



Best practice and tools for use and re-use [of scientific data in archaeology and heritage]

Dr Scott Allan Orr

Lecturer in Heritage Data Science

SSHOC/SEADDA/E-RIHS Workshop on
Use and Re-Use of Scientific Data in Archaeology and Heritage

2 April 2020

Outline

- Data in heritage and archaeological contexts
- Re-purposing data
- Best practice
- 5-star Open [Linked] Data
- Analysis

Data stewardship versus re-use

- Inherently related
- Data management for **re-use** oriented towards comprehension
 - Accurate
 - Efficient
 - Transparent
- ‘Think about your future self...’

What is specific to scientific data?

- Produced by computer, but interpreted by humans and computers
- Wide range of formats, usually structured
- Can be of staggering scale
- Power is in interconnectivity
- *Analysis*

Data in heritage and archaeology

- Scientific imaging
- Analytical techniques
- Environmental monitoring
- Collection records
- Born-digital assets
- n-D representations

Data in heritage and archaeology

- Metadata
 - Structural, administrative, rights, management
- Paradata
 - Process by which the data was collected
 - Can be integral to informing re-use by understanding intentions, bias, and gaps

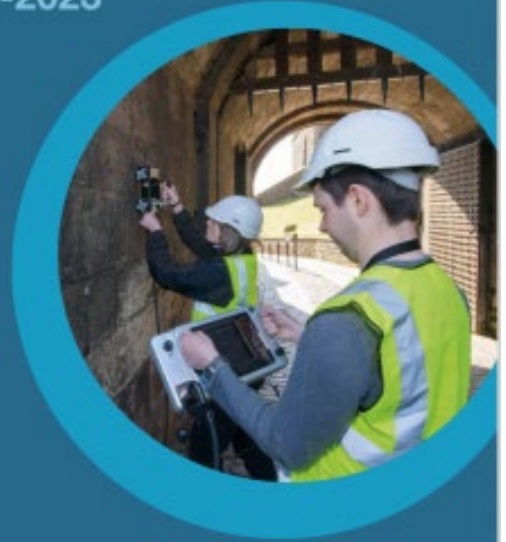


Unidentified Filing Object

Potential re-uses of heritage and archaeological scientific data

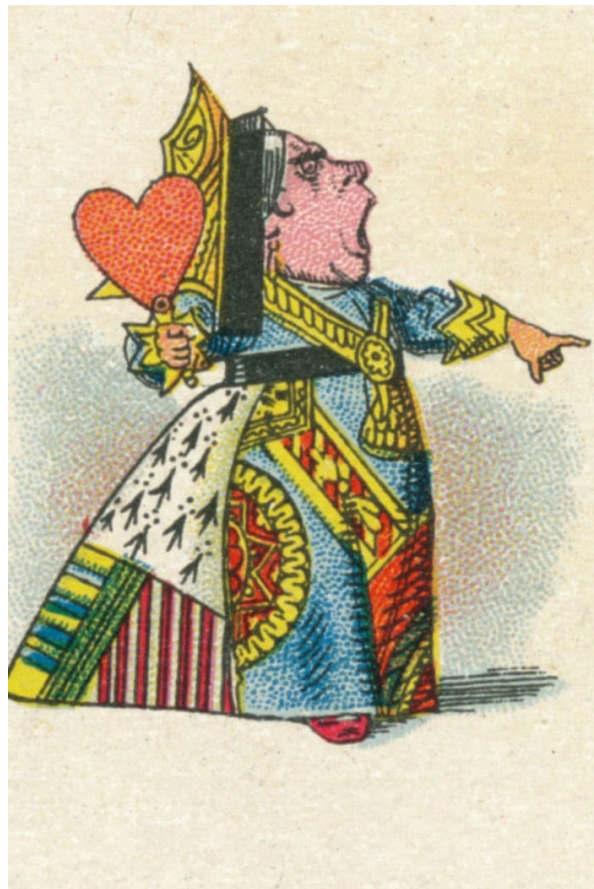
- Conservation and management
- Interpretation
- Engagement

Strategic Framework for Heritage Science in the UK 2018-2023



Examples of best practice

- File formats (proprietary versus plain)
- Headers
- Eternal columns
 - Is there a reason to keep them?
Software, human assumption, etc.



Examples of best practice

- Archive unprocessed data
 - Enables the widest variety of types of re-use

Naming conventions

- Dates: YYYYMMDD
 - Chronological ordering
- Leading zeroes
- Alternatives to spaces
 - file_underscore
 - CamelCase or camelCase
 - file-dashed

0.1-dedication.tex

0.2-acknowledge...

0.3-abstract.tex

1-introduction.tex

10-closing.tex

2-origins.tex

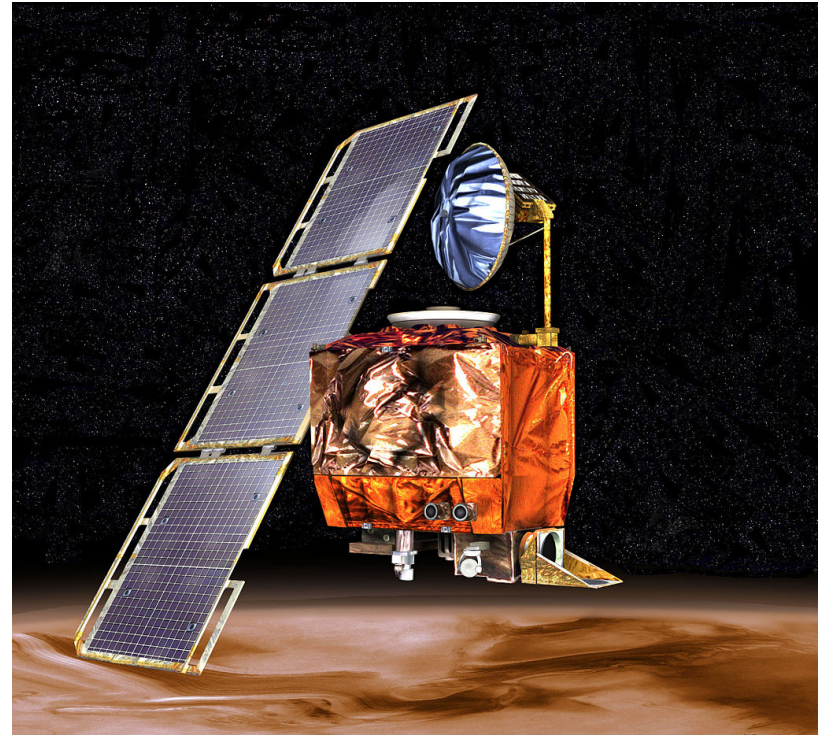
3-propaganda.tex

4-carillonart.tex

5-posthumou... ▼

Examples of best practice

- Floating points and units
 - Cultural differences
 - Disciplinary variation



Examples of best practice: file formats

- *Lossless* file formats

Data type	Lossless	Lossy
Image	RAW, TIFF (usually), PNG (reversible), BMP (proprietary)	JPG, JPEG, GIF (for more than 256 colours)
Audio	WAV	MP3, AAC, OGG
Video	Few, due to size	H.264, H.265, MPEG4

5-star Open Data


- ★ make your stuff available on the Web (whatever format) under an open license¹
- ★★ make it available as structured data (e.g., Excel instead of image scan of a table)²
- ★★★ make it available in a non-proprietary open format (e.g., CSV instead of Excel)³
- ★★★★ use URIs to denote things, so that people can point at your stuff⁴
- ★★★★★ link your data to other data to provide context⁵

XML Data Formats on the Web

- XML implementation for the OGC and ISO Observations and Measurements (O&M) conceptual model (OGC Observations and Measurements v2.0 also published as ISO/DIS 19156)

← → ↻ ⏠ ⓘ Not secure | schemas.opengis.net/om/2.0/examples/

Index of /om/2.0/examples

	<u>Name</u>	<u>Last modified</u>	<u>Size</u>	<u>Description</u>
	Parent Directory			-
	CategoryObservation1.xml	2011-03-22 16:02	1.4K	
	CategoryObservation1a.xml	2011-03-22 16:02	1.2K	
	CategoryObservation1b.xml	2011-03-22 16:02	1.3K	
	CountObservation.xml	2011-03-22 16:02	1.2K	
	DCObservation1.xml	2011-03-22 16:02	2.4K	
	GeometryObservation2shape.xml	2011-03-22 16:02	1.6K	
	SWEArrayObservation1.xml	2011-03-22 16:02	2.4K	
	SWEArrayObservation2.xml	2011-03-22 16:02	4.4K	
	TemporalObservation2.xml	2011-03-22 16:02	1.4K	
	TruthObservation.xml	2011-03-22 16:02	1.1K	
	collection1.xml	2011-03-22 16:02	2.1K	
	collection2.xml	2011-03-22 16:02	2.2K	
	complexObservation3.xml	2011-03-22 16:02	2.4K	
	coverageObservation1.xml	2011-03-22 16:02	4.1K	
	dataObservation4.xml	2011-03-22 16:02	1.4K	
	measurement1.xml	2011-03-22 16:02	1.7K	

Example: True or false?

```
<om:OM_Observation xmlns:om="http://www.opengis.net/om/2.0"
  xmlns:xs="http://www.w3.org/2001/XMLSchema"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xmlns:xlink="http://www.w3.org/1999/xlink"
  xmlns:gml="http://www.opengis.net/gml/3.2" gml:id="obsTest2t"
  xsi:schemaLocation="http://www.opengis.net/om/2.0
    http://schemas.opengis.net/om/2.0/observation.xsd">
  <gml:description>Observation test instance: truth test</gml:description>
  <gml:name>Observation test 2t</gml:name>
  <om:type xlink:href="http://www.opengis.net/def/observationType/OGC-OM/2.0/OM_TruthObservation"/>
  <om:phenomenonTime>
    <gml:TimeInstant gml:id="ot2t">
      <gml:timePosition>2005-01-11T17:22:25.00</gml:timePosition>
    </gml:TimeInstant>
  </om:phenomenonTime>

  <om:resultTime xlink:href="#ot2t"/>
  <om:procedure xlink:href="http://www.example.org/register/party/abc99"/>
  <om:observedProperty xlink:href="urn:example:Truth"/>
  <om:featureOfInterest xlink:href="http://wfs.example.org?request=getFeature featureid=Statement37f"/>
  <om:result xsi:type="xs:boolean">false</om:result>
</om:OM_Observation>
```

Example: True or false?

Product that has a price of \$200” The following would be valid XML:

```
<item>
  <title>AWorkOfArt</title>
  <creator>NameOfArtist</creator>
</item>
```

Another valid XML could be:

```
<item title="AWorkOfArt">
  <creator>NameOfArtist</creator>
</item>
```

Modeling this same data in RDF would only have one way of representing it:

```
ex:item1 rdf:type ex:artwork .
ex:item1 ex:title "AWorkOfArt" .
ex:item1 ex:creator "NameOfArtist" .
```

 URI

Analysis: code, algorithms

- Code and algorithms
- Algorithms are not just code – also procedures
 - Analog and digital
 - Manual and automatic

Code documentation

- Sample use

```
1 # READ.MOIST350B
2 # SCOTT A ORR
3 # 19 MARCH 2017
4
5 # dir <- "/Users/orrscott/Desktop"
6 # filename <- "GraniteWalls_Set5b2.hfs"
7 #
8 # setwd(dir)
9 # filepath <- paste0(dir, "/", filename)
10 # d <- readLines(filepath)
11 #
12 # arrs.t <- hfRead(d)
13
14 ▾ #####
15
16 ▾ # FUNCTIONS #####
17
18 ▸ hfRead <- function(filepath) {←}
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41 ▸ makeArray <- function(line.ref, layer.starts) {←}
42
43
44
45
46
47
48
49
50
51
52 ▸ makeLayer <- function(line.ref, dims, y.vals, x.vals) {←}
53
54
55
56
57
58
59
60
61
62
63
64 ▸ norm <- function(matrix) {←}
65
66
67
68
69
```

Code documentation

- Markup

6 Layer Reflectance Calculator			
n1_re	<input type="text" value="1"/>	n1_im	<input type="text" value="0"/>
n2_re	<input type="text" value="3.50"/>	n2_im	<input type="text" value="-0.35"/>
		t2	<input type="text" value="0.222"/>
n3_re	<input type="text" value="1"/>	n3_im	<input type="text" value="0"/>
		t3	<input type="text" value="0"/>
n4_re	<input type="text" value="1"/>	n4_im	<input type="text" value="0"/>
		t4	<input type="text" value="0"/>
n5_re	<input type="text" value="1"/>	n5_im	<input type="text" value="0"/>
		t5	<input type="text" value="0"/>
n6_re	<input type="text" value="1"/>	n6_im	<input type="text" value="0"/>
		t6	<input type="text" value="0"/>
n7_re	<input type="text" value="1"/>	n7_im	<input type="text" value="0"/>
		t7	<input type="text" value="0"/>
n8_re	<input type="text" value="1.5"/>	n8_im	<input type="text" value="0"/>
θ (deg)	<input type="text" value="15"/>	λ (μm)	<input type="text" value="1.55"/>
Rs	<input type="text"/>	rs	<input type="text"/>
		ts	<input type="text"/>
Rp	<input type="text"/>	rp	<input type="text"/>
		tp	<input type="text"/>
Get Reflectance			
Clear			
Defaults			

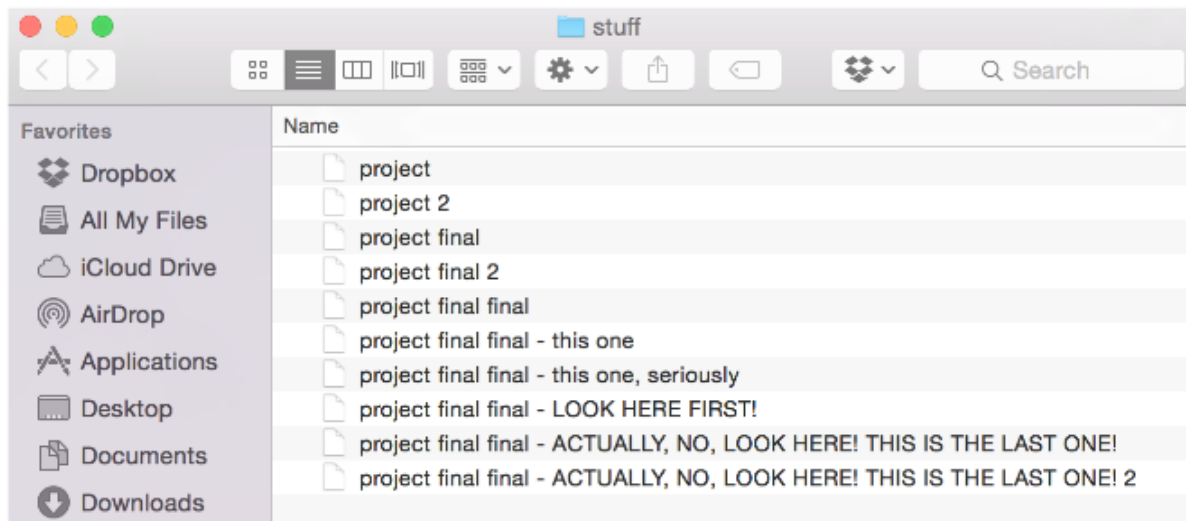
```

187 var Nsqr = [0, 0, 0, 0, 0, 0, 0, 0] ;
188 var Nsqi = [0, 0, 0, 0, 0, 0, 0, 0] ;
189 var ksqr = [0, 0, 0, 0, 0, 0, 0, 0] ; // initialize various arrays for 8 regions
190 var ksqi = [0, 0, 0, 0, 0, 0, 0, 0] ;
191 var kxr = [0, 0, 0, 0, 0, 0, 0, 0] ;
192 var kxi = [0, 0, 0, 0, 0, 0, 0, 0] ;
193
194
195 k0 = 2*Math.PI/lam; // in um-1
196 k0sq = k0*k0;
197 //kz = k0*N1r*Math.sin(theta); // propogation constant along layers (same in each
layer .... basically Snell's law of refraction) .. real for real angle of incidence
198 //kzsq = kz*kz;
199 kz_r = k0*N1r*Math.sin(theta); // kz may be complex if incident medium is complex
200 kz_i = k0*N1i*Math.sin(theta);
201 kzsqr = Csqr(kz_r, kz_i );
202 kzsqi = Csqi(kz_r, kz_i );
203
204
205
206 for(i=0; i<8; i++) {
207   Nsqr[i] = Csqr(Nr[i], Ni[i]) ;
208   Nsqi[i] = Csqi(Nr[i], Ni[i]) ;
209   ksqr[i] = k0sq*Nsqr[i] ;
210   ksqi[i] = k0sq*Nsqi[i] ;
211   kxr[i] = Crootr(ksqr[i] - kzsqr, ksqi[i] - kzsqi);
212   kxi[i] = Crooti(ksqr[i] - kzsqr, ksqi[i] - kzsqi);
213   if (Ni[0] == 0 && kxi[i] > 0) { // if incident medium is lossless, normal
components of propog. const. must all be negative; otherwise select other root
214     kxr[i] = -kxr[i] ;
215     kxi[i] = -kxi[i];
216   }
217 }

```

Version control

- Version control, Git (e.g. Github, etc.)



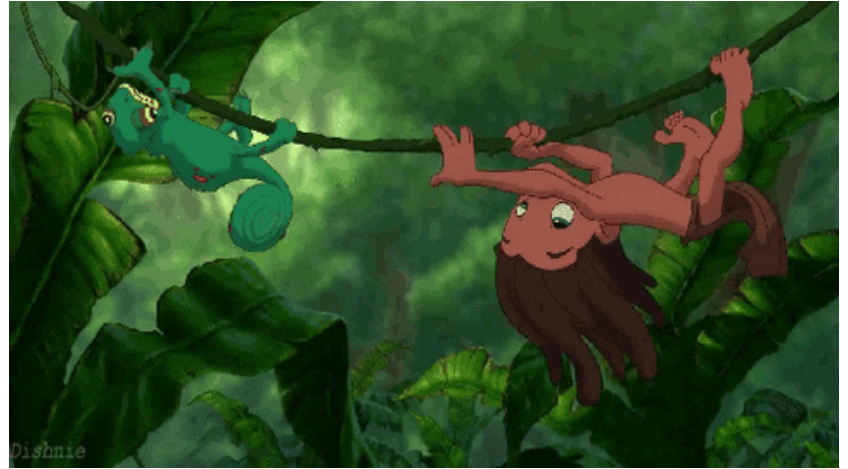
Structure

- Each script is one step, operating on a single subject
- One master file to run them analysis
- Keep multi-use components in a central place
 - All config set-ups go in here

```
1 # READ.MOIST350B
2 # SCOTT A ORR
3 # 19 MARCH 2017
4
5 # dir <- "/Users/orrscott/Desktop"
6 # filename <- "GraniteWalls_Set5b2.hfs"
7 #
8 # setwd(dir)
9 # filepath <- paste0(dir, "/", filename)
10 # d <- readLines(filepath)
11 #
12 # arrs.t <- hfRead(d)
13
14 ▶ #####
15
16 ▶ # FUNCTIONS #####
17
18 ▶ hfRead <- function(filepath) {  
40  
41 ▶ makeArray <- function(line.ref, layer.starts) {  
51  
52 ▶ makeLayer <- function(line.ref, dims, y.vals, x.vals) {  
63  
64 ▶ norm <- function(matrix) {  
69
```

“Put your faith in what you most believe in...”












- Don't trust someone else's code
- Never trust your own code !



Additional tips for code

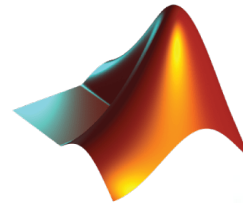
- Save the intermediate result of each script**
 - Useful if the processing takes a long time
 - Not necessarily that important if analysis is quick
 - Storage and memory limitations apply as well
- Plot everything

Delete superfluous files

 .DS_Store	17/08/2016 11:38	DS_STORE File	7 KB
 .Rhistory	01/08/2019 12:54	RHISTORY File	1 KB
 FLIR1760.csv	16/08/2016 16:31	Microsoft Excel Co...	515 KB
 FLIR1760.jpg	16/08/2016 16:20	JPG File	722 KB
 FLIR1760.txt	16/08/2016 16:36	Text Document	516 KB
 FLIR1814- photo.jpg	16/08/2016 16:21	JPG File	74 KB
 FLIR1814.csv	16/08/2016 17:09	Microsoft Excel Co...	515 KB
 FLIR1814.jpg	16/08/2016 16:21	JPG File	673 KB
 IRFit_fromT.R	16/08/2016 18:53	R File	6 KB
 Untitled.R	16/08/2016 17:42	R File	6 KB
 Untitled2.R	17/08/2016 12:11	R File	7 KB

Types of analytical tools

- Consider a wide range of users (or your future self)
- Open source (Python vs. MATLAB)
 - Need to balance with existing packages and tools



MATLAB



Final points

- Disclaimer: do as I say, not as I do
- Best practice and tools for re-use
 - Meet funder requirements
 - Foster future use (you would hope!)
 - Improve the robustness of current work
 - Support collaboration
 - Save time (now, and in future; and for all involved)
 - Avoid awkward (and sometimes serious) mistakes



Best practice and tools for use and re-use [of scientific data in archaeology and heritage]

Dr Scott Allan Orr

Lecturer in Heritage Data Science

SSHOC/SEADDA/E-RIHS Workshop on
Use and Re-Use of Scientific Data in Archaeology and Heritage

2 April 2020