

# A standardized data layout for archiving



Allan Reese and Grant Stentiford  
Centre for Environment, Fisheries and Aquaculture, Weymouth  
cefas.co.uk

## Why archive? Why standardize?

Data are expensive to collect but have lasting value. Climate change, chronic effects of pollution, and biodiversity change all prompt re-use of data and samples collected years earlier for other purposes. Data if not archived are lost to science, but re-used data may be misinterpreted when methods and assumptions are not described in sufficient detail - it's easy to assume "everyone knows that." Problems arise from acronyms, local jargon, made-up codes, implied dittos, and mixing data with subtotals etc.

We must all manage storage of data and samples more effectively. Digital storage is replacing visible storage of physical notebooks. Constraints imposed by the immediate IT environment may be mistaken for features essential for long-term storage. Files can be lost in cyber-space.

Adopting a standard layout eases all the tasks of depositing data, cataloguing files so they don't get lost, and understanding the content for appropriate re-use.

	A	B	C	D	E	F	G	H	I	J	K	L
1	RA	DATE	IDENTITY	STAT NO	AREA	LAT	LONG	E/W	SPECIES	LGTH	WGT	LV WGT
2	RA040399	22/06/2004	1	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	57	N80
3	RA040399	22/06/2004	2	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	60	1.1
4	RA040399	22/06/2004	3	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	57	0.8
5	RA040399	22/06/2004	4	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	63	N80
6	RA040399	22/06/2004	5	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	56	0.8
7	RA040399	22/06/2004	6	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	79	0.7
8	RA040399	22/06/2004	7	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	60	0.8
9	RA040399	22/06/2004	8	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	71	47	0.8
10	RA040399	22/06/2004	9	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	71	52	0.8
11	RA040399	22/06/2004	10	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	70	0.8
12	RA040399	22/06/2004	11	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	20	90	1.1
13	RA040399	22/06/2004	12	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	60	4.6
14	RA040399	22/06/2004	13	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	58	56	N80
15	RA040399	22/06/2004	14	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	21	93	1.2
16	RA040399	22/06/2004	15	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	71	47	N80
17	RA040399	22/06/2004	16	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	58	50	0.8
18	RA040399	22/06/2004	17	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	26	87	0.9
19	RA040399	22/06/2004	18	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	71	48	0.8
20	RA040399	22/06/2004	19	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	46	0.4
21	RA040399	22/06/2004	20	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	60	0.8
22	RA040399	22/06/2004	21	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	24	83	0.9
23	RA040399	22/06/2004	22	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	26	88	0.9
24	etc											

1) **Raw data sheet.** Table starts at top left (cell A1). No random text or other workings round or in-between the data. All cells are values, not formulae. First row used for one-word unique column headings (names), supplied by the user. Data table is rectangular, so repeat values must be copied onto each row and missing values flagged with non-blank value. The Manager should refer the data back to get the blank "LV weights" changed. Excel allows mixing numbers and text cells within a column; most statistics programs will not.

## Standards for Metadata

Metadata is the key to data, but standards from computer science describe storage structures rather than meaning. The "Dublin Core" (NISO Standard Z39.85-2001), for example, lists required fields such as "title" and "description", but leaves their content to the individual. This risks omitting information essential to new users outside the research group.

Social scientists are more used to data archiving. The ESRC requires grant holders to deposit their data with the UKDA (www.data-archive.ac.uk). Archive staff help to ensure that "data are deposited to a standard that enables them to be used by a third party, including the provision of adequate documentation. Consent, confidentiality, ethical and legal issues are considered, included in the project management plan and addressed before data collection starts."

Our own project clarified roles for data **managers** and data **collectors**. Managers are IT specialists who guarantee the physical security and continued readability of files, but take on trust that the contents make sense. Collectors - scientists or administrators - are perhaps too close to their data to notice the assumptions and feel more pressure to meet their own objectives than think of alternative uses.

1	RA	DATE	IDENTITY	STAT NO	AREA	LAT	LONG	E/W	SPECIES	LGTH	WGT	LV WGT	14	15	16	17	18	19	20	21	22	23	24	
2	RA040399	22/06/2004	1	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	57	N80												
3	RA040399	22/06/2004	2	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	60	1.1												
4	RA040399	22/06/2004	3	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	57	0.8												
5	RA040399	22/06/2004	4	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	63	N80												
6	RA040399	22/06/2004	5	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	56	0.8												
7	RA040399	22/06/2004	6	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	79	0.7												
8	RA040399	22/06/2004	7	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	60	0.8												
9	RA040399	22/06/2004	8	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	71	47	0.8												
10	RA040399	22/06/2004	9	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	71	52	0.8												
11	RA040399	22/06/2004	10	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	70	0.8												
12	RA040399	22/06/2004	11	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	20	90	1.1												
13	RA040399	22/06/2004	12	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	60	4.6												
14	RA040399	22/06/2004	13	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	58	56	N80												
15	RA040399	22/06/2004	14	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	21	93	1.2												
16	RA040399	22/06/2004	15	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	71	47	N80												
17	RA040399	22/06/2004	16	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	58	50	0.8												
18	RA040399	22/06/2004	17	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	26	87	0.9												
19	RA040399	22/06/2004	18	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	71	48	0.8												
20	RA040399	22/06/2004	19	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	46	0.4												
21	RA040399	22/06/2004	20	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	59	60	0.8												
22	RA040399	22/06/2004	21	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	24	83	0.9												
23	RA040399	22/06/2004	22	1	1 Inner Loughgan bog	52 17.6	4 17.0	W	DAB	26	88	0.9												
24	etc																							

2) **Codebook sheet.** Our standard column headings in the template set out what is needed to expand each column heading in the raw data. The name, which may be cryptic, is explained in a label and has the type, storage and meaning of values explained. Details of acceptable values or ranges can be applied to the data collection process, so that values are checked as acceptable and consistent (no pregnant males unless sea-horses) immediately on entry. The reference is vital, so that studies that use the same name with different definitions are not naively merged.

## A standard layout for data collectors

Like it or not, Microsoft Excel is widely used for data handling. Its help guidelines on LISTS are clear and relevant. Data relating to more than one level may imply an Access database, but the extra complication is often not justified. Scientific data tend to be stable once checked, and transaction processing is not a good model. However, both programs are so flexible that idiosyncratic use is the norm, making data structures hard to recognise.

We propose a **three-step procedure** for data submission. **Raw data** in columns form a single table (or view from a relational database); the **codebook** describes each column; and the **metadata** description uses fields from the Dublin Core. These three documents could be submitted as text files (eg CSV ASCII) but we use the familiarity of Excel to provide a template of three sheets in a workbook, as shown. This example from fish-disease monitoring is at the stage where it should be checked by the data manager or a third party.

Notes explaining the system in more detail as a standard operating procedure (SOP) are available from the first author.

1	A	B
2	RA	DATE
3	RA040399	22/06/2004
4	IDENTITY	STAT NO
5	1	1
6	AREA	LAT
7	1 Inner Loughgan bog	52 17.6
8	LONG	E/W
9	4 17.0	W
10	SPECIES	LGTH
11	DAB	59
12	WGT	LV WGT
13	57	N80
14	14	15
15	16	17
16	18	19
17	20	21
18	22	23
19	24	25
20	26	27
21	28	29
22	30	31
23	32	33
24	34	35
25	36	37
26	38	39
27	40	41
28	42	43
29	44	45
30	46	47
31	48	49
32	50	51
33	52	53
34	54	55
35	56	57
36	58	59
37	60	61
38	62	63
39	64	65
40	66	67
41	68	69
42	70	71
43	72	73
44	74	75
45	76	77
46	78	79
47	80	81
48	82	83
49	84	85
50	86	87
51	88	89
52	90	91
53	92	93
54	94	95
55	96	97
56	98	99
57	100	101
58	102	103
59	104	105
60	106	107
61	108	109
62	110	111
63	112	113
64	114	115
65	116	117
66	118	119
67	120	121
68	122	123
69	124	125
70	126	127
71	128	129
72	130	131
73	132	133
74	134	135
75	136	137
76	138	139
77	140	141
78	142	143
79	144	145
80	146	147
81	148	149
82	150	151
83	152	153
84	154	155
85	156	157
86	158	159
87	160	161
88	162	163
89	164	165