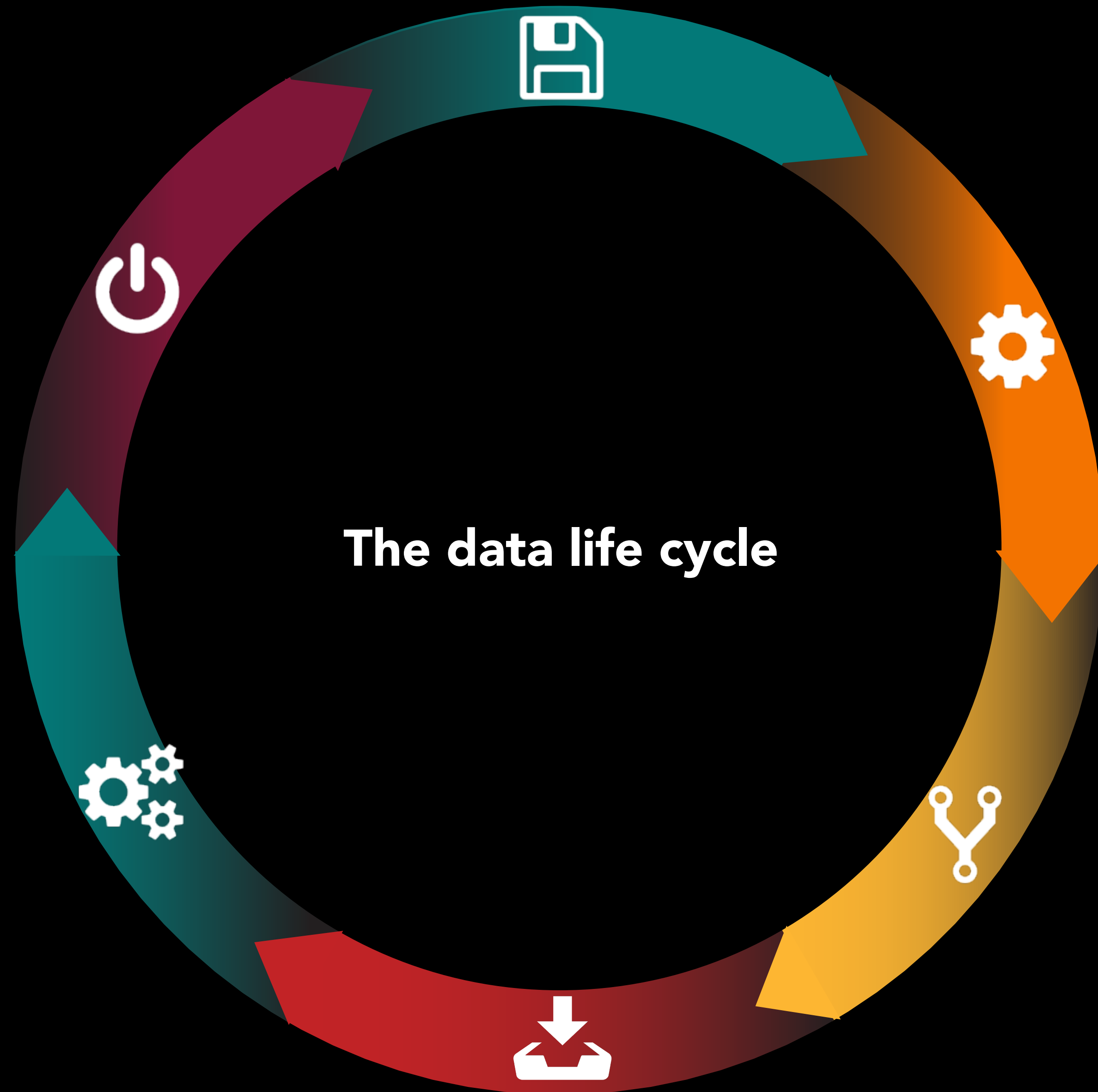


My Data Journey

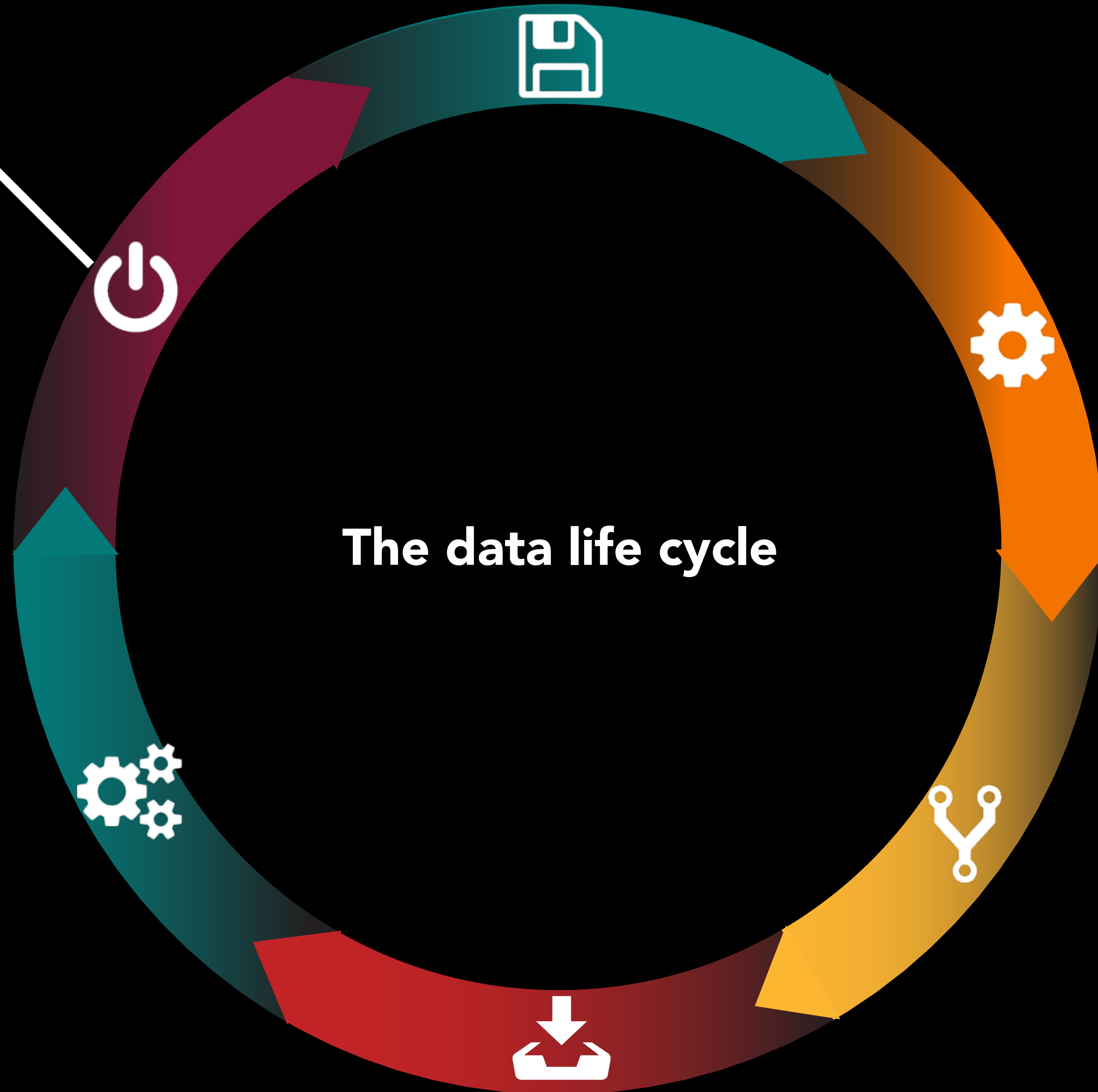


LinkedEarth

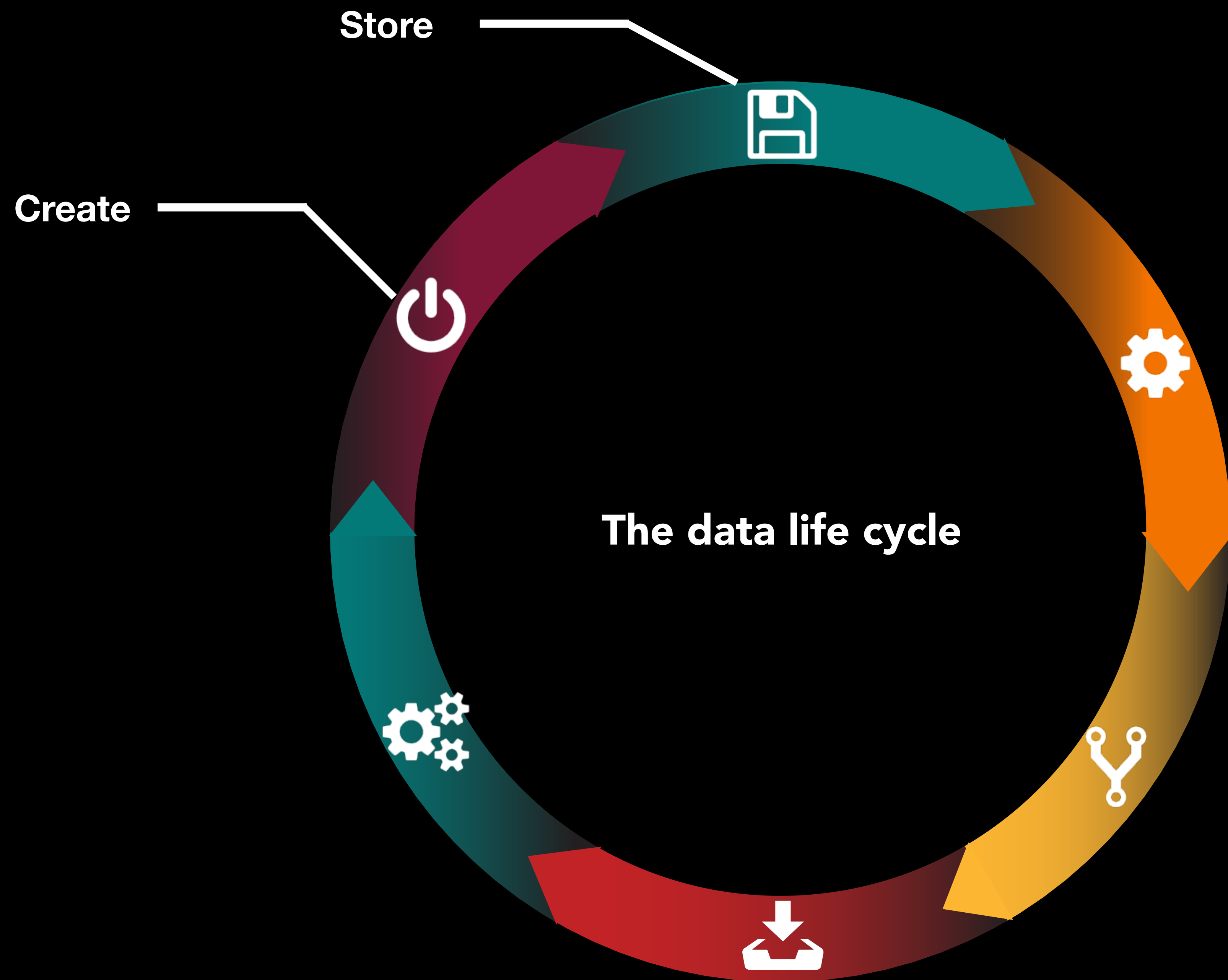


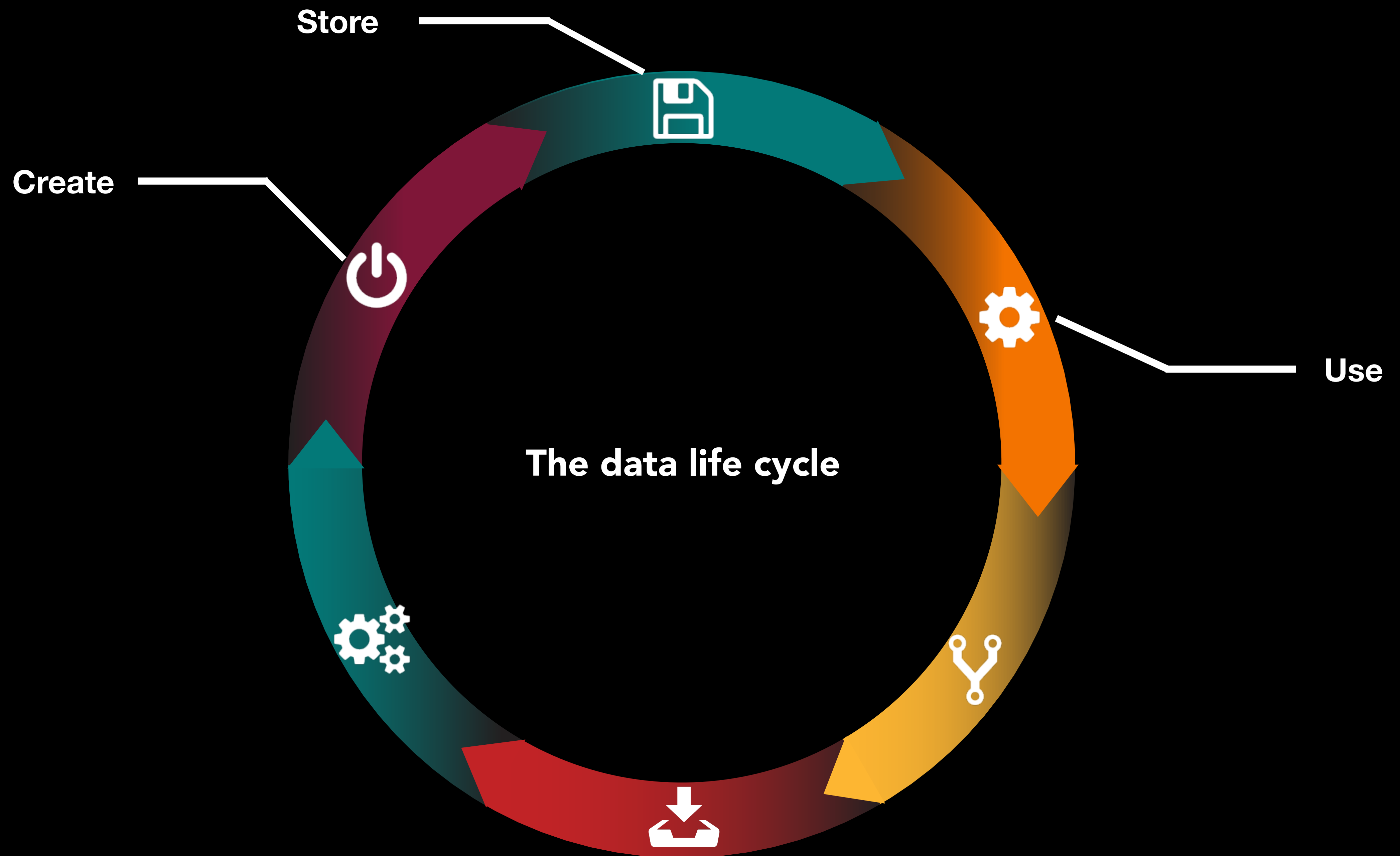


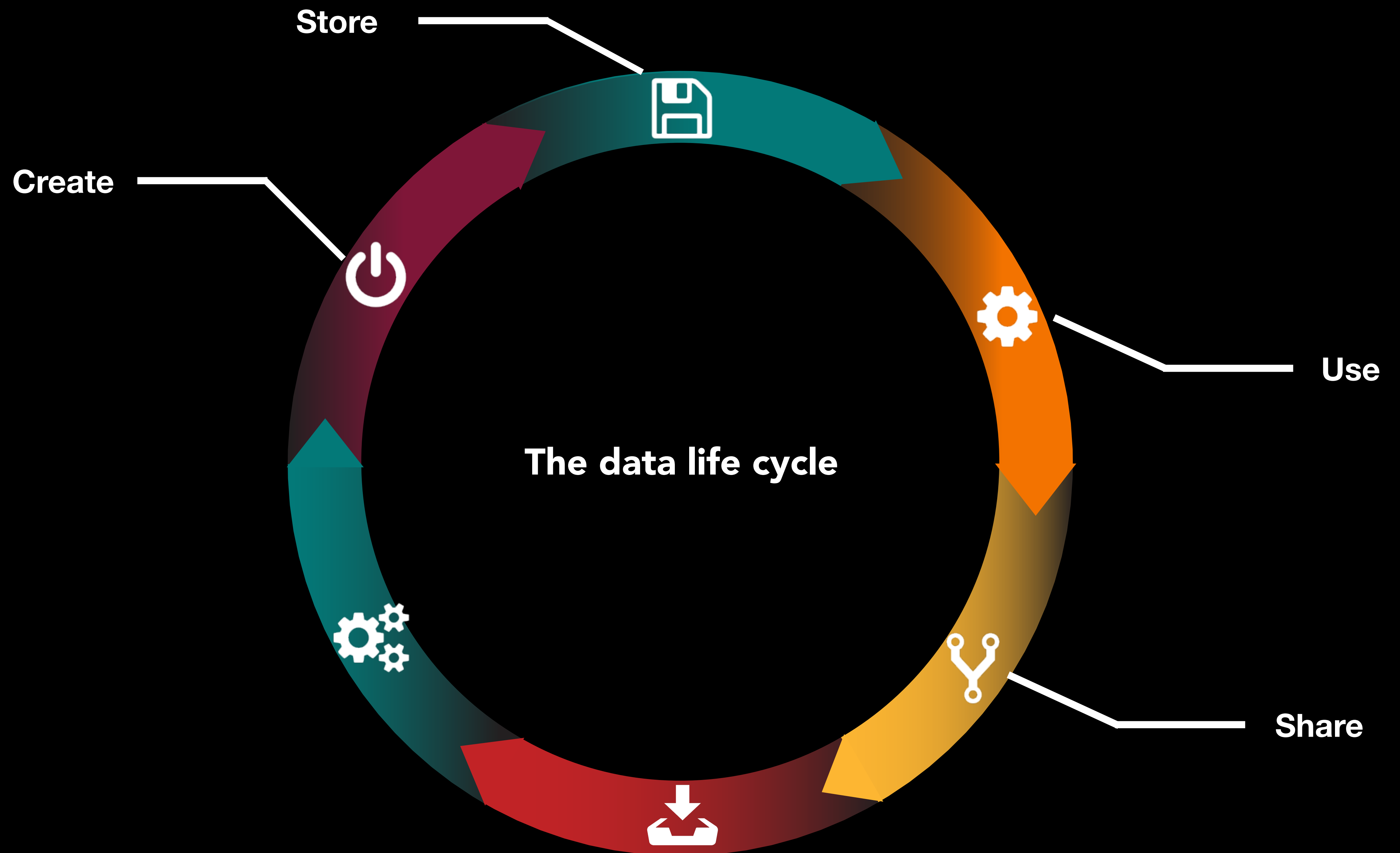
Create

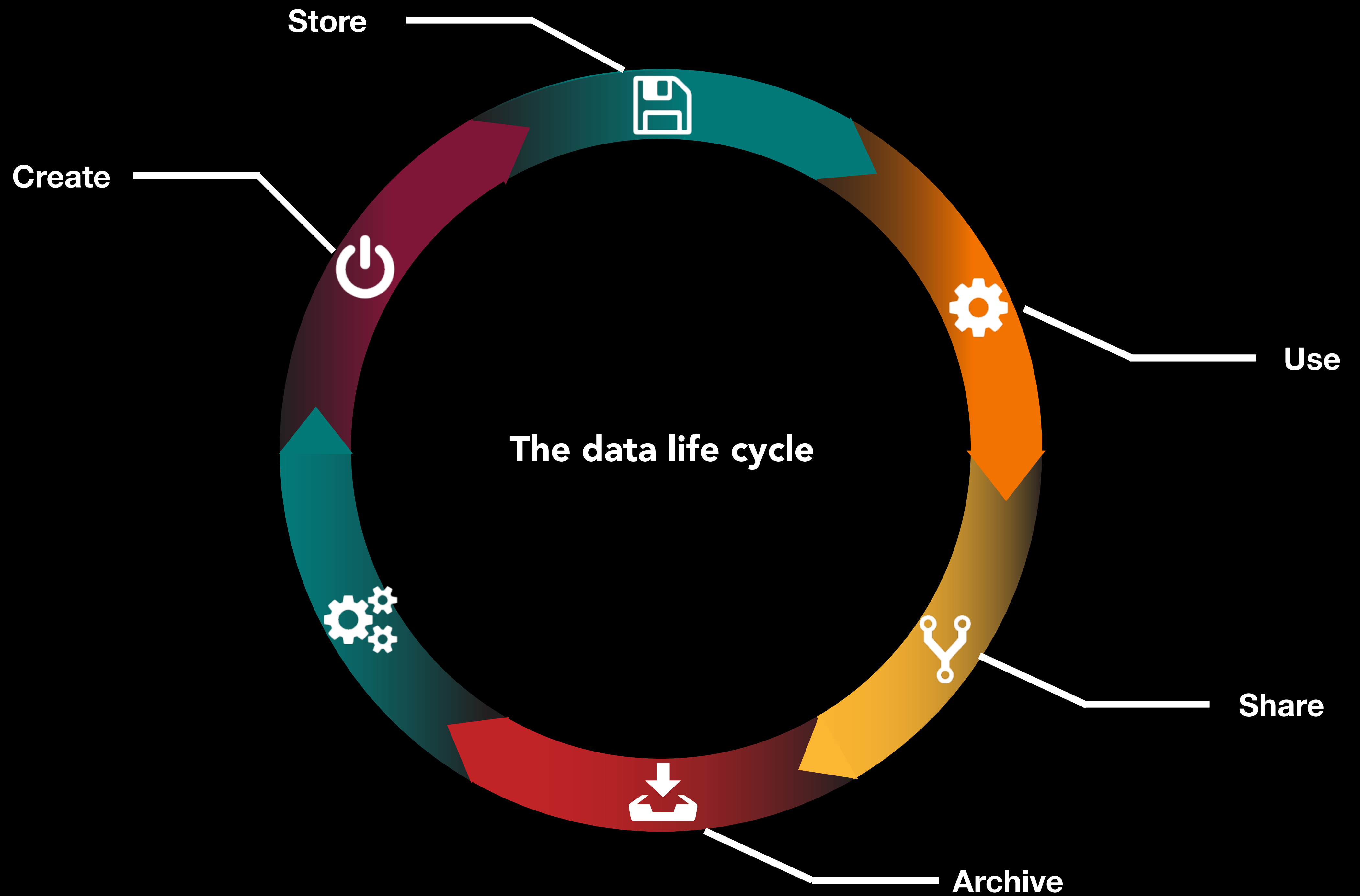


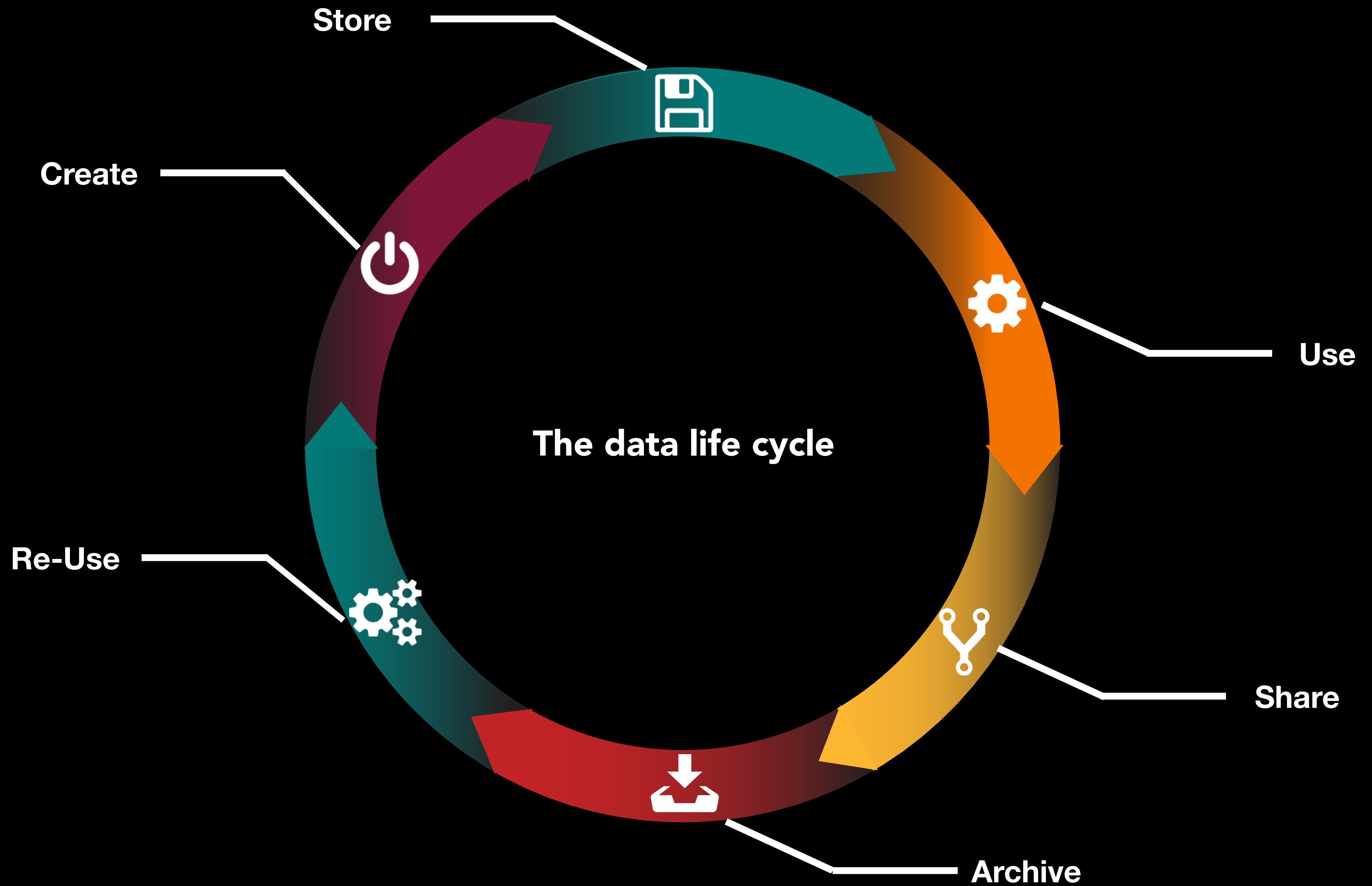
The data life cycle









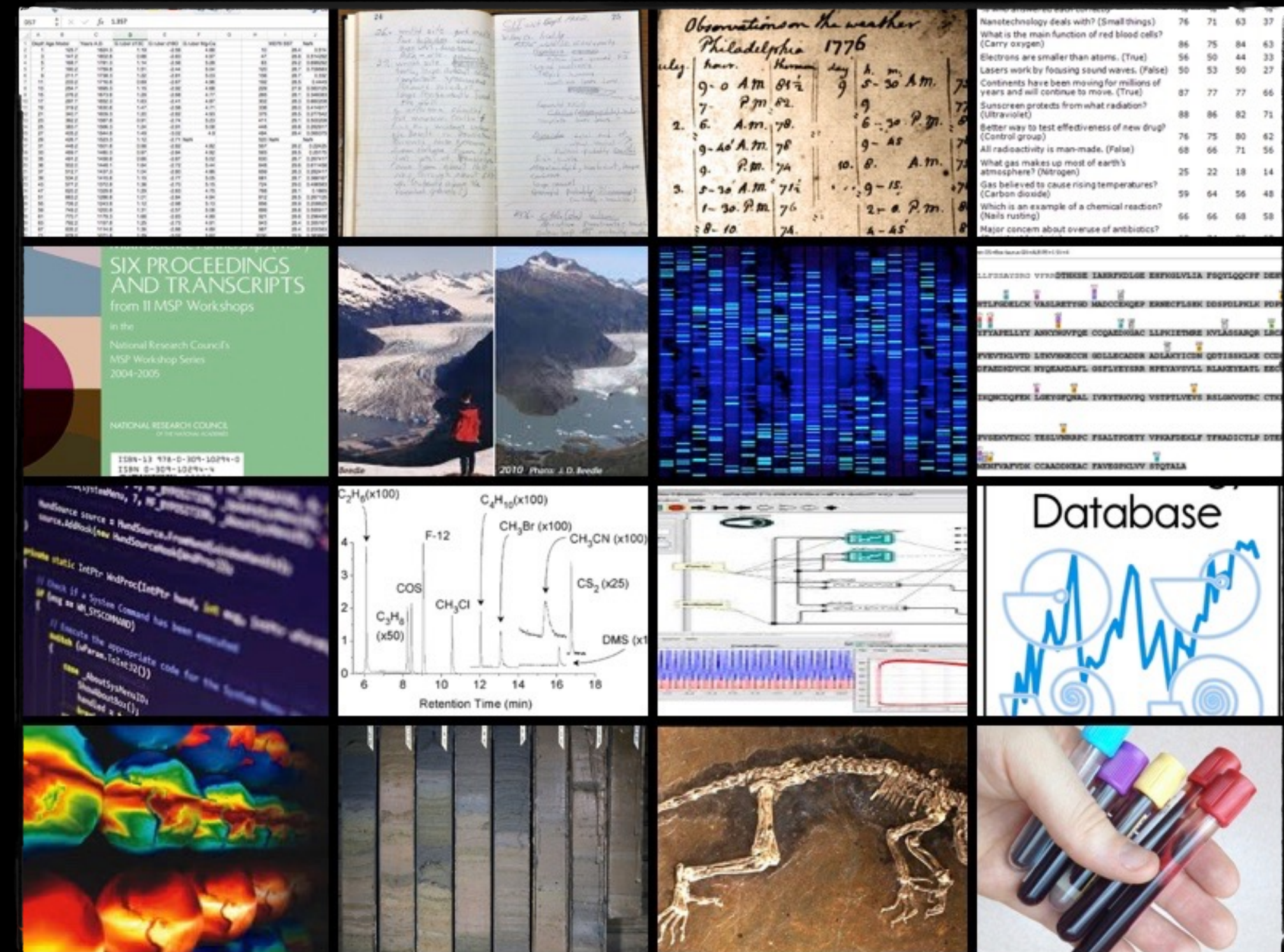


DATA:

"...THE RECORDED FACTUAL
MATERIAL COMMONLY ACCEPTED
IN THE SCIENTIFIC COMMUNITY AS
NECESSARY TO **VALIDATE**
RESEARCH FINDINGS."

DATA:

"...THE RECORDED FACTUAL MATERIAL COMMONLY ACCEPTED IN THE SCIENTIFIC COMMUNITY AS NECESSARY TO **VALIDATE** RESEARCH FINDINGS."



DATA:

"...THE RECORDED FACTUAL MATERIAL COMMONLY ACCEPTED IN THE SCIENTIFIC COMMUNITY AS NECESSARY TO **VALIDATE** RESEARCH FINDINGS."

METADATA:

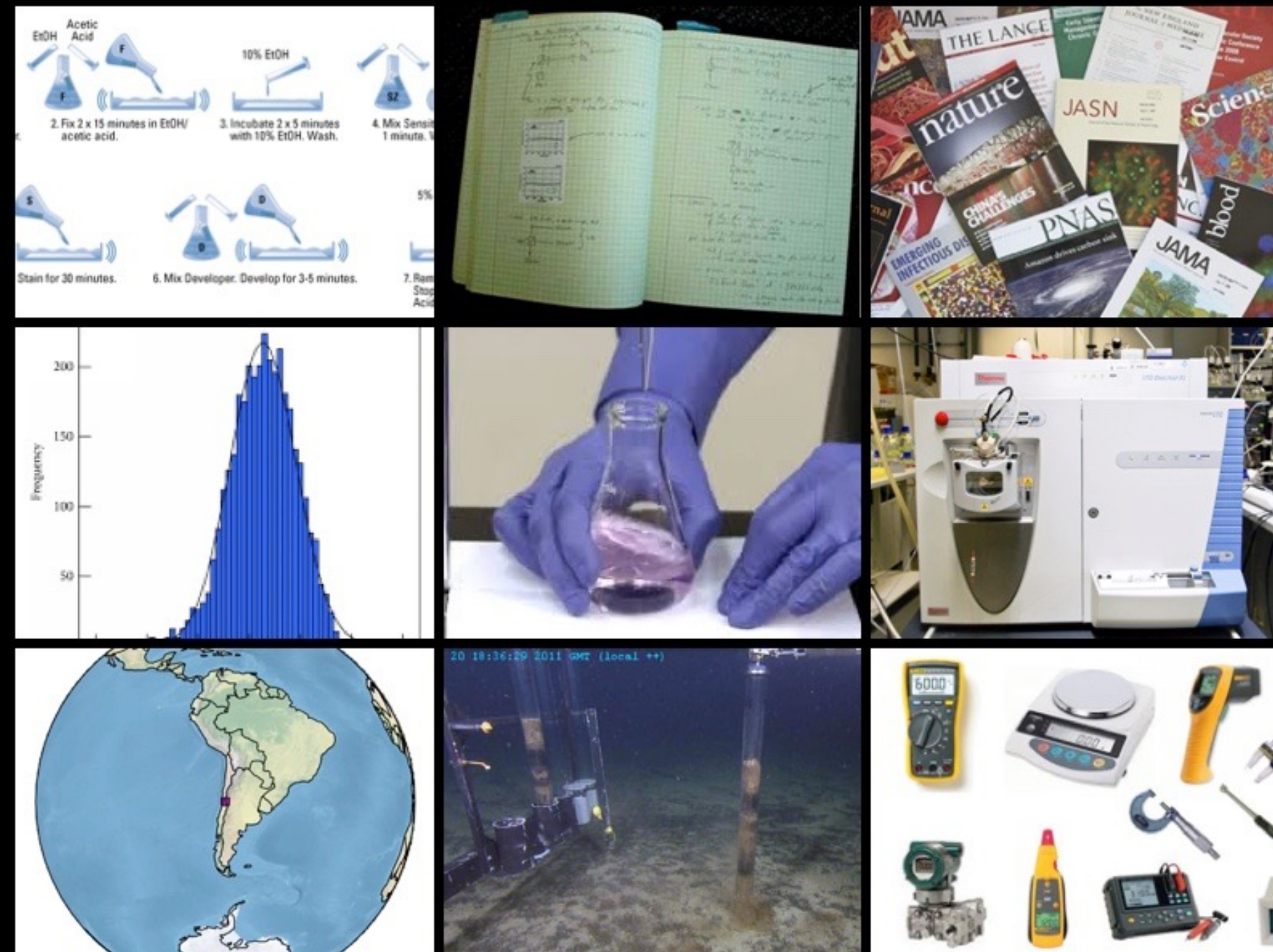
METADATA IS INFORMATION ABOUT THE DATA THAT PROVIDES CONTEXT KEY TO **UNDERSTAND** WHAT THE DATA REPRESENTS.

DATA:

"...THE RECORDED FACTUAL MATERIAL COMMONLY ACCEPTED IN THE SCIENTIFIC COMMUNITY AS NECESSARY TO **VALIDATE** RESEARCH FINDINGS."

METADATA:

METADATA IS INFORMATION ABOUT THE DATA THAT PROVIDES CONTEXT KEY TO **UNDERSTAND** WHAT THE DATA REPRESENTS.



DATA:

"...THE RECORDED FACTUAL MATERIAL COMMONLY ACCEPTED IN THE SCIENTIFIC COMMUNITY AS NECESSARY TO **VALIDATE** RESEARCH FINDINGS."

METADATA:

METADATA IS INFORMATION ABOUT THE DATA THAT PROVIDES CONTEXT KEY TO **UNDERSTAND** WHAT THE DATA REPRESENTS.

DATA MANAGEMENT:

ACTIONS THAT CONTRIBUTE TO EFFECTIVE **STORAGE, PRESERVATION,** AND **REUSE** OF DATA AND METADATA THROUGHOUT THE RESEARCH LIFECYCLE.

Why do you need to know about data management?

- Scientists are changing

Why do you need to know about data management?

- Scientists are changing

Open Data



Open Access



Open Publications

Open Source



Why do you need to know about data management?

- Scientists are changing
- Publishers are changing

Data & Software for Authors

WHAT IS NEEDED?

AGU requires that the underlying data needed to understand, evaluate, and build upon the reported research be available at the time of peer review and publication. Additionally, authors should make available software that has a significant impact on the research. This entails:

1. Depositing the data and software in a trusted repository, as appropriate, and preferably with a DOI
2. Including an [Availability Statement](#) as a separate paragraph in the Open Research section explaining to the reader where and how to access the data and software
3. And including [citation\(s\)](#) to the deposited data and software, in the Reference Section.

Click on the headings below for detailed information on:

- [Models & Simulations](#)
- [Journal-Specific Data Guidance](#)
- [International Geo Sample Numbers](#)

WHAT DATA NEEDS TO BE AVAILABLE?

Primary and processed data used for your research should be preserved and made available.

Generally, the underlying data are considered to be the types of data usually preserved in domain repositories for each discipline. These may include raw data, but are usually the processed or refined data that support and lead to the described results and allow other readers to assess your conclusions and build off your work.

In your paper, cite these data, as well as any data you used from other sources, and include information about access to the data in the availability statement. For model or simulation data, follow [journal specific guidance](#) on prioritizing preserved output; in general, availability of software is most important.

Very large data (greater than 1 terabyte or TB) can be a challenge to preserve as there often fees and additional resources required. One option to consider, institutions often offer solutions for data preservation and compliance. Again, refer to the [journal specific guidance](#) for more information or email DataHelp@agu.org.

Why do you need to know about data management?

- Scientists are changing
- Publishers are changing
- Funders are changing



The screenshot shows the White House website header with the text "the WHITE HOUSE PRESIDENT BARACK OBAMA" and navigation links "Contact Us" and "Get Email Updates". Below the header is a blue navigation bar with the White House seal and links: "BRIEFING ROOM", "ISSUES", "THE ADMINISTRATION", "PARTICIPATE", and "1600 PENN". A search bar is on the right. The main content area has a breadcrumb "HOME · BLOG" and a large title "Expanding Public Access to the Results of Federally Funded Research". Below the title is the date and author "FEBRUARY 22, 2013 AT 12:04 PM ET BY MICHAEL STEBBINS" and three social media icons (Twitter, Facebook, Email). A horizontal line separates the header from the summary text.

the WHITE HOUSE PRESIDENT BARACK OBAMA

Contact Us Get Email Updates

BRIEFING ROOM ISSUES THE ADMINISTRATION PARTICIPATE 1600 PENN

HOME · BLOG


Expanding Public Access to the Results of Federally Funded Research

FEBRUARY 22, 2013 AT 12:04 PM ET BY MICHAEL STEBBINS

Summary: The Obama Administration is committed to the proposition that citizens deserve easy access to the results of research their tax dollars have paid for. That's why, in a policy memorandum released today, OSTP Director John Holdren has directed Federal agencies with more than \$100M in R&D expenditures to develop plans to make the results of federally funded research freely available to the public—generally within one year of publication.

Why do you need to know about data management?

- Scientists are changing
- Publishers are changing
- Funders are changing
- Universities are changing

 UNIVERSITY OF CALIFORNIA

UC Publication Management

Manage your publications.
Participate in the UC Open Access Policy.
Increase the impact of your work.

Select your campus to get started:

[UC Berkeley >>](#)[UC Riverside >>](#)





[UC Davis >>](#)[UC Santa Barbara >>](#)

[UC Irvine >>](#)[UC Santa Cruz >>](#)

[UCLA >>](#)[UC San Diego >>](#)

[UC Merced >>](#)[UCSF >>](#)

Managing your publications

-  We'll scan the web for publications you've authored.
-  Log in (at left) to review what we've found.
-  Claim publications that are yours; reject those that aren't.
-  Upload your manuscript for public display on [eScholarship](#).

Resources and support

Learn more about the [UC Open Access Policy](#).

Get answers to [Frequently Asked Questions](#).

Find out who to contact for [additional support](#).


Logout notice

To protect your accounts from unauthorized access, please lock your workstation or exit your browser after logging out of this site.

Why do you need to know about data management?

- Sharing your data is not only required but it also helps you!

Sharing Detailed Research Data Is Associated with Increased Citation Rate

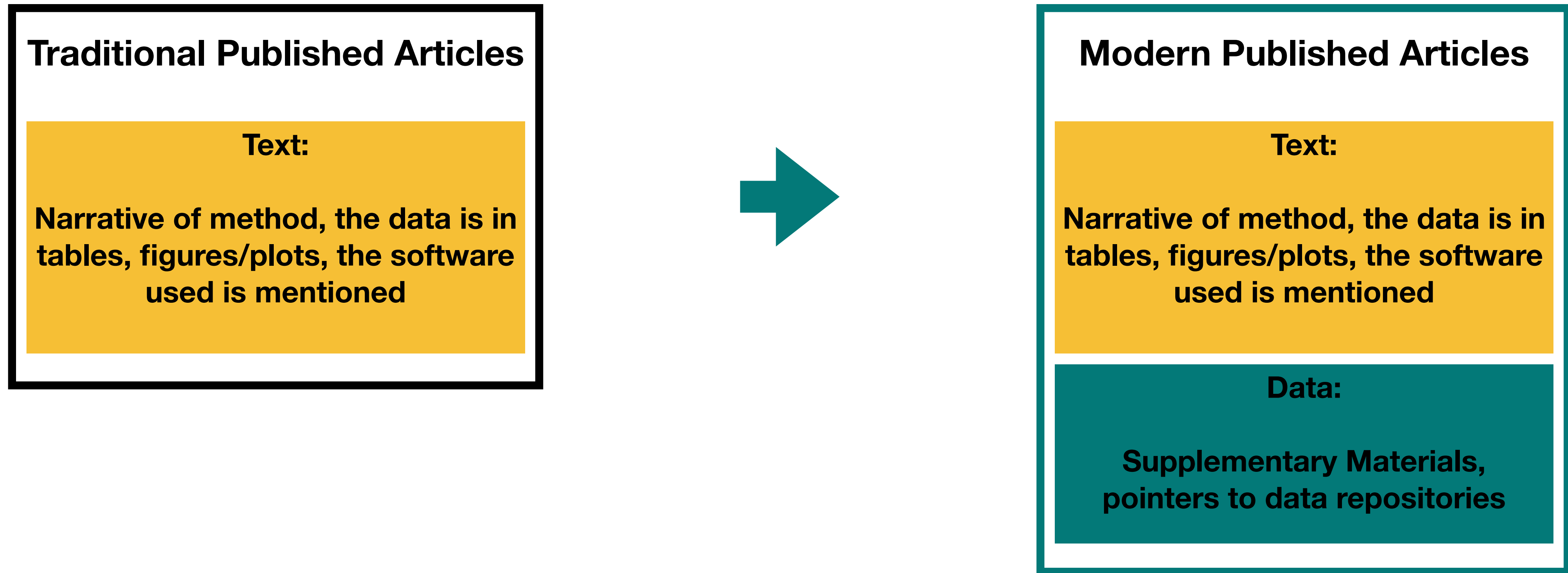
Heather A. Piwowar , Roger S. Day, Douglas B. Fridsma

Published: March 21, 2007 • DOI: 10.1371/journal.pone.0000308 • Featured in PLOS Collections

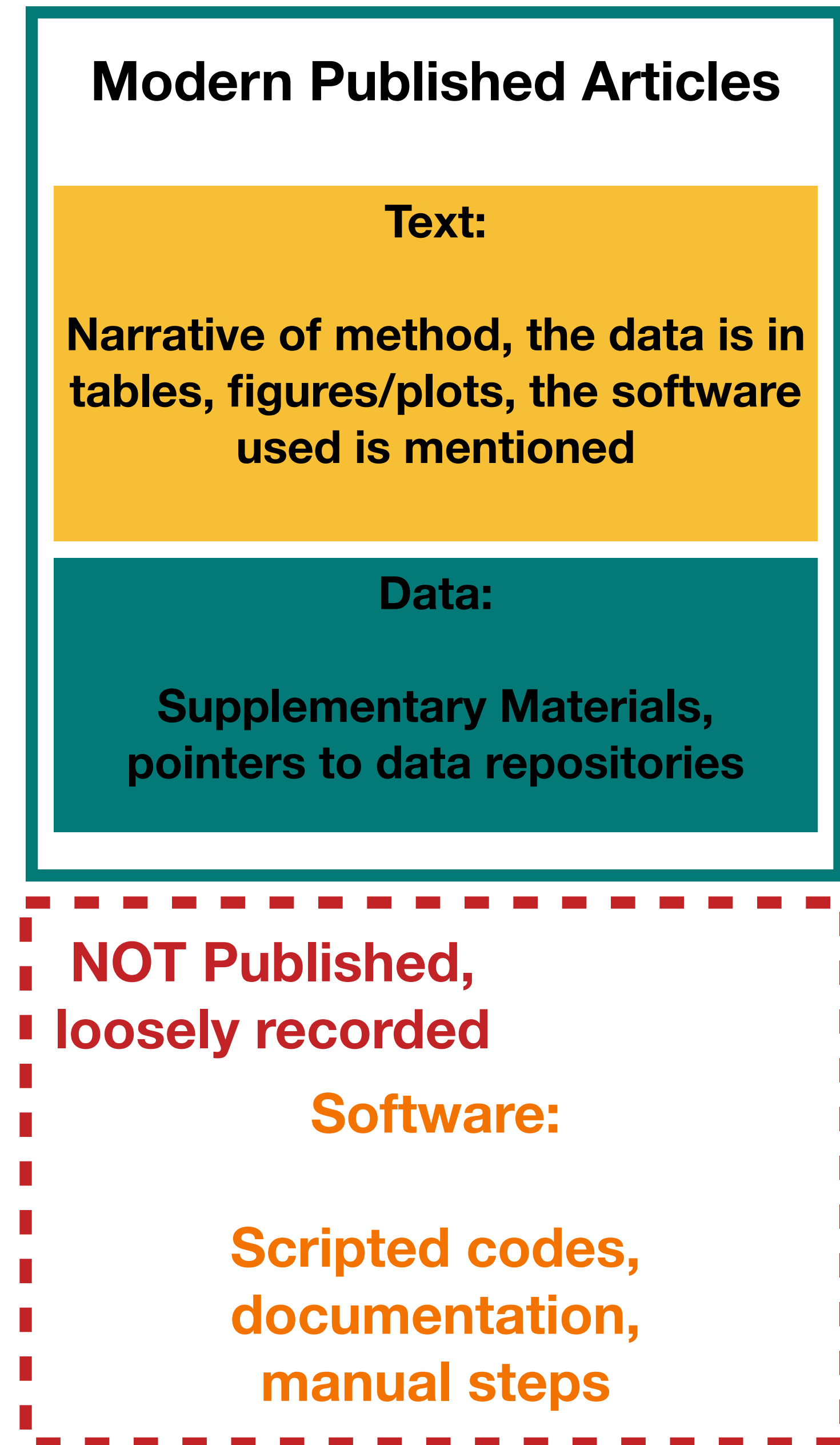
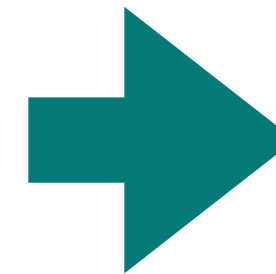
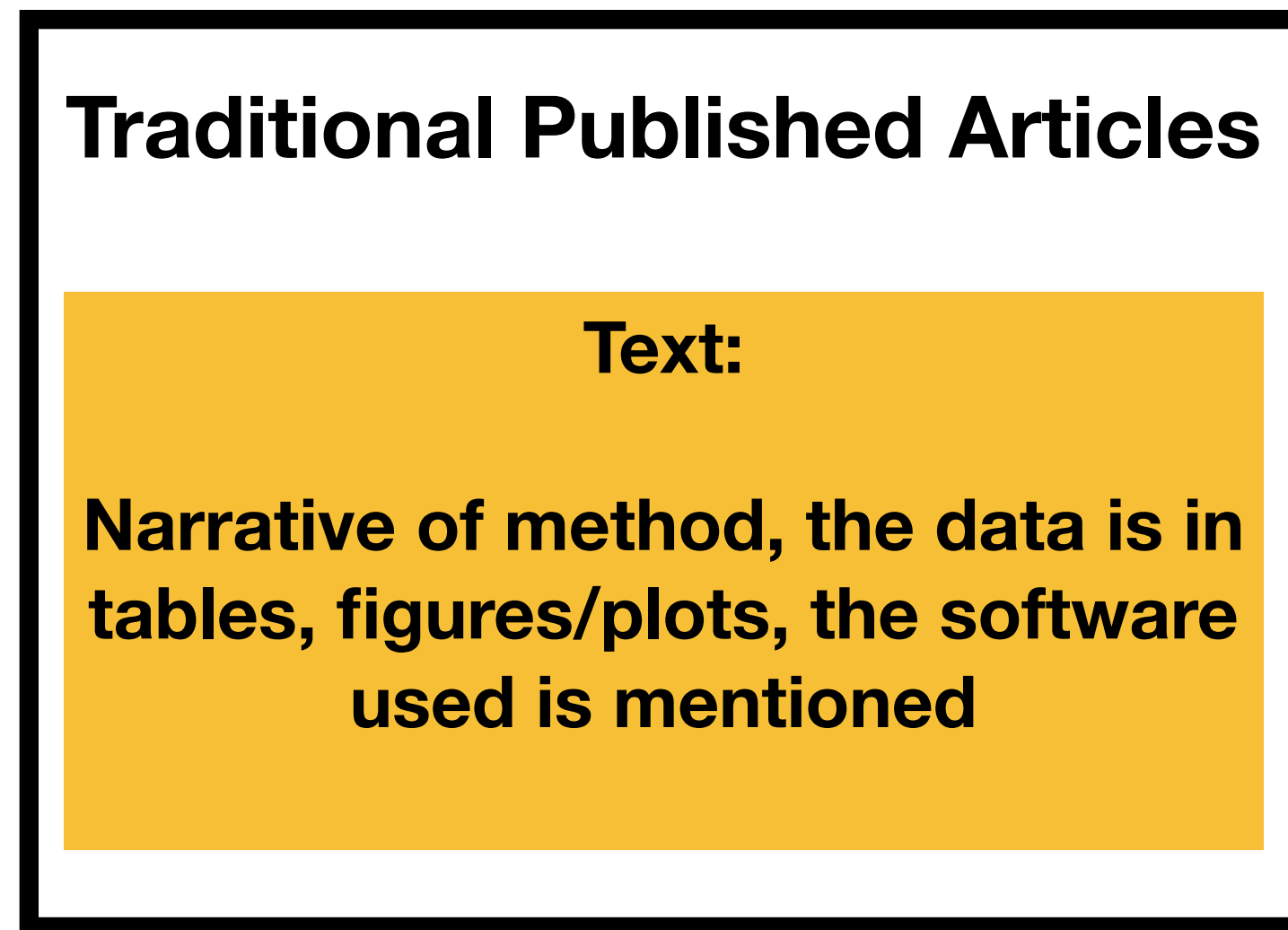
Principal Findings

We examined the citation history of 85 cancer microarray clinical trial publications with respect to the availability of their data. The 48% of trials with publicly available microarray data received 85% of the aggregate citations. Publicly available data was significantly ($p = 0.006$) associated with a 69% increase in citations, independently of journal impact factor, date of publication, and author country of origin using linear regression.

Modern Scientific Articles



Modern Scientific Articles



Scientific Paper of the Future

Modern Paper

Text:

Narrative of the method, some data is in tables, figures/plots, and the software used is mentioned

Data:

Include data as supplementary materials and pointers to data repositories

Reproducible Publication

Software:

For data preparation, data analysis, and visualization

Provenance and methods:

Workflow/scripts specifying dataflow, codes, configuration files, parameter settings, and runtime dependencies

Open Science

Sharing:

Deposit data and software (and provenance/workflow) in publicly shared repositories

Open licenses:

Open source licenses for data and software (and provenance/workflow)

Metadata:

Structured descriptions of the characteristics of data and software (and provenance/workflow)

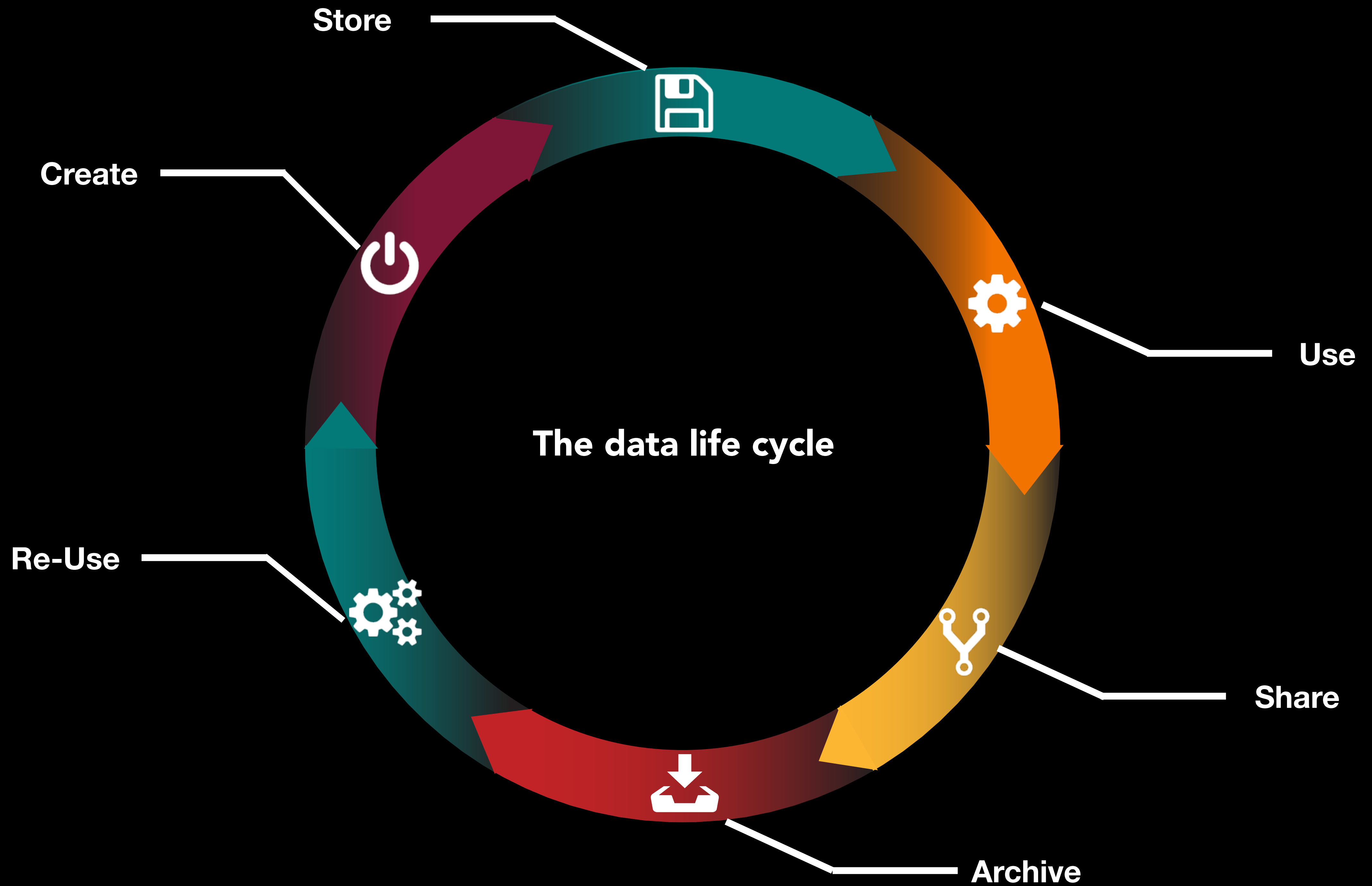
Digital Scholarship

Persistent identifiers:

For data, software, and authors (and provenance/workflow)

Citations:

Citations for data and software (and provenance/workflow)





Store





AND





AND





AND





Store



AND





Store



AND

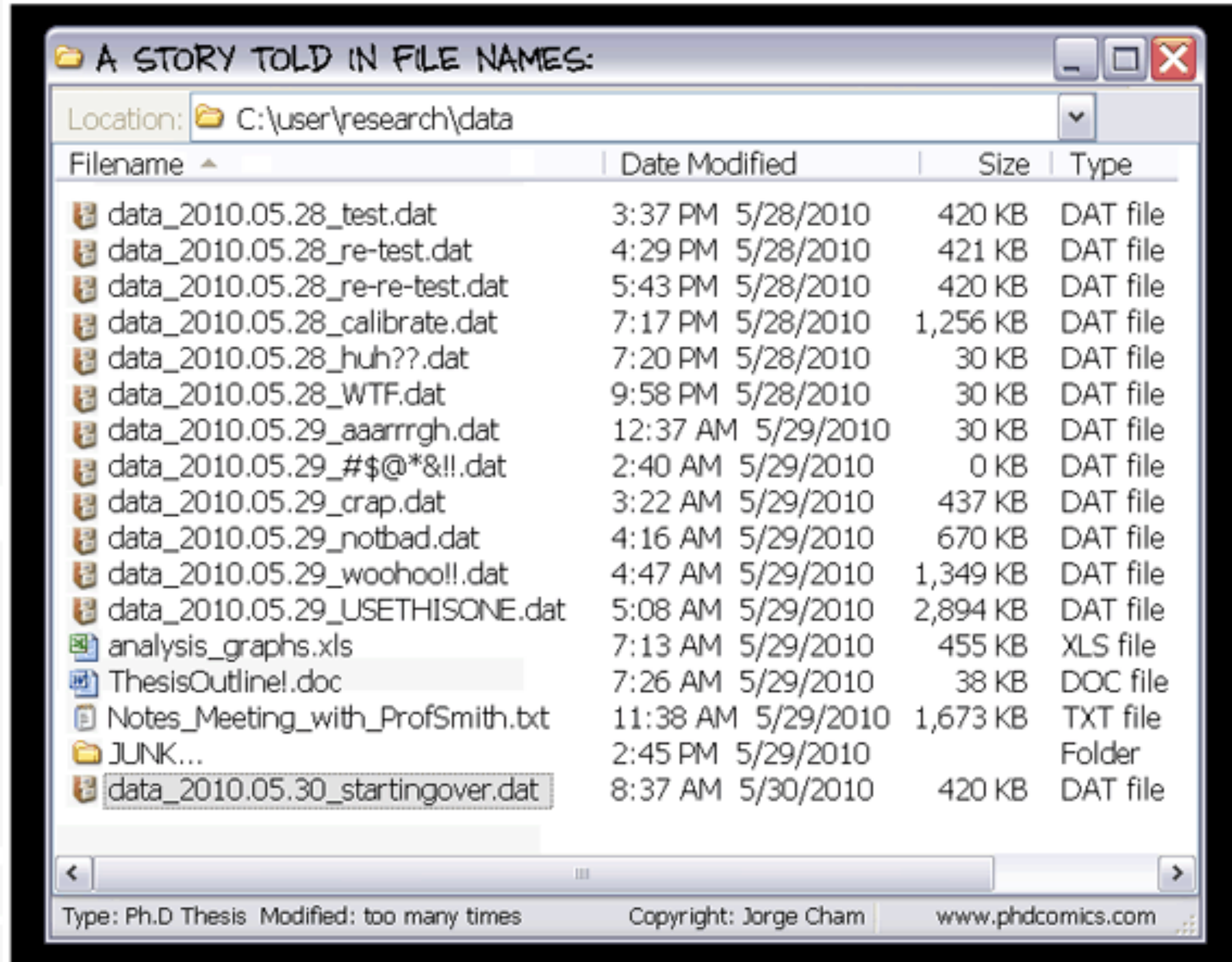


Keep in mind:

1. Some data backup is better than none
2. Automated backups are better than manual
3. Your data is only as safe as the last backup

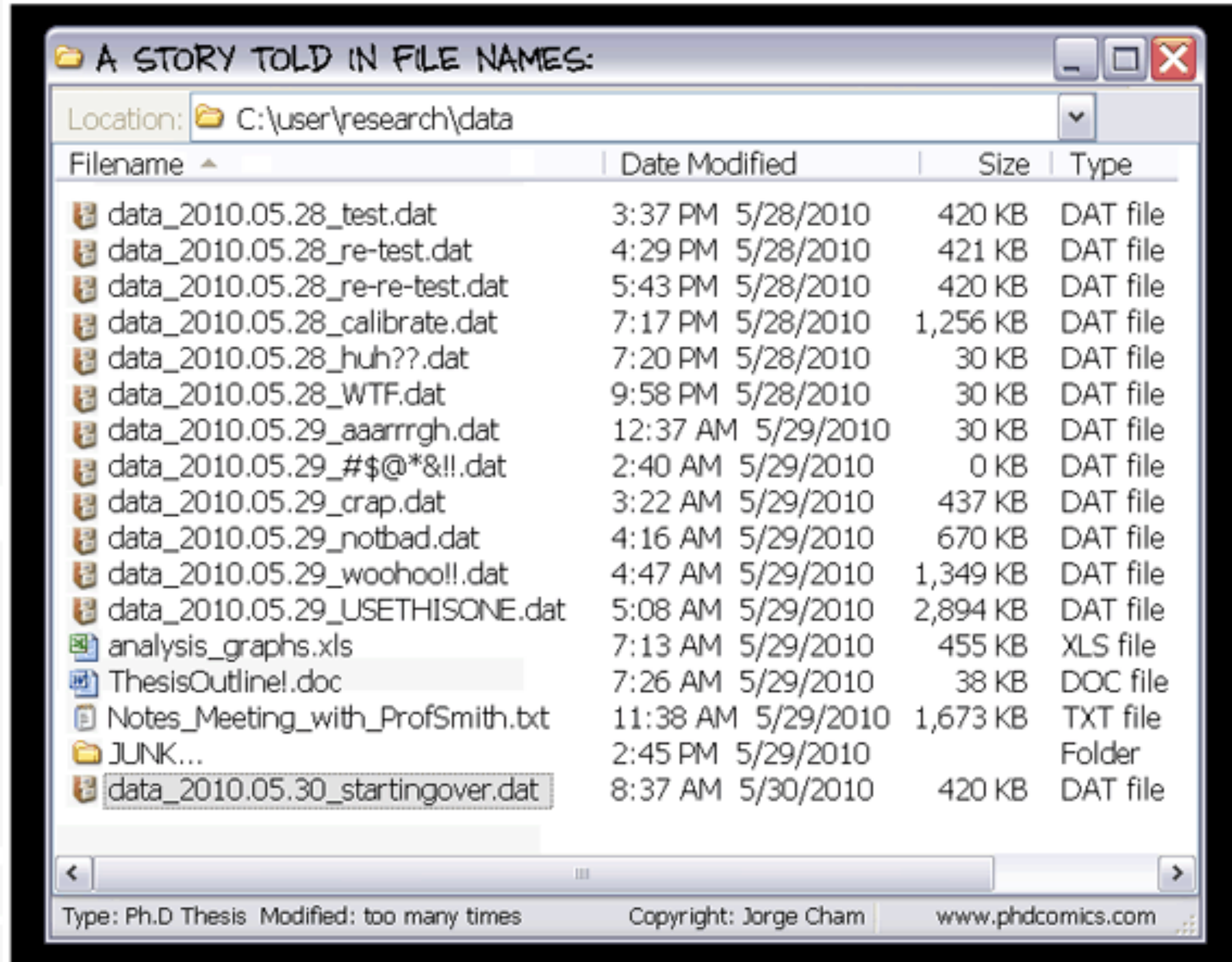


File Name and Organization





File Name and Organization



Morgan Edwards

@mangoedwards

[Follow](#)

I can't send you the original data because I don't remember what my excel file names mean anymore [#overlyhonestmethods](#)

9:11 AM - 8 Jan 2013



130



77



Project_Date_Description

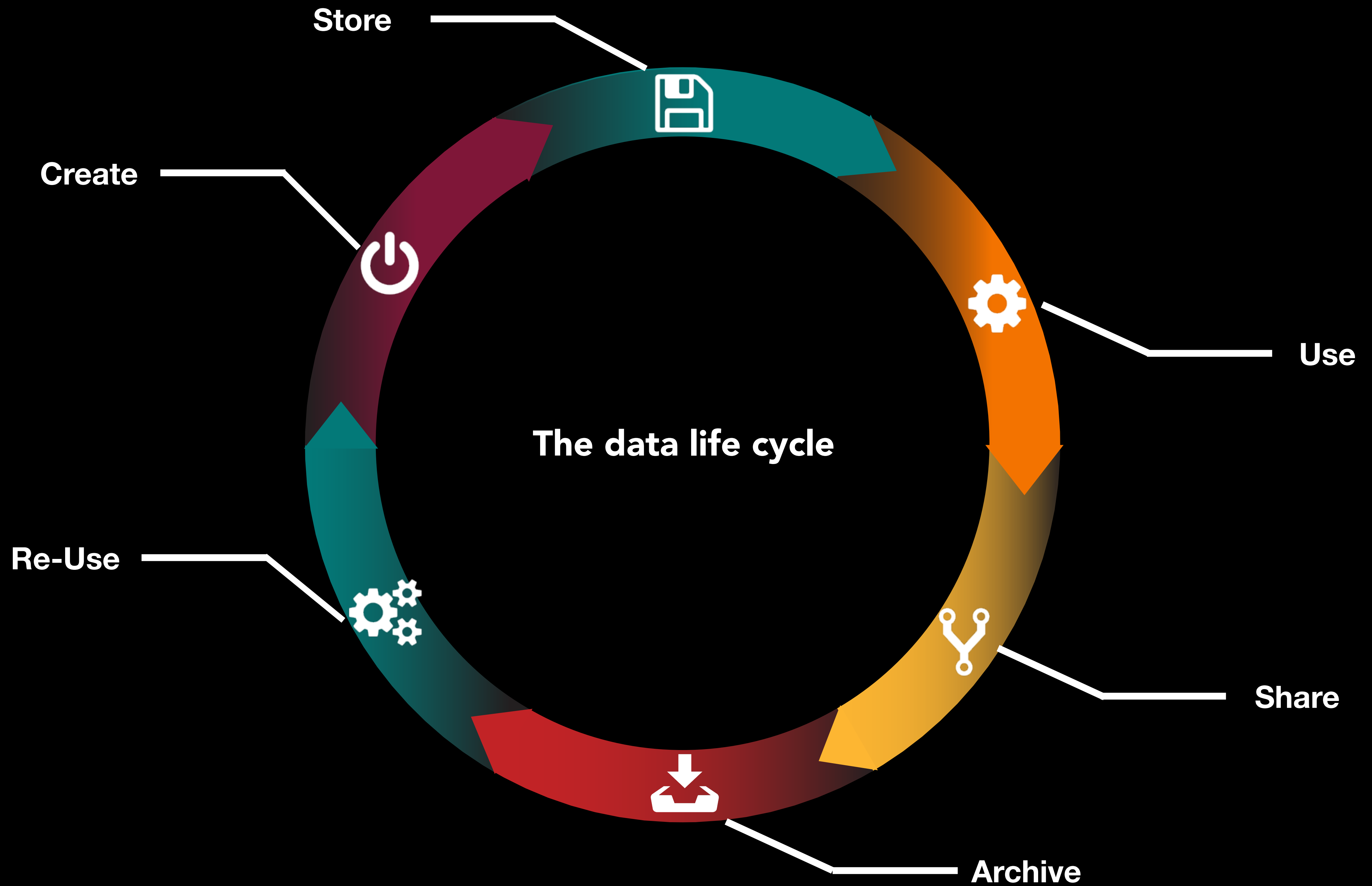


File Name and Organization



KEEP A FILE ABOUT YOUR FILES

► **ReadMe:** Description of what the files/folders contain





KEEP YOUR RAW DATA RAW!

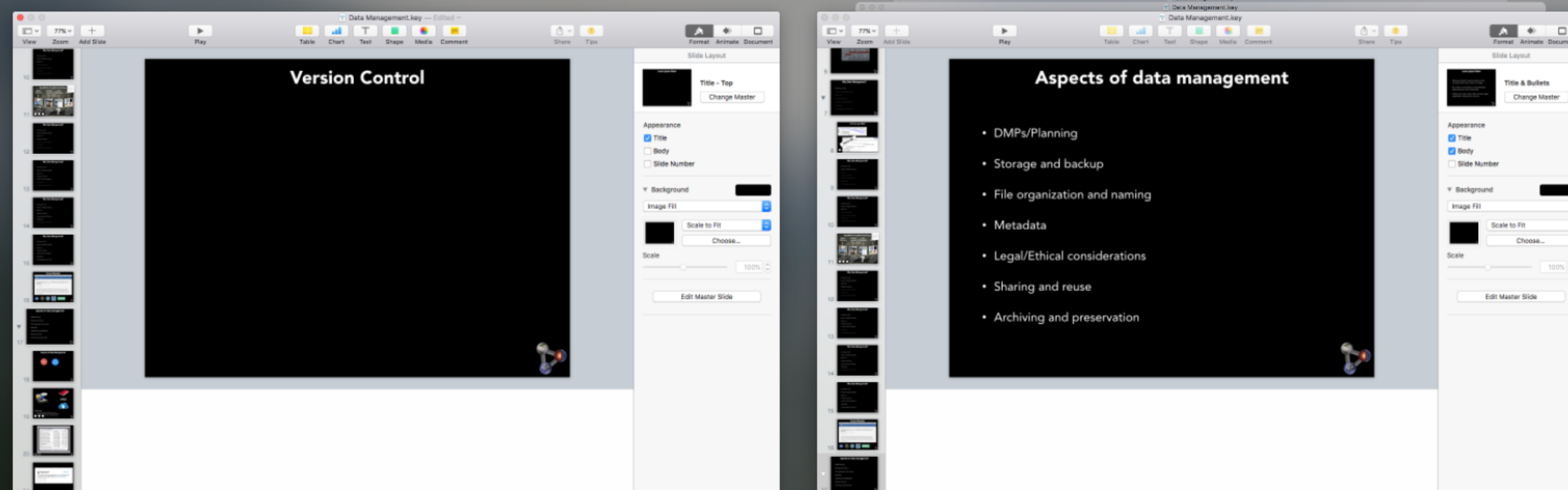
- ▶ Always keep the original data as raw as possible. Create new versions of dataset if you perform any data cleaning.
- ▶ Even more important when calibrating datasets or applying normalization, filters,...



Version Control



Version Control



Current Document

Done

Restore

Today at 2:16 PM

May 2016

Today

Built into Pages, Numbers, Keynote on the Mac...



Version Control

The image shows a Google Docs interface with the 'File' menu open and the 'Revision history' sidebar visible. The document is titled 'Charter'.

File Menu Options:

- Share...
- New
- Open... ⌘O
- Rename...
- Make a copy...
- Organize...
- Move to trash
- See revision history ⌘+Option+Shift+H
- See new changes
- Language
- Download as
- Publish to the web...
- Email collaborators...
- Email as attachment...
- Page setup...
- Print preview
- Print ⌘P

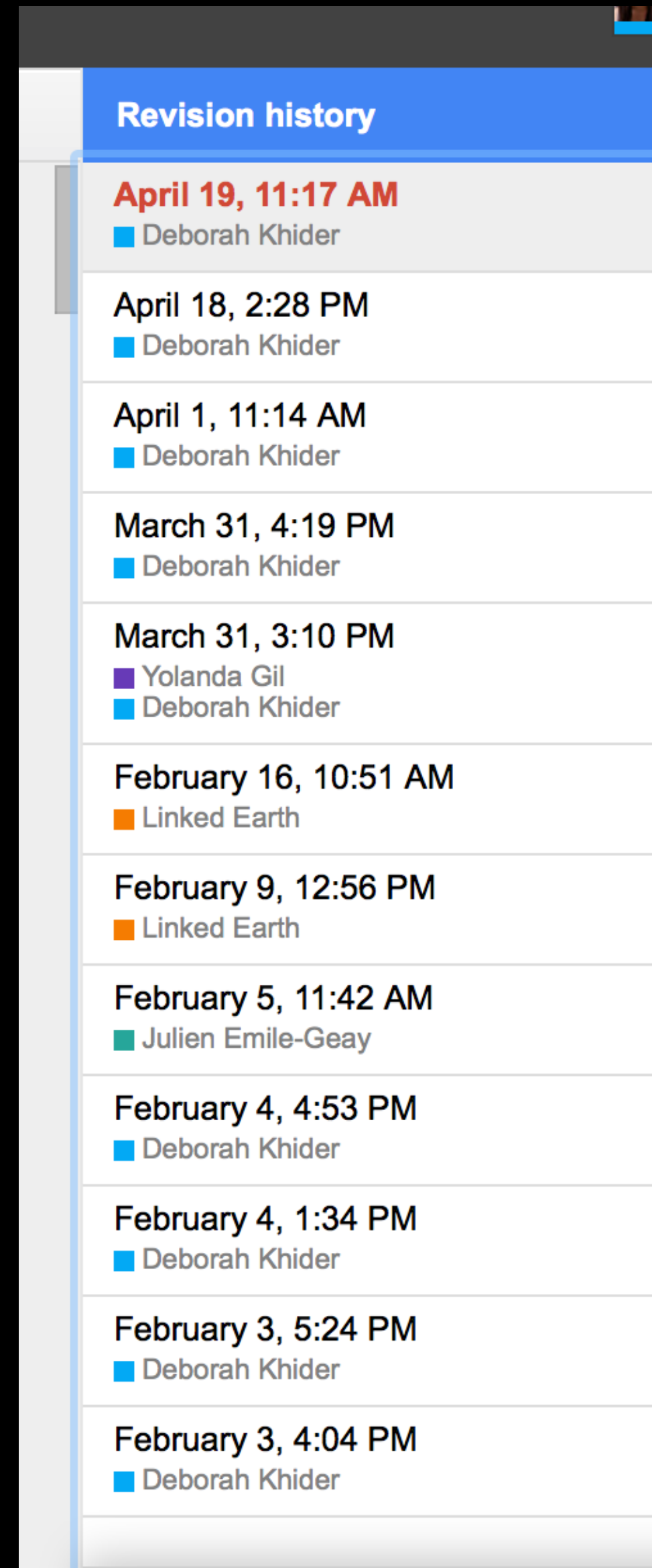
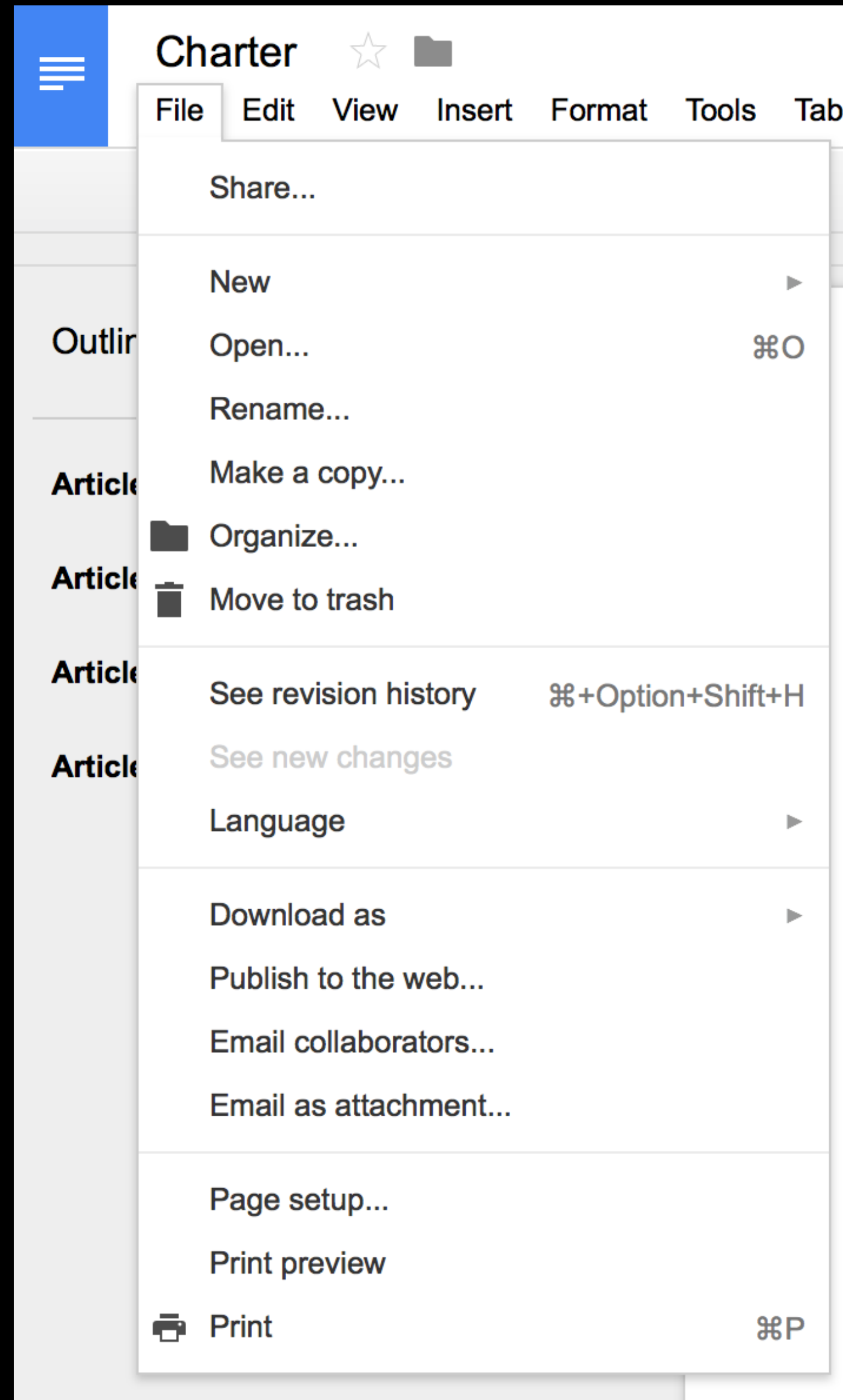
Revision history:

Revision	Author
April 19, 11:17 AM	Deborah Khider
April 18, 2:28 PM	Deborah Khider
April 1, 11:14 AM	Deborah Khider
March 31, 4:19 PM	Deborah Khider
March 31, 3:10 PM	Yolanda Gil, Deborah Khider
February 16, 10:51 AM	Linked Earth
February 9, 12:56 PM	Linked Earth
February 5, 11:42 AM	Julien Emile-Geay
February 4, 4:53 PM	Deborah Khider
February 4, 1:34 PM	Deborah Khider
February 3, 5:24 PM	Deborah Khider
February 3, 4:04 PM	Deborah Khider

... as well as Google apps



Version Control



... as well as Google apps



Version Control

Current Repository
datatransformation_regrid

Current Branch
main

Fetch origin
Last fetched just now

Changes 7

History New

acf1c2a4-25bd-460b-8d48-e738b907fbe8.yaml

New +

7 changed files

acf1c2a4-25bd-460b-8d48-e738b907fbe8.yaml +

ECMWF_regridded.nc +

ModelAnnotatedValues.yml +

RegriddingTransform.cwl +

RegriddingTransform_mic.yaml +

spec.yaml +

values.yml +

Summary (required)

Description

+ Commit to main

@@ -0,0 +1,209 @@

1 +!!python/object:modelcatalog.models.model_configuration.ModelConfiguration

2 +_author: null

3 +_citation: null

4 +_compatible_visualization_software: null

5 +_contributor: null

6 +_copyright_holder: null

7 +_date_created: null

8 +_date_published: null

9 +_description: null

10 +_doi: null

11 +_had_primary_source: null

12 +_has_assumption: null

13 +_has_build_file: null

14 +_has_causal_diagram: null

15 +_has_component_location:

16 +- https://publisher.mint.isi.edu/10PwcE/RegriddingTransform.cwl

17 +_has_constraint: null

18 +_has_contact_person: null

19 +_has_documentation: null

20 +_has_download_instructions: null

21 +_has_download_url: null

22 +_has_equation: null

23 +_has_example: null

24 +_has_executable_instructions: null

25 +_has_executable_notebook: null

26 +_has_execution_command: null

27



Use GitHub (code)

Don't Forget the Metadata!!!

TYPES OF METADATA

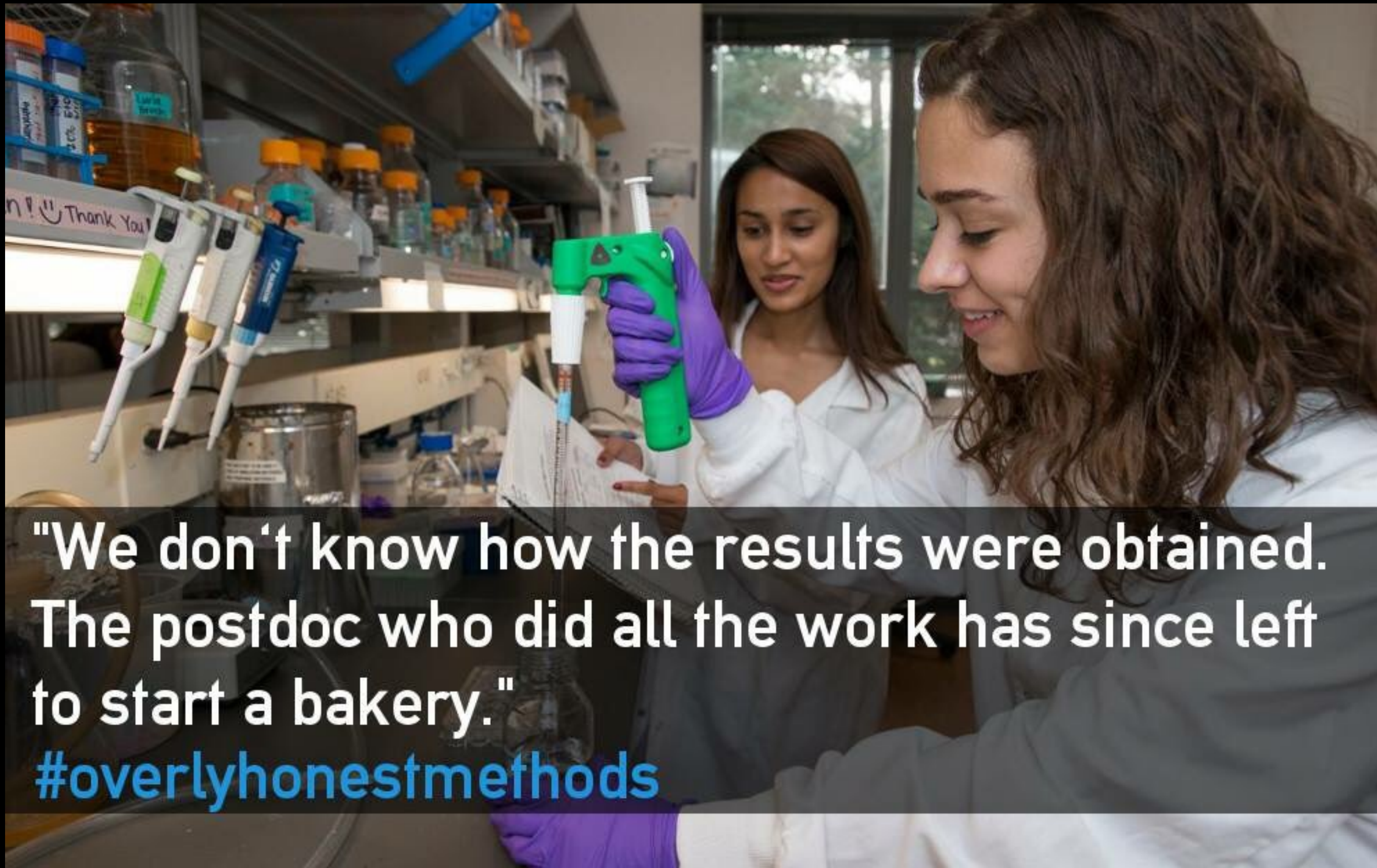
- ▶ **Descriptive Metadata:** Location, collection frequency, object, etc...
- ▶ **Data Characteristics:** Size, statistical properties,...
- ▶ **Provenance metadata:** Instrument, Method/Software, Parameters...



Don't Forget the Metadata!!!

USE OF METADATA

- **Facilitate** reuse by others





Don't Forget the Metadata!!!

USE OF METADATA

- ▶ **Facilitate** reuse by others
- ▶ Support **queries** on data repository

```
In [4]: import json
import requests

url = "http://wiki.linked.earth/store/ds/query"

query = """PREFIX core: <http://linked.earth/ontology#>
PREFIX wiki: <http://wiki.linked.earth/Special:URIResolver/>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
SELECT distinct ?a
WHERE {
{
    ?dataset wiki:Property-3AArchiveType ?a.
}UNION
{
    ?w core:proxyArchiveType ?t.
    ?t rdfs:label ?a
}
}"""

response = requests.post(url, data = {'query': query})
res = json.loads(response.text)

print("The following archive types are available on the wiki:")
for item in res['results']['bindings']:
    print ("*" + item['a']['value'])
```

```
The following archive types are available on the wiki:
*marine sediment
*coral
*lake sediment
*glacier ice
*tree
*documents
*speleothem
*sclerosponge
*borehole
*hybrid
*bivalve
*Rock
```

Don't Forget the Metadata!!!

USE OF METADATA

- ▶ **Facilitate** reuse by others
- ▶ Support **queries** on data repository
- ▶ **Explain a data analysis** by providing context for the data





Don't Forget the Metadata!!!

USE OF METADATA

- ▶ **Facilitate** reuse by others
- ▶ Support **queries** on data repository
- ▶ **Explain a data analysis** by providing context for the data
- ▶ Enable **automated data integration**



Use

THE BIGGEST LIE I TELL
MYSELF IS "I DON'T
NEED TO WRITE THAT
DOWN I'LL REMEMBER."

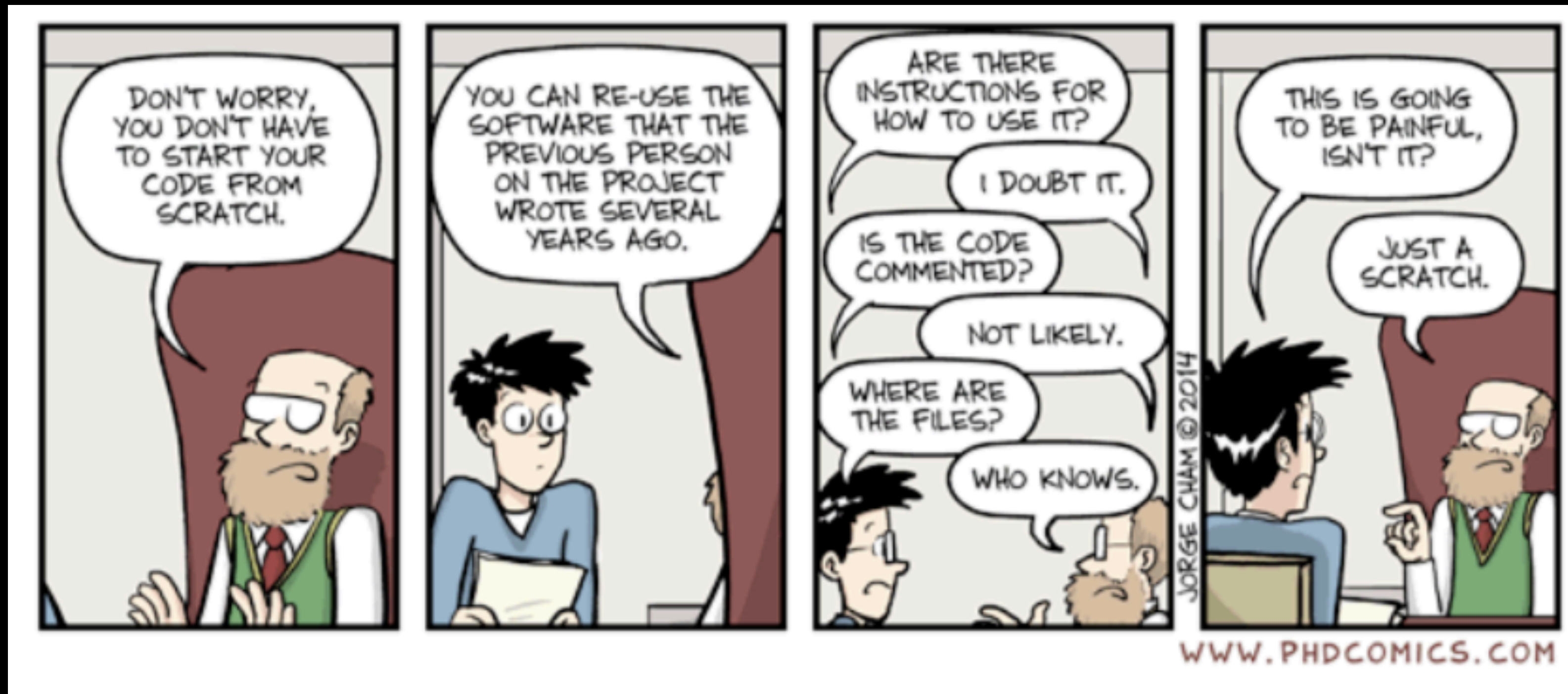
Everyone at some point in their life.



Use

WRITE IT DOWN!

- **Methods:** Laboratory, statistics, data cleaning...
- **Comment your code:** 1 line of code = 1 line of comment
- **Meeting notes**





Keep code and explanation together!



Notebook

jupyter spectrogram (autosaved)



File Edit View Insert Cell Kernel Help

Python 3

Save + Cut Copy Paste Undo Redo Markdown CellToolbar

Simple spectral analysis

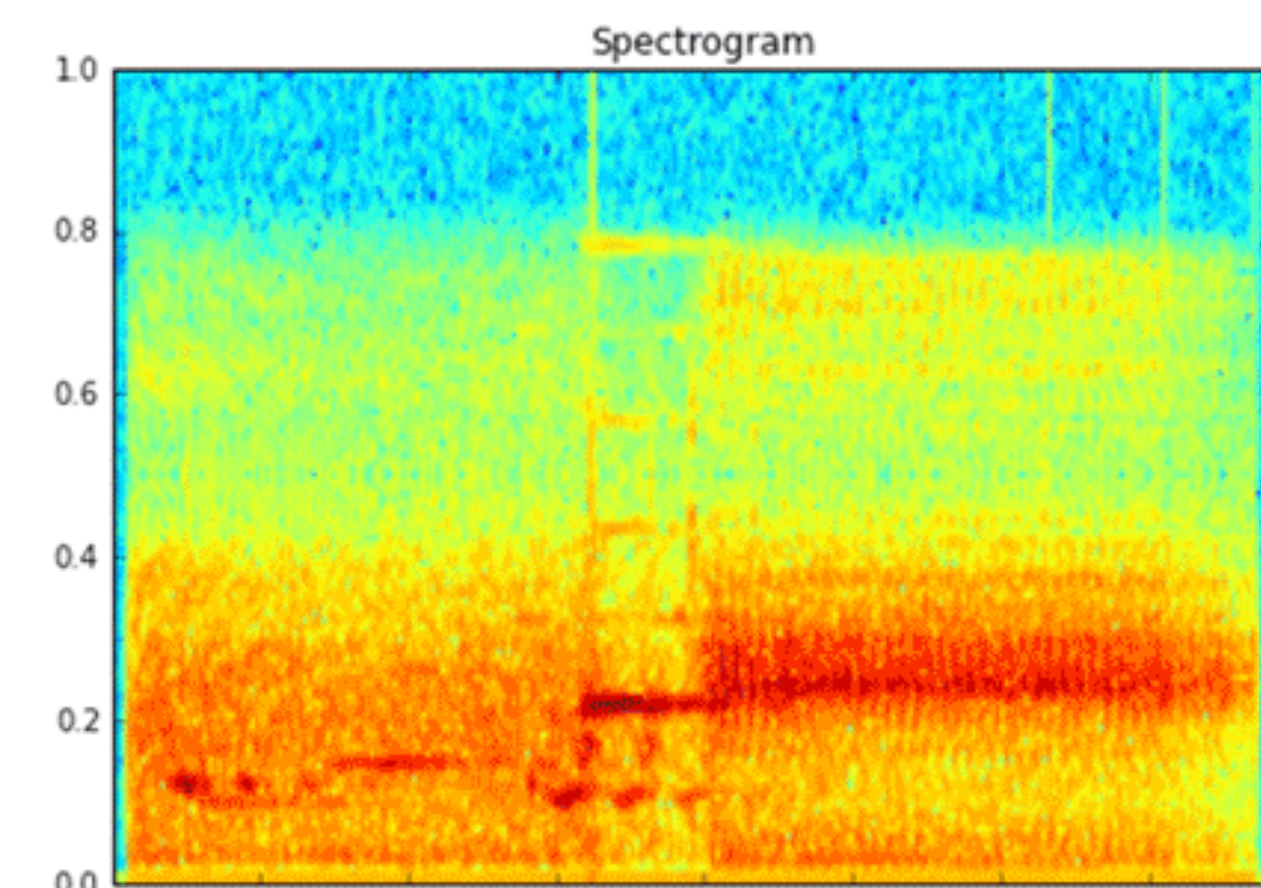
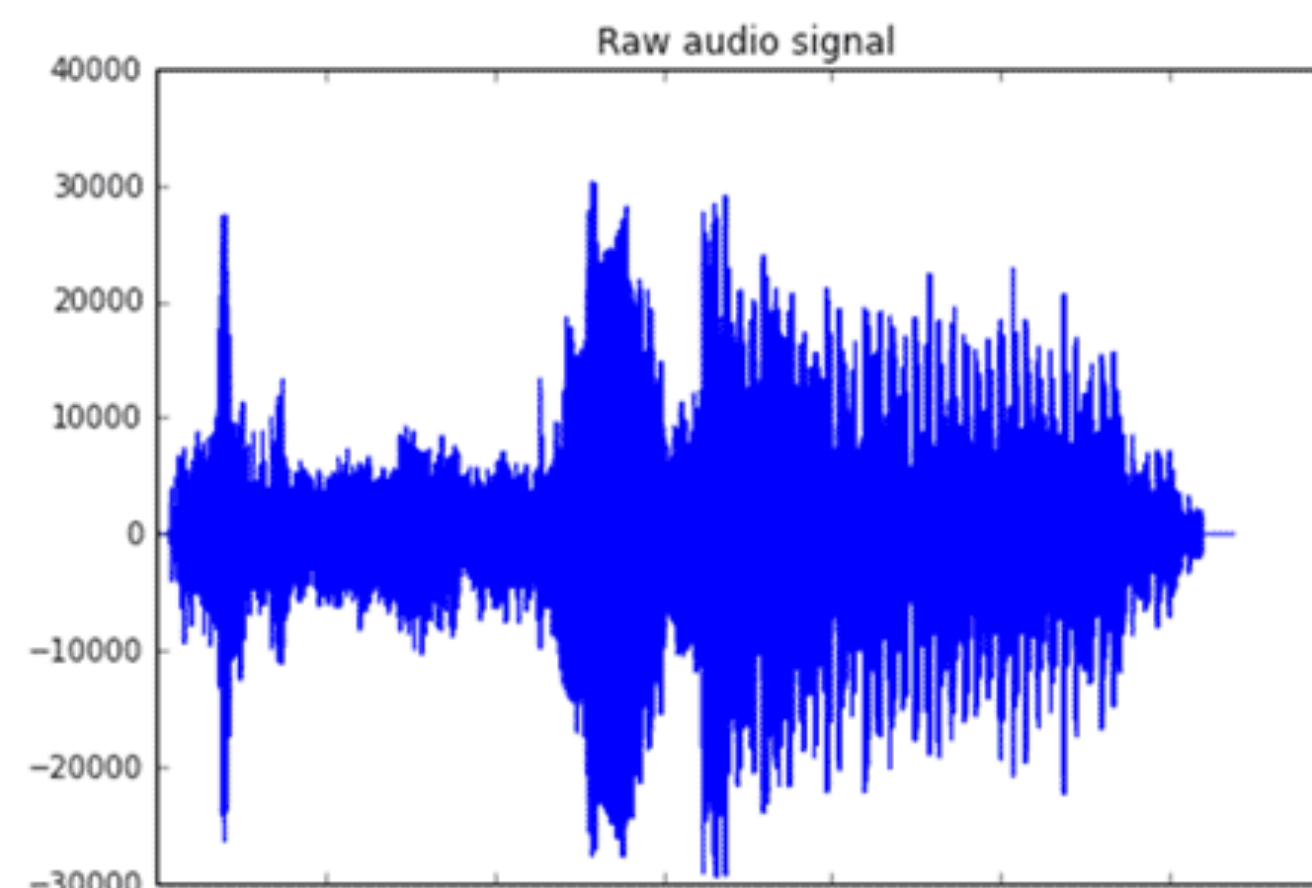
An illustration of the [Discrete Fourier Transform](#)

$$X_k = \sum_{n=0}^{N-1} x_n \exp\left(\frac{-2\pi i}{N} kn\right) \quad k = 0, \dots, N-1$$

```
In [2]: from scipy.io import wavfile
rate, x = wavfile.read('test_mono.wav')
```

And we can easily view it's spectral structure using matplotlib's builtin specgram routine:

```
In [5]: fig, (ax1, ax2) = plt.subplots(1,2,figsize(16,5))
ax1.plot(x); ax1.set_title('Raw audio signal')
ax2.specgram(x); ax2.set_title('Spectrogram');
```





Keep code and explanation together!



R Markdown

from R Studio

```
1 ---  
2 title: "Exploring the BRFSS data"  
3 output:  
4   html_document:  
5     fig_height: 4  
6     highlight: pygments  
7     theme: spacelab  
8 ---  
9  
10 ## Setup  
11  
12 ### Load packages  
13  
14 ```{r load-packages, message = FALSE}  
15 library(ggplot2)  
16 library(dplyr)  
17 ```  
18  
19 ### Load data  
20 Load the BRFSS data into the workspace.  
21  
22 ```{r load-data}  
23 load("brfss2013.RData")  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835  
1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997  
1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051  
2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105  
2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159  
2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193
```




Keep code and explanation together!



R Markdown

from R Studio

```
1 ---
2 title: "Exploring the BRFSS data"
3 output:
4   html_document:
5     fig_height: 4
6     highlight: pygments
7     theme: spacelab
8 ---
9
10 ## Setup
11
12 ### Load packages
13
14 ```{r load-packages, message = FALSE}
15 library(ggplot2)
16 library(dplyr)
17 ```
18
19 ### Load data
20 Load the BRFSS data into the workspace.
21
22 ```{r load-data}
23 load("brfss2013.RData")
```

1:1 # Exploring the BRFSS data R Markdown

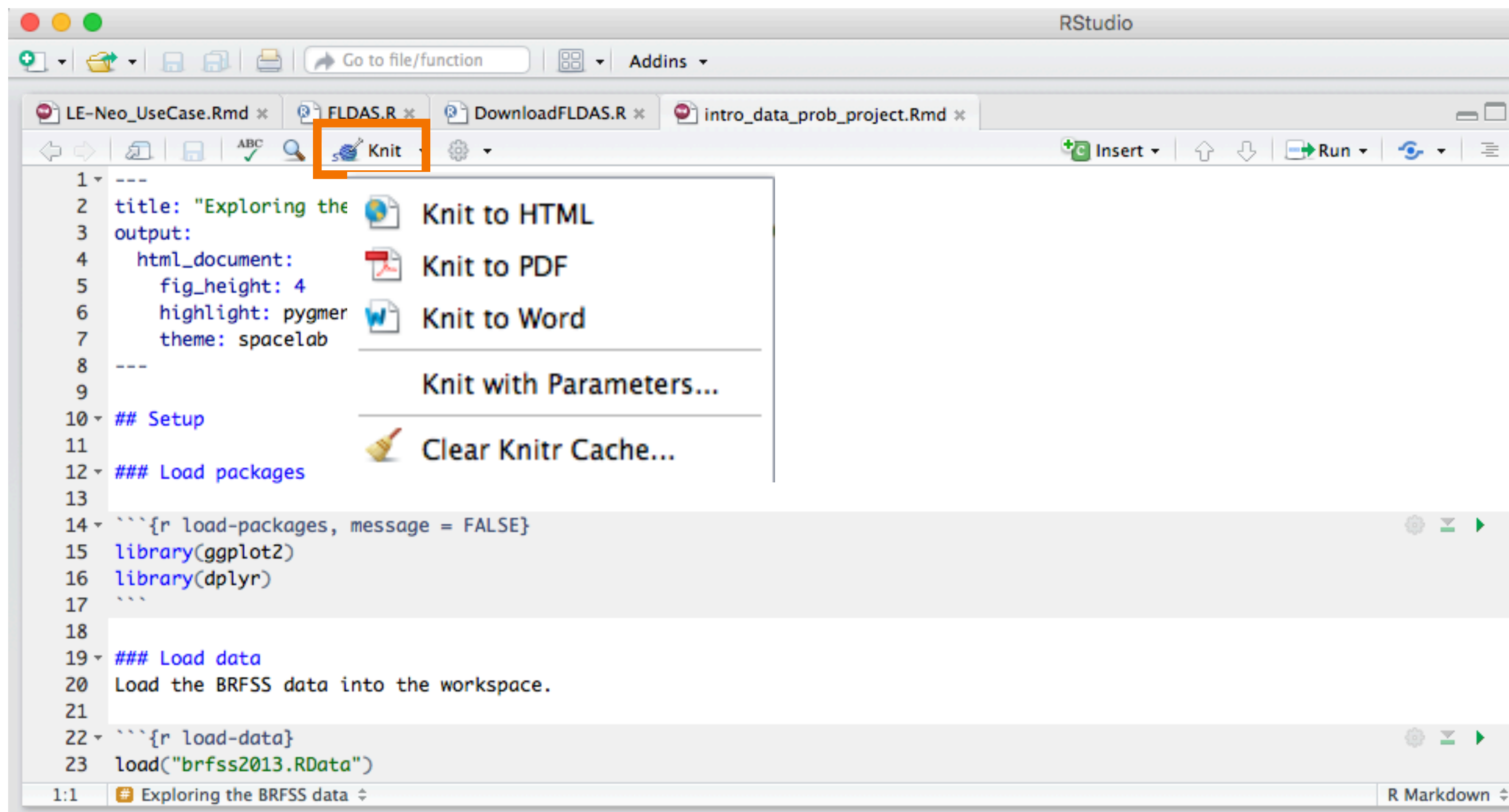


Keep code and explanation together!



R Markdown

from R Studio





Keep code and explanation together!



Notebook

R Markdown

from  Studio

Part 3: Exploratory data analysis

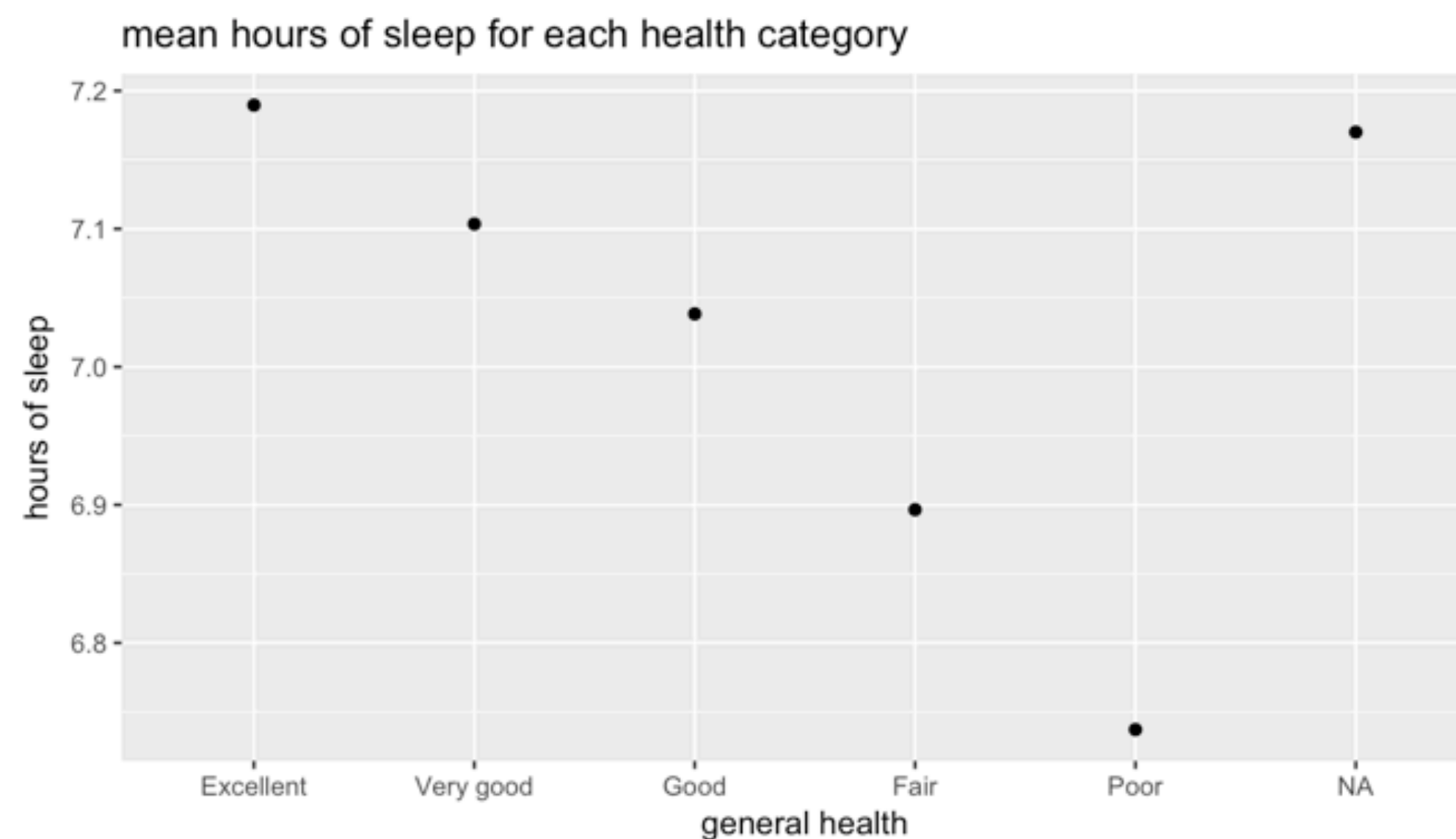
Research question 1: Relationship between sleep and general health

```
#Remove outliers (can't really sleep more than 24hours)
sleeptime <- filter(brfss2013, sleptiml<=24)
```

```
## Warning: package 'bindrcpp' was built under R version 3.2.5
```

```
# Summarize the data
healthSleep <- sleeptime %>%
  group_by(genhlth) %>%
  summarise(meanSleep = mean(sleptiml))

#Plot the general health/mean of sleep time
ggplot(healthSleep, aes(genhlth, meanSleep)) + geom_point(aes(genhlth, meanSleep)) +
  labs(title="mean hours of sleep for each health category", x = 'general health', y='hours of sleep')
```



From the plot above, there appears

to be a general correlation between health and the amount of sleep.



Keep code and explanation together!



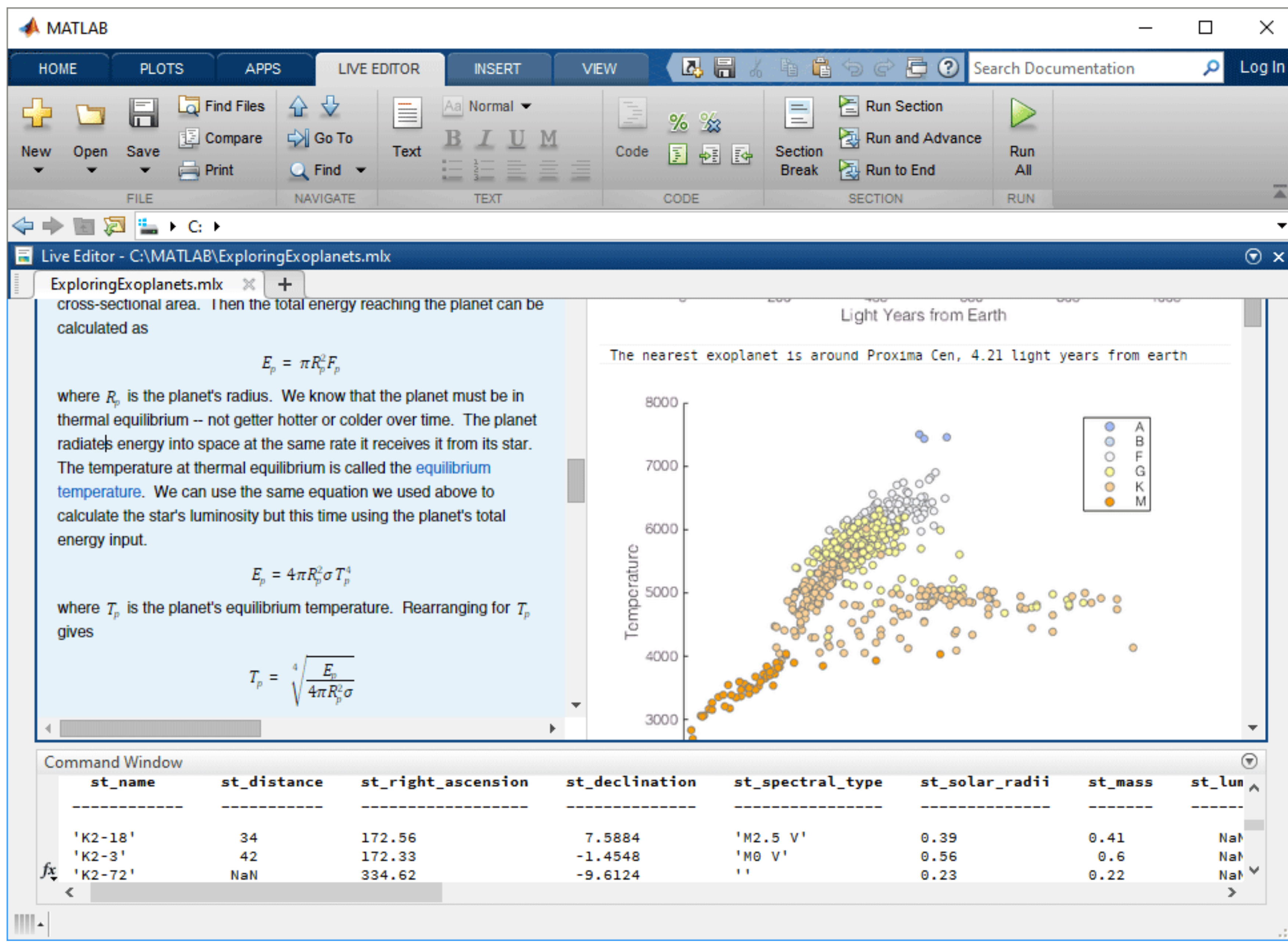
Notebook

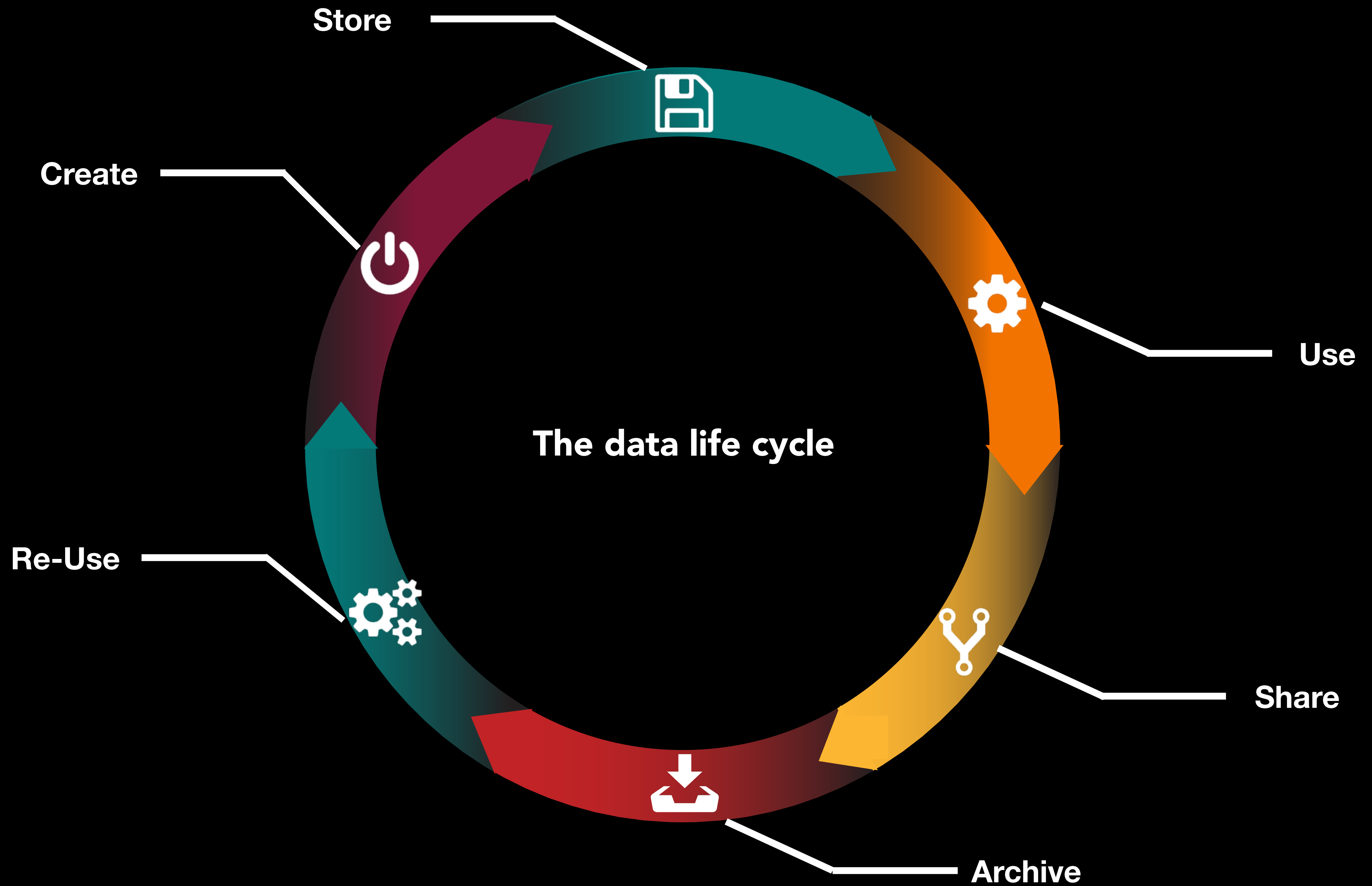
R Markdown

from R Studio



MATLAB Live Editor







INTERNATIONAL GEO SAMPLE NUMBER: IGSN

- ▶ Globally unique and persistent identifier for physical samples in the Earth Sciences
- ▶ To obtain a number, go to <http://www.geosamples.org/>
- ▶ Record and register quality metadata for your samples
 - ▶ At a minimum: Location, contact, access restrictions, lithology
- ▶ Use IGSNs in your publications: text, data tables,...

IGSN: DSR00050U



IGSN:	DSR00050U
Sample Name:	VM28-328A
Other Name(s):	
Sample Type:	Core
Parent IGSN:	Not Provided



Share: Data

1. Publication in a shared repository

figshare.com/

Highly connected drug file

Tretinoin	257	46	Rv1155, aroG, Rv1264, mscL, thyX, gmk, glnA1, Rv1264, thyX, glnA1, trpD, leuA, blaI, ethR, Rv0223c, lipJ, Rv1264, ephG, blaI, ethR, sigC, b
Levothyroxine	173	36	25 cyp130, Rv1264, lppX, gpml, ligA, nirA
Methotrexate	156	32	TB31.7, Rv1264, mscL, lppX, coaA, pcaA, Rv3676, t
4-Hydroxytamoxifen	115	115	fabG1,
Estradiol	98	20	mmaA4, bphD, Rv1264, mscL, thyX, lppX, mmaA2, ptl
Amantadine	79	1	TB31.7, cyp130, aroG, Rv1264, secA1, trpD, nirA
Rifampin	78	13	pth, ethR, clpP, glbN, inhA,
Raloxifene	75	18	pknD, lipJ, fabH, Rv1941, Rv3361c, Rv1264, lppX
Propofol	54	5	mmaA4, Rv1264, groEL, lppX, secA1, glmU, nusA, R
Indinavir	51	14	mmaA4, Rv1264, thyX, lppX, secA1, serA1, Rv3529c
Penicillamine	44	10	pepD, Rv1264, thyX, ethR, trxB2,
Daunorubicin	44	12	pknD, pepD, fabH, Rv1941, devB, ppp, ftsZ, cyp12
Triclosan	42	5	
Darunavir	40	15	

[Enlarge to see the rest of the document](#)

[Enlarge](#) [Download](#)

Cite this: Garijo, Daniel; Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li; Kinnings, Sarah; Bourne, Phil (2013): Highly connected drug file. figshare.
<http://dx.doi.org/10.6084/m9.figshare.776887>
Retrieved 08:56, Feb 20, 2015 (GMT)

Description

Highly connected drug file obtained as a result of the TB-Drugome Workflow.

Links

- <http://purl.org/net/tb-drugome-run>

Categories

- Computational Biology

Authors


Daniel Garijo
Lei Xie
Yinliang Zhang
Yolanda Gil
Li Xie
Sarah Kinnings
Phil Bourne

Tags

- results
- tb-drugome

License (what's this?)

CC-BY





Share: Data



Open
Core Data



HYDROSHARE

Magic



Neotoma




OpenTopography

CSDCO

BALTO



Share: Data

 [figshare.com/](#)

Highly connected drug file

Tretinoin	257	46	Rv1155, aroG, Rv1264, mscL, thyX, gmk, glnA1, Rv1264, thyX, glnA1, trpD, leuA, blaI, ethR, Rv0223c, lipJ, Rv1264, ephG, blaI, ethR, sigC, b
Levothyroxine	173	36	25 cyp130, Rv1264, lppX, gpml, ligA, nirA
Methotrexate	156	32	TB31.7, Rv1264, mscL, lppX, coaA, pcaA, Rv3676, fabG1,
4-Hydroxytamoxifen	115	115	mmaA4, bphD, Rv1264, mscL, thyX, lppX, mmaA2, ptl
Estradiol	98	20	TB31.7, cyp130, aroG, Rv1264, secA1, trpD, nirA
Amantadine	79	1	pth, ethR, clpP, glbN, inhA,
Rifampin	78	13	pknD, lipJ, fabH, Rv1941, Rv3361c, Rv1264, lppX
Raloxifene	75	18	mmaA4, Rv1264, groEL, lppX, secA1, glmU, nusA, R
Propofol	54	5	mmaA4, Rv1264, thyX, lppX, secA1, serA1, Rv3529c
Indinavir	51	14	pepD, Rv1264, thyX, ethR, trxB2,
Penicillamine	44	10	pknD, pepD, fabH, Rv1941, devB, ppp, ftsZ, cyp12
Daunorubicin	44	12	
Triclosan	42	5	
Darunavir	40	15	

[Enlarge to see the rest of the document](#)

[Enlarge](#) [Download](#)

Cite this: Garijo, Daniel; Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li; Kinnings, Sarah; Bourne, Phil (2013): Highly connected drug file. [figshare](#).
<http://dx.doi.org/10.6084/m9.figshare.776887>
Retrieved 08:56, Feb 20, 2015 (GMT)

Description

Highly connected drug file obtained as a result of the TB-Drugome Workflow.

Links

- <http://purl.org/net/tb-drugome-run>

Categories

- Computational Biology

Authors


Daniel Garijo
Lei Xie
Yinliang Zhang
Yolanda Gil
Li Xie
Sarah Kinnings
Phil Bourne

Tags

[results](#) [tb-drugome](#)

License (what's this?)

CC-BY



1. Publication in a shared repository

2. General and domain metadata



Share: Data

GENERAL

- ▶ Dataset name/title
- ▶ Description
- ▶ Creator(s)
- ▶ Publication date
- ▶ License
- ▶ Publisher/contact
- ▶ Version
- ▶ Resource type
- ▶ Location of the data

DOMAIN SPECIFIC

- ▶ Categories
- ▶ Keywords/tags
- ▶ Related Links
- ▶ A data repository in a given discipline may request metadata using accepted standards



Share: Data

GENERAL

- ▶ Dataset name/title
- ▶ Description
- ▶ Creator(s)
- ▶ Publication date
- ▶ License
- ▶ Publisher/contact
- ▶ Version
- ▶ Resource type
- ▶ Location of the data

DOMAIN SPECIFIC

- ▶ Categories
- ▶ Keywords/tags
- ▶ Related Links
- ▶ A data repository in a given discipline may request metadata using accepted standards



Share: Data

Choose a License

Creative Commons Corporation creativecommons.org/choose/

YG WINGS WINGS-Portal ODS DII EC ECC ISD ISI

creative commons About Licenses Public Domain Support CC Projects News

License Features

Your choices on this panel will update the other panels on this page.

Allow adaptations of your work to be shared?

☒ Yes ☐ No



☐ Yes, as long as others share alike

Allow commercial uses of your work?


☒ Yes ☐ No

Selected License

Attribution 4.0 International

This is a Free Culture License!



Help others attribute you!

This part is optional, but filling it out will add machine-readable metadata to the suggested HTML!

Title of work

Attribute work to name

Attribute work to URL


Source work URL

More permissions URL

Format of work

License mark

Have a web page?



This work is licensed under a [Creative Commons Attribution 4.0 International License](http://creativecommons.org/licenses/by/4.0/).

Copy this code to let your visitors know!

```
<a rel="license"
href="http://creativecommons.org/licenses/by/4.0/">
</a><br />This work is licensed under a <a rel="license"
href="http://creativecommons.org/licenses/by/4.0/">Creativ
a Commons Attribution 4.0 International License</a>
```

☒ Normal Icon ☐ Compact Icon

Recommended: CC-BY and CC0



Attribution CC BY

This license lets others distribute, remix, tweak, and build upon your work, even commercially, as long as they credit you for the original creation. This is the most accommodating of licenses offered. Recommended for maximum dissemination and use of licensed materials.

CC0 (datasets)



CC0 can be particularly important for the sharing of data and databases, since it otherwise may be unclear whether highly factual data and databases are restricted by copyright or other rights.

Databases may contain facts that, in and of themselves, are not protected by copyright law.

CC0 is recommended for data and databases and is used by hundreds of organizations. It is especially recommended for scientific data. Although CC0 doesn't legally require users of the data to cite the source, it does not take away the moral responsibility to give attribution, as is common in scientific research.

<http://creativecommons.org/licenses/>



Share: Data

figshare.com/

Highly connected drug file

Tretinoin	257	46	Rv1155, aroG, Rv1264, mscL, thyX, gmk, glnA1, Rv1264, thyX, glnA1, trpD, leuA, blaI, ethR, Rv0223c, lipJ, Rv1264, ephG, blaI, ethR, sigC, b
Levothyroxine	173	36	25 cyp130, Rv1264, lppX, gpml, ligA, nirA
Methotrexate	156	32	TB31.7, Rv1264, mscL, lppX, coaA, pcaA, Rv3676, fabG1,
4-Hydroxytamoxifen	115	115	mmaA4, bphD, Rv1264, mscL, thyX, lppX, mmaA2, ptl
Estradiol	98	20	TB31.7, cyp130, aroG, Rv1264, secA1, trpD, nirA
Amantadine	79	1	pth, ethR, clpP, glbN, inhA,
Rifampin	78	13	pknD, lipJ, fabH, Rv1941, Rv3361c, Rv1264, lppX
Raloxifene	75	18	mmaA4, Rv1264, groEL, lppX, secA1, glmU, nusA, R
Propofol	54	5	mmaA4, Rv1264, thyX, lppX, secA1, serA1, Rv3529c
Indinavir	51	14	pepD, Rv1264, thyX, ethR, trxB2,
Penicillamine	44	10	pknD, pepD, fabH, Rv1941, devB, ppp, ftsZ, cyp12!
Daunorubicin	44	12	
Triclosan	42	5	
Darunavir	40	15	

[Enlarge to see the rest of the document](#)

[Enlarge](#) [Download](#)

Categories

- Computational Biology

Authors

Daniel Garijo
Lei Xie
Yinliang Zhang
Yolanda Gil
Li Xie
Sarah Kinnings
Phil Bourne

Tags

- results
- tb-drugome

License [\(what's this?\)](#)

CC-BY


Cite this: Garijo, Daniel; Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li; Kinnings, Sarah; Bourne, Phil (2013): Highly connected drug file. figshare. <http://dx.doi.org/10.6084/m9.figshare.776887>
Retrieved 08:56, Feb 20, 2015 (GMT)

Description

Highly connected drug file obtained as a result of the TB-Drugome Workflow.

Links

- <http://purl.org/net/tb-drugome-run>



1. Publication in a shared repository

2. General and domain metadata

3. Accessibility of data (domain and machine)



Share: Data

figshare.com/

Highly connected drug file

Tretinoin	257	46	Rv1155, aroG, Rv1264, mscL, thyX, gmk, glnA1, Rv1
Levothyroxine	173	36	icl, Rv1264, thyX, glnA1, trpD, leuA, blaI, ethR
Methotrexate	156	32	Rv0223c, lipJ, Rv1264, ephG, blaI, ethR, sigC, b
4-Hydroxytamoxifen		115	25 cyp130, Rv1264, lppX, gpml, ligA, nirA
Estradiol	98	20	TB31.7, Rv1264, mscL, lppX, coaA, pcaA, Rv3676, t
Amantadine	79	1	fabG1,
Rifampin	78	13	mmaA4, bphD, Rv1264, mscL, thyX, lppX, mmaA2, ptl
Raloxifene	75	18	TB31.7, cyp130, aroG, Rv1264, secA1, trpD, nirA
Propofol	54	5	pth, ethR, clpP, glbN, inhA,
Indinavir	51	14	pknD, lipJ, fabH, Rv1941, Rv3361c, Rv1264, lppX
Penicillamine	44	10	mmaA4, Rv1264, groEL, lppX, secA1, glmU, nusA, R
Daunorubicin	44	12	mmaA4, Rv1264, thyX, lppX, secA1, serA1, Rv3529c
Triclosan	42	5	pepD, Rv1264, thyX, ethR, trxB2,
Darunavir	40	15	pknD, pepD, fabH, Rv1941, devB, ppp, ftsZ, cyp12!

[Enlarge to see the rest of the document](#)

[Enlarge](#) [Download](#)

[results](#) [tb-drugome](#)

Cite this: Garijo, Daniel; Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li; Kinnings, Sarah; Bourne, Phil (2013): Highly connected drug file. figshare.
<http://dx.doi.org/10.6084/m9.figshare.776887>
Retrieved 08:56, Feb 20, 2015 (GMT)

Description

Highly connected drug file obtained as a result of the TB-Drugome Workflow.

Links

- <http://purl.org/net/tb-drugome-run>

Categories

- Computational Biology


Authors

Daniel Garijo
Lei Xie
Yinliang Zhang
Yolanda Gil
Li Xie
Sarah Kinnings
Phil Bourne

Tags

[results](#) [tb-drugome](#)

License (what's this?)
CC-BY



1. Publication in a shared repository

2. General and domain metadata

3. Accessibility of data (domain and machine)

4. Unique Persistent Identifier



Share: Data

1. Publication in a shared repository

2. General and domain metadata

3. Accessibility of data (domain and machine)

4. Unique Persistent Identifier

4. Citation preference

figshare.com/

Highly connected drug file

Tretinoin	257	46	Rv1155, aroG, Rv1264, mscL, thyX, gmk, glnA1, Rv1264, thyX, glnA1, trpD, leuA, blaI, ethR
Levothyroxine	173	36	Rv0223c, lipJ, Rv1264, ephG, blaI, ethR, sigC, b
Methotrexate	156	32	25 cyp130, Rv1264, lppX, gpml, ligA, nirA
4-Hydroxytamoxifen	115	115	TB31.7, Rv1264, mscL, lppX, coaA, pcaA, Rv3676, fabG1,
Estradiol	98	20	mmaA4, bphD, Rv1264, mscL, thyX, lppX, mmaA2, ptl
Amantadine	79	1	TB31.7, cyp130, aroG, Rv1264, secA1, trpD, nirA
Rifampin	78	13	pth, ethR, clpP, glbN, inhA,
Raloxifene	75	18	pknD, lipJ, fabH, Rv1941, Rv3361c, Rv1264, lppX
Propofol	54	5	mmaA4, Rv1264, groEL, lppX, secA1, glmU, nusA, R
Indinavir	51	14	mmaA4, Rv1264, thyX, lppX, secA1, serA1, Rv3529c
Penicillamine	44	10	pepD, Rv1264, thyX, ethR, trxB2,
Daunorubicin	44	12	pknD, pepD, fabH, Rv1941, devB, ppp, ftsZ, cyp12
Triclosan	42	5	
Darunavir	40	15	

[Enlarge to see the rest of the document](#)

[Enlarge](#) [Download](#)

Categories

- Computational Biology

Authors

Daniel Garijo
Lei Xie
Yinliang Zhang
Yolanda Gil
Li Xie
Sarah Kinnings
Phil Bourne

Tags

results tb-drugome

License (what's this?)
CC-BY

Cite this: Garijo, Daniel; Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li; Kinnings, Sarah; Bourne, Phil (2013): Highly connected drug file. figshare.
<http://dx.doi.org/10.6084/m9.figshare.776887>
Retrieved 08:56, Feb 20, 2015 (GMT)

Description

Highly connected drug file obtained as a result of the TB-Drugome Workflow.

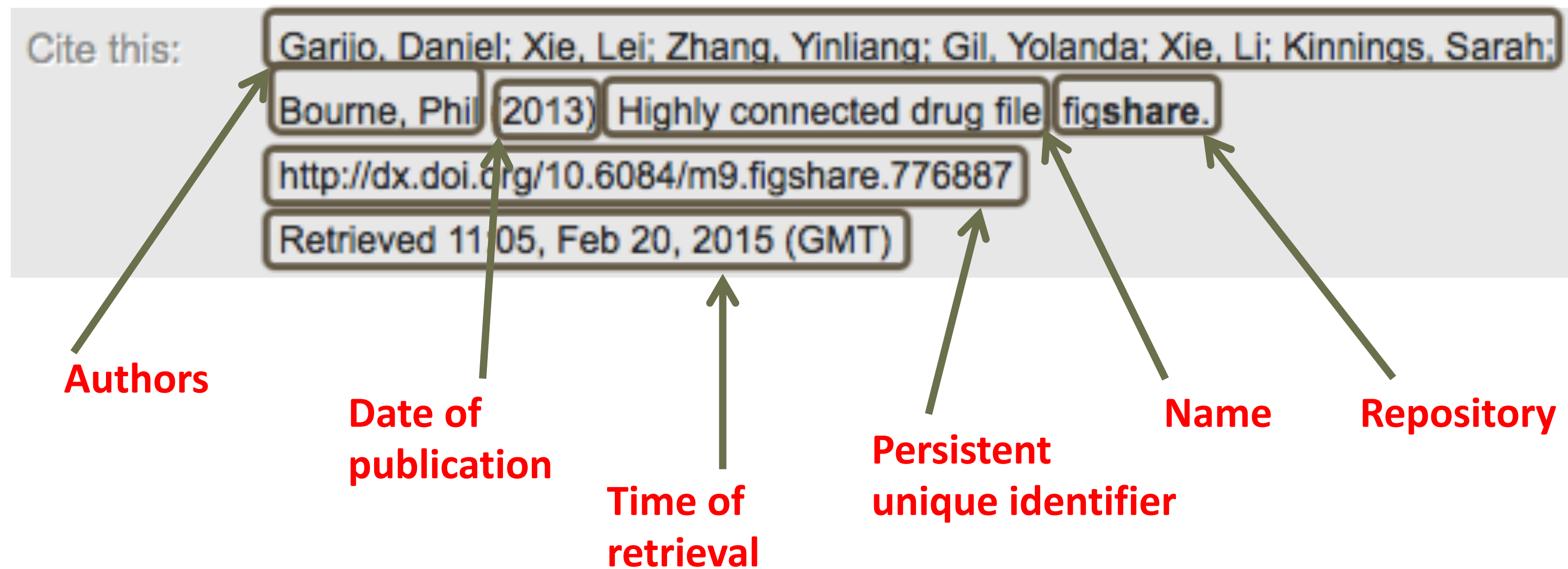
Links

- <http://purl.org/net/tb-drugome-run>



Share: Data

Citing Data: Data repositories and journals often specify how to cite data





WHY IS SCIENTIFIC SOFTWARE NOT SHARED?

- ▶ “No one would use my code if I shared it”
- ▶ “My code is really bad”
- ▶ “My code is not ready to be shared”
- ▶ “Sharing my software will take a lot of time”
- ▶ “I won’t get anything out of sharing my software”
- ▶ “I’ve shared software before, bad things happened”
- ▶ “I work for the government”
- ▶ “I want to commercialize my software”
- ▶ “I don’t want anyone to commercialize my software”
- ▶ “I don’t know where to start”



WHY IS SCIENTIFIC SOFTWARE NOT SHARED?

- ▶ “No one would use my code if I shared it”
- ▶ “My code is really bad”
- ▶ “My code is not ready to be shared”
- ▶ “Sharing my software will take a lot of time”
- ▶ “I won’t get anything out of sharing my software”
- ▶ “I’ve shared software before, bad things happened”
- ▶ “I work for the government”
- ▶ “I want to commercialize my software”
- ▶ “I don’t want anyone to commercialize my software”
- ▶ “I don’t know where to start”





Share: Code

BEST PRACTICES FOR MAKING SOFTWARE AVAILABLE

- ▶ Accessible from a public location
 - ▶ Software repository



Share: Code

BEST PRACTICES FOR MAKING SOFTWARE AVAILABLE

- ▶ Accessible from a public location
 - ▶ Software repository





BEST PRACTICES FOR MAKING SOFTWARE AVAILABLE

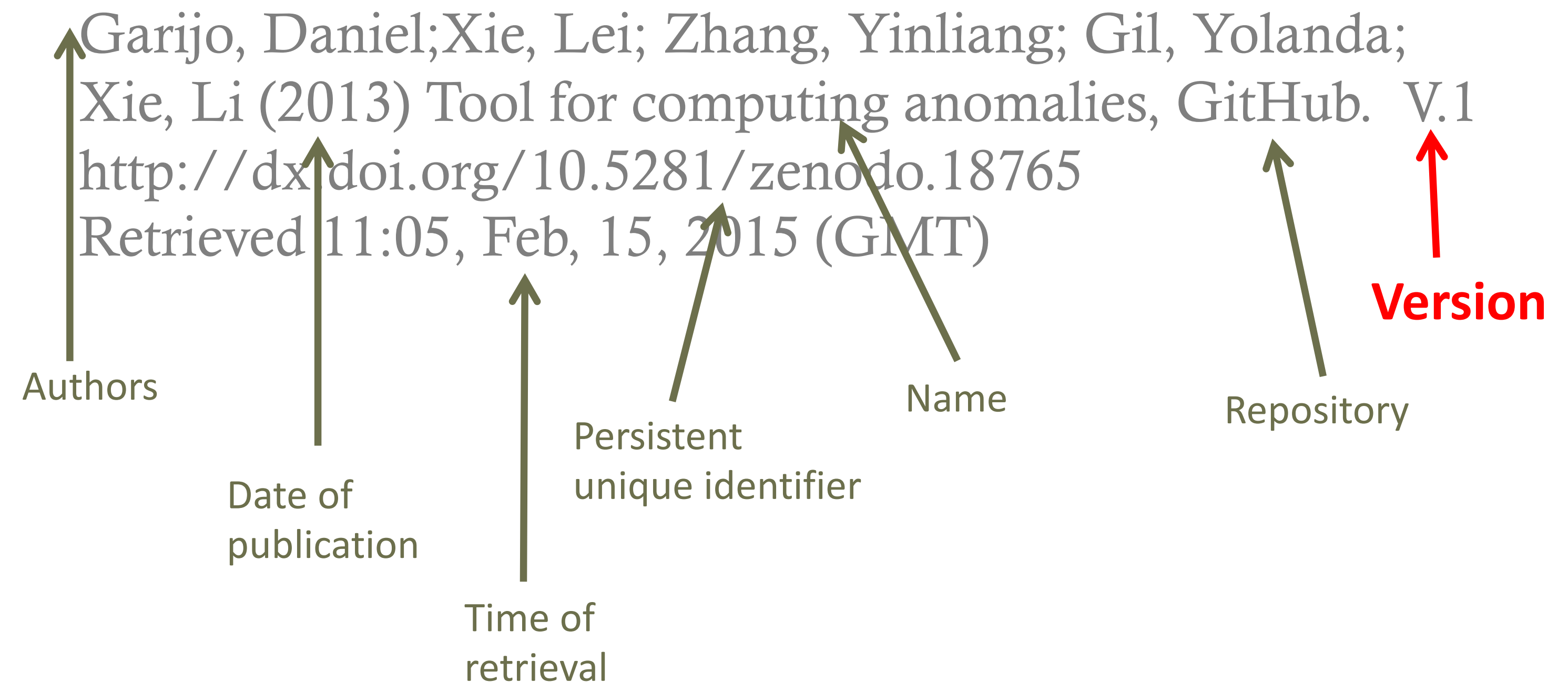
- ▶ Accessible from a public location
- ▶ License
 - ▶ Open source license: reduce constraints and enable software developers to make their source code available to the public
 - ▶ “Copyleft” license (GNU General Public License)
 - ▶ “Permissive” license



Share: Code

BEST PRACTICES FOR MAKING SOFTWARE AVAILABLE

- ▶ Accessible from a public location
- ▶ License
- ▶ Citation



<https://zenodo.org>



Share: Code

BEST PRACTICES FOR MAKING SOFTWARE AVAILABLE



- ▶ Accessible from a public location
- ▶ License
- ▶ Citation
- ▶ Executable via a zero-install environment (in the cloud)

Turn a Git repo into a collection of interactive notebooks

How it works

1

Enter your repository information

Provide in the above form a URL or a GitHub repository that contains Jupyter notebooks, as well as a branch, tag, or commit hash. Launch will build your Binder repository. If you specify a path to a notebook file, the notebook will be opened in your browser after building.

2

We build a Docker image of your repository

Binder will search for a dependency file, such as requirements.txt or environment.yml, in the repository's root directory ([more details on more complex dependencies in documentation](#)). The dependency files will be used to build a Docker image. If an image has already been built for the given repository, it will not be rebuilt. If a new commit has been made, the image will automatically be rebuilt.

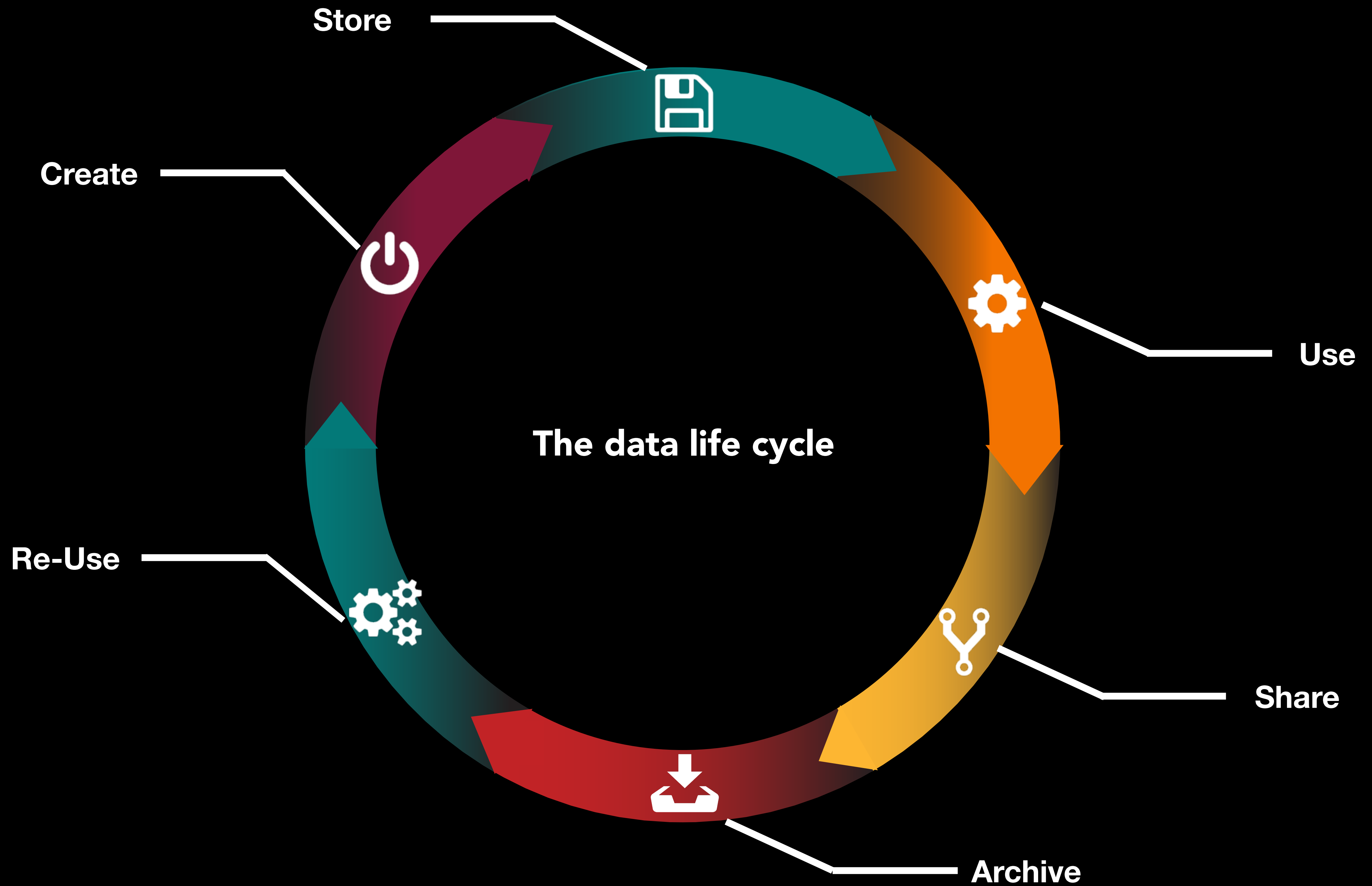
3

Interact with your notebooks in a live environment!

A [JupyterHub](#) server will host your repository's contents. We offer you a reusable link and badge to your live repository that you can easily share with others.

<https://mybinder.org>

<https://github.com/LinkedEarth/paleoHackathon>



References

- Geoscience Paper of the Future:
 - Gil, Y., & . (Ed .) .. (2016, April 17). The Geoscience Paper of the Future: OntoSoft Training (Version 9). figshare. <https://doi.org/10.6084/m9.figshare.1586773.v9>

<http://www.scientificpaperofthefuture.org/gpf/>

Slides Availability

doi: 10.6084/m9.figshare.6510305

<https://figshare.com/s/999787b6f9f6416266b1>

License: CC BY 4.0

