

R: A Hitchhikers Guide to Reproducible Research

- My favourite mistake

Brendan Palmer,

Clinical Research Facility - Cork &

School of Public Health

 @B_A_Palmer

Fundamental problem



I'm not in the office at the moment. Send any work to be translated

Beware of default settings

Ziemann et al. *Genome Biology* (2016) 17:177
DOI 10.1186/s13059-016-1044-7

Genome Biology

COMMENT

Open Access



Gene name errors are widespread in the scientific literature

Mark Ziemann¹, Yotam Eren^{1,2} and Assam El-Osta^{1,3*}

Abstract

The spreadsheet software Microsoft Excel, when used with default settings, is known to convert gene names to dates and floating-point numbers. A programmatic scan of leading genomics journals reveals that approximately one-fifth of papers with supplementary Excel gene lists contain erroneous gene name conversions.

frequently reused. Our aim here is to raise awareness of the problem.

We downloaded and screened supplementary files from 18 journals published between 2005 and 2015 using a suite of shell scripts. Excel files (.xls and .xlsx suffixes) were converted to tabular separated files (tsv) with ssconvert (v1.12.9). Each sheet within the Excel file was converted to a separate tsv file. Each column of data in the tsv file was screened for the presence of gene sym-

NCBI LocusLink

Search LocusLink Display Brief Organism: All

Query: Go Clear

View Hs NEDD5 One of 1 Loci Save All Loci

Click to Display tRNA-Genomic Alignments (spanning 38716 bps)

PUB OMIM REVIEW UNIGENE MAP VAR HOMOL GDB

ef UCSC

Homo sapiens Official Gene Symbol and Name (HGNC)

NEDD5: neural precursor cell expressed, developmentally down-regulated 5

LocusID: 4735

Overview Submit GeneRIF ?

Locus Type: gene with protein product, function known or inferred

Product: neural precursor cell expressed, developmentally down-regulated 5

Alternate Symbols: DIFF6, SEPT2, hNedd5, KIAA0158

Relationships ?

Mouse Homology Maps:

NCBI vs. MGD	1 cM	2-Sep	Hs Mm
UCSC vs. MGD	1 cM	Sept2	Hs Mm
UCSC vs. Hudson et al.	1 1319.34 cR	AW208991	Hs Mm

Map Information ?

Take small steps to big changes

THE AMERICAN STATISTICIAN
2018, VOL. 72, NO. 1, 2–10
<https://doi.org/10.1080/00031305.2017.1375989>



 OPEN ACCESS 

Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA

ABSTRACT

Spreadsheets are widely used software tools for data entry, storage, analysis, and visualization. Focusing on the data entry and storage aspects, this article offers practical recommendations for organizing spreadsheet data to reduce errors and ease later analyses. The basic principles are: be consistent, write dates like YYYY-MM-DD, do not leave any cells empty, put just one thing in a cell, organize the data as a single rectangle (with subjects as rows and variables as columns, and with a single header row), create a data dictionary, do not include calculations in the raw data files, do not use font color or highlighting as data, choose good names for things, make backups, use data validation to avoid data entry errors, and save the data in plain text files.

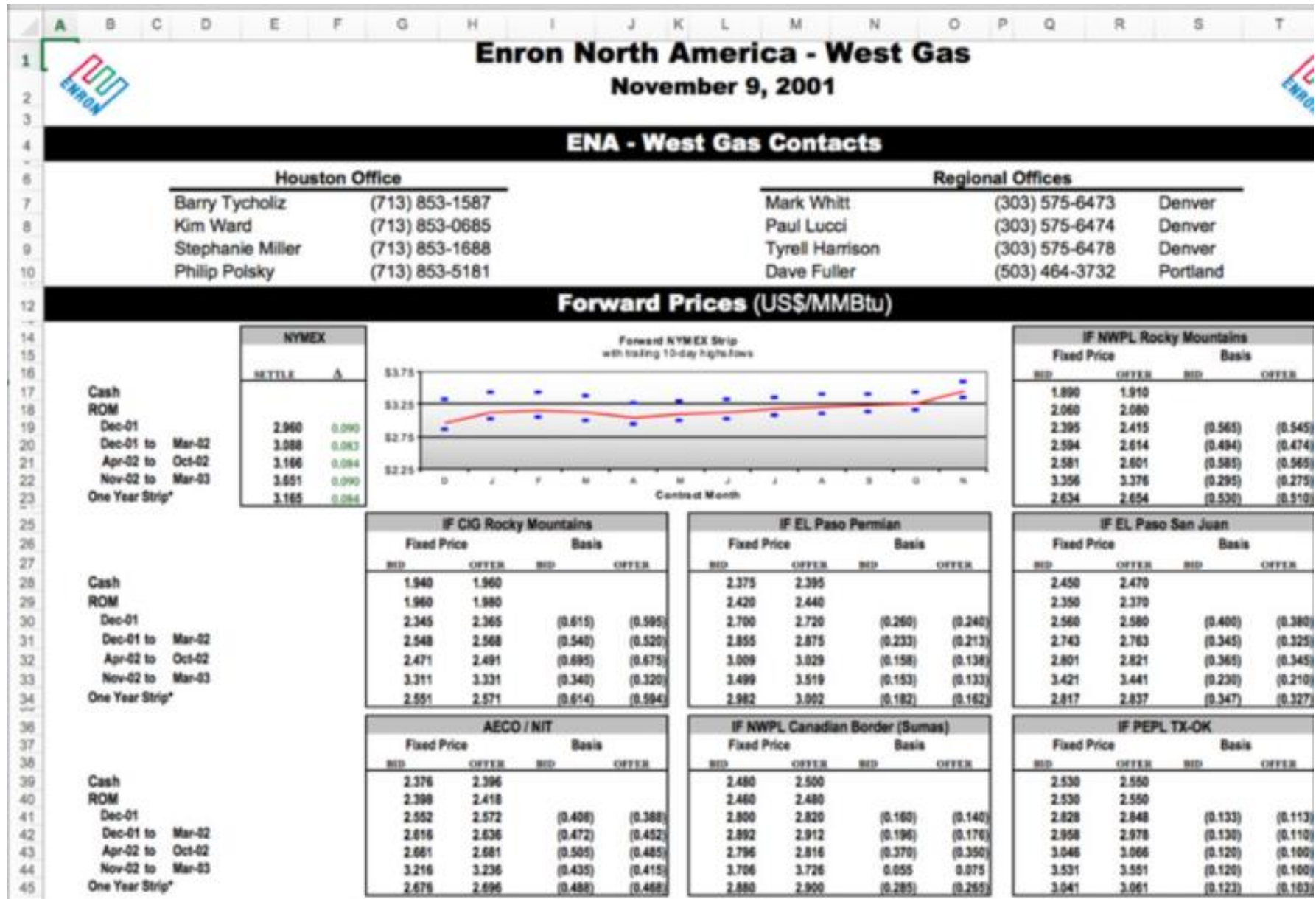
ARTICLE HISTORY

Received June 2017
Revised August 2017

KEYWORDS

Data management; Data organization; Microsoft Excel; Spreadsheets

This is a big problem



Less stress, more success

	A	B	C	D	E	F	G	H	I	J	K	L
1	id	week_no	filter_name	treatment	replicate_no	flavonoids	biomass	variety	date	investigator		
2	1	0	ptp	nofilter	1	1.061	0.39	cos	2019/04/01	Darren Dahly		
3	2	0	ptp	nofilter	2	1.1805	0.42	cos	2019/04/01	Darren Dahly		
4	3	0	ptp	nofilter	3	1.0345	0.62	cos	2019/04/01	Darren Dahly		
5	4	0	ptp	nofilter	4	1.094	0.63	cos	2019/04/01	Brendan Palmer		
6	1	0	my	nofilter	1	1.061	0.39	cos	2019/04/01	Brendan Palmer		
7	2	0	my	nofilter	2	1.1805	0.42	cos	2019/04/01	Brendan Palmer		
8	3	0	my	nofilter	3	1.0345	0.62	cos	2019/04/01	Brendan Palmer		
9	4	0	my	nofilter	4	1.094	0.63	cos	2019/04/01	Brendan Palmer		
10	1	0	ca	nofilter	1	1.061	0.39	cos	2019/04/01	Brendan Palmer		
11	2	0	ca	nofilter	2	1.1805	0.42	cos	2019/04/01	Brendan Palmer		
12	3	0	ca	nofilter	3	1.0345	0.62	cos	2019/04/01	Brendan Palmer		
13	4	0	ca	nofilter	4	1.094	0.63	cos	2019/04/01	Darren Dahly		
14	5	1	ptp	filter	1	0.87	0.76	cos	2019/04/08	Darren Dahly		
15	6	1	ptp	filter	2	0.847	0.95	cos	2019/04/08	Darren Dahly		
16	7	1	ptp	filter	3	1.022	0.95	cos	2019/04/08	Darren Dahly		
17	8	1	ptp	filter	4	0.916	0.95	cos	2019/04/08	Darren Dahly		
18	9	1	my	filter	1	1.119	1.55	cos	2019/04/08	Darren Dahly		
19	10	1	my	filter	2	0.845	3.16	cos	2019/04/08	Darren Dahly		
20	11	1	my	filter	3	1.299	4.9	cos	2019/04/08	Brendan Palmer		
21	12	1	my	filter	4	1.149	5.5	cos	2019/04/08	Brendan Palmer		
22	13	1	ca	filter	1	0.716	5.5	cos	2019/04/08	Brendan Palmer		
23	14	1	ca	filter	2	0.881	7.94	cos	2019/04/08	Brendan Palmer		
24	15	1	ca	filter	3	0.586	8.71	cos	2019/04/08	Brendan Palmer		
25	16	1	ca	filter	4	0.561	8.71	cos	2019/04/08	Brendan Palmer		
26	17	2	ptp	filter	1	0	14.45	cos	2019/04/15	Brendan Palmer		
27	18	2	ptp	filter	2	1.006	2.14	cos	2019/04/15	Brendan Palmer		
28	19	2	ptp	filter	3	1.236	1.86	cos	2019/04/15	Brendan Palmer		
29	20	2	ptp	filter	4	1.206	1.2	cos	2019/04/15	Brendan Palmer		
30	21	2	mv	filter	1	1.545	2.45	cos	2019/04/15	Brendan Palmer		
		data	dictionary	values								

Less stress, more success

	A	B	C	D	E	F	G	H	I	J	K	L
1	id	week_no	filter_name	treatment	replicate_no	flavonoids	biomass	variety	date	investigator		
2	1	0	ptp	nofilter	1	1.061	0.39	cos	2019/04/01	Darren Dahly		
3	2	0	ptp									
4	3	0	ptp	1	field_name	data_type	data_format	example	standard_units	description		
5	4	0	ptp	2	id	numeric	integer	23	NA	Unique identifier applied to each observation		
6	1	0	my	3	week_no	numeric	integer	1	NA	Week number, 1 = 7 days exposure, 2 = 14 days exposure		
7	2	0	my	4	filter_name	character	NA	my	NA	3 filter types; 'ptp' = polytunnel plastic blocks all UV light		
8	3	0	my	5	treatment	character	NA	filter	NA	Presence or absence of a filter at the time of sampling		
9	4	0	my	6	replicate_no	numeric	integer	1	NA	The number of replicates in each treatment		
10	1	0	ca	7	flavonoids	numeric	double	0.3421	parts per million (ppm)	Leaf disc taken from the tip of the most mature leaf at th		
11	2	0	ca	8	biomass	numeric	double		gram (g)	Above ground biomass on the day of harvest		
12	3	0	ca	9	variety	character	NA	cos	NA	3 commerical varieties of red lettuce used; 'cos' = Cos Di		
13	4	0	ca	10	date	date	YYYY/MM/DD	2019/06/28	ISO 8601	Experiment date		
14	5	1	ptp	11	investigator	character	Firstname Lastname	Aoife Coffey	NA	Primary researcher who performed the experiment		
15	6	1	ptp	12								
16	7	1	ptp	13								
17	8	1	ptp	14								
18	9	1	my	15								
19	10	1	my	16								
20	11	1	my	17								
21	12	1	my	18								
22	13	1	ca	19								
23	14	1	ca	20								
24	15	1	ca	21								
25	16	1	ca	22								
26	17	2	ptp	23								
27	18	2	ptp	24								
28	19	2	ptp	25								
29	20	2	ptp	26								
30	21	2	mv	27								
		data	dictionary	28								
				29								
				30								
					data	dictionary	values					

Less stress, more success

[illegible]

Step by step guide

← → ↺ 🏠 <https://www.youtube.com/watch?v=Ry2xjTBtNFE>



Search



Book1 - Excel (Product Activation Failed)

File Home Insert Page Layout Formulas Data Review View Tell me what you want to do... Dahly, Darren Share

Clipboard Font Alignment Number Styles Cells Editing

Calibri 11 A A Bold Italic Underline Wrap Text Merge & Center General Conditional Formatting Table Normal Bad Good Neutral Calculation Check Cell Insert Delete Format AutoSum Fill Clear Sort & Find & Filter Select

D2

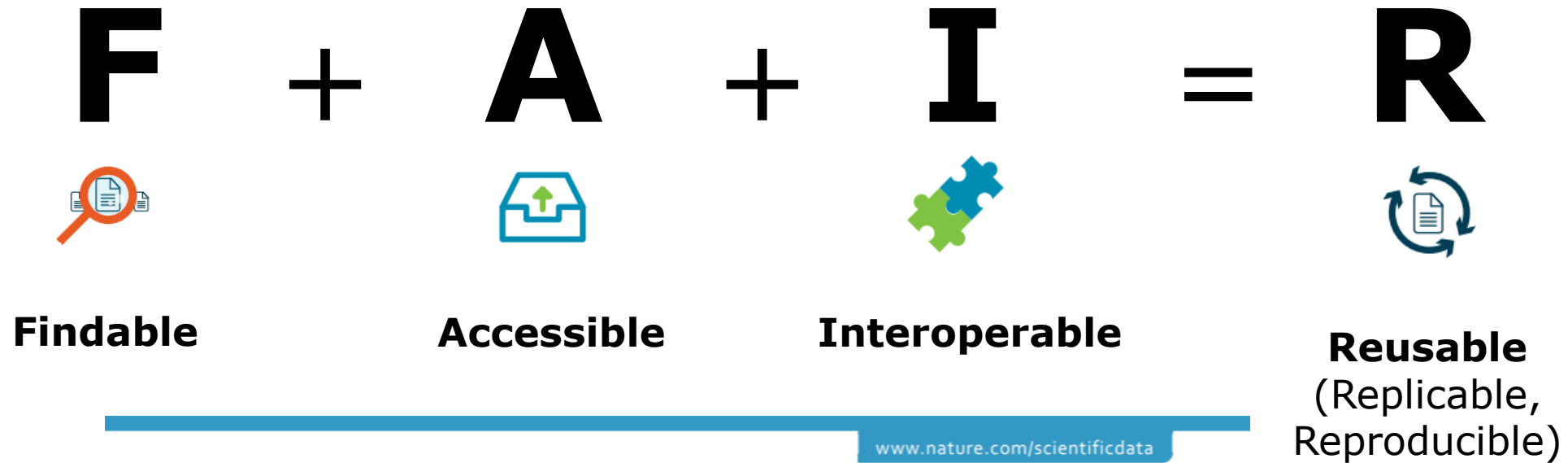
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1	id	gender	gender_other	age	nationality	year_program																
2																						
3																						
4																						
5																						
6																						
7																						
8																						
9																						
10																						
11																						
12																						
13																						
14																						
15																						
16																						
17																						
18																						
19																						
20																						
21																						
22																						
23																						
24																						
25																						
26																						
27																						
28																						

Sheet1 Sheet2 Sheet3

Ready

Type here to search

The movement towards FAIR data



SCIENTIFIC DATA 

Amended: Addendum

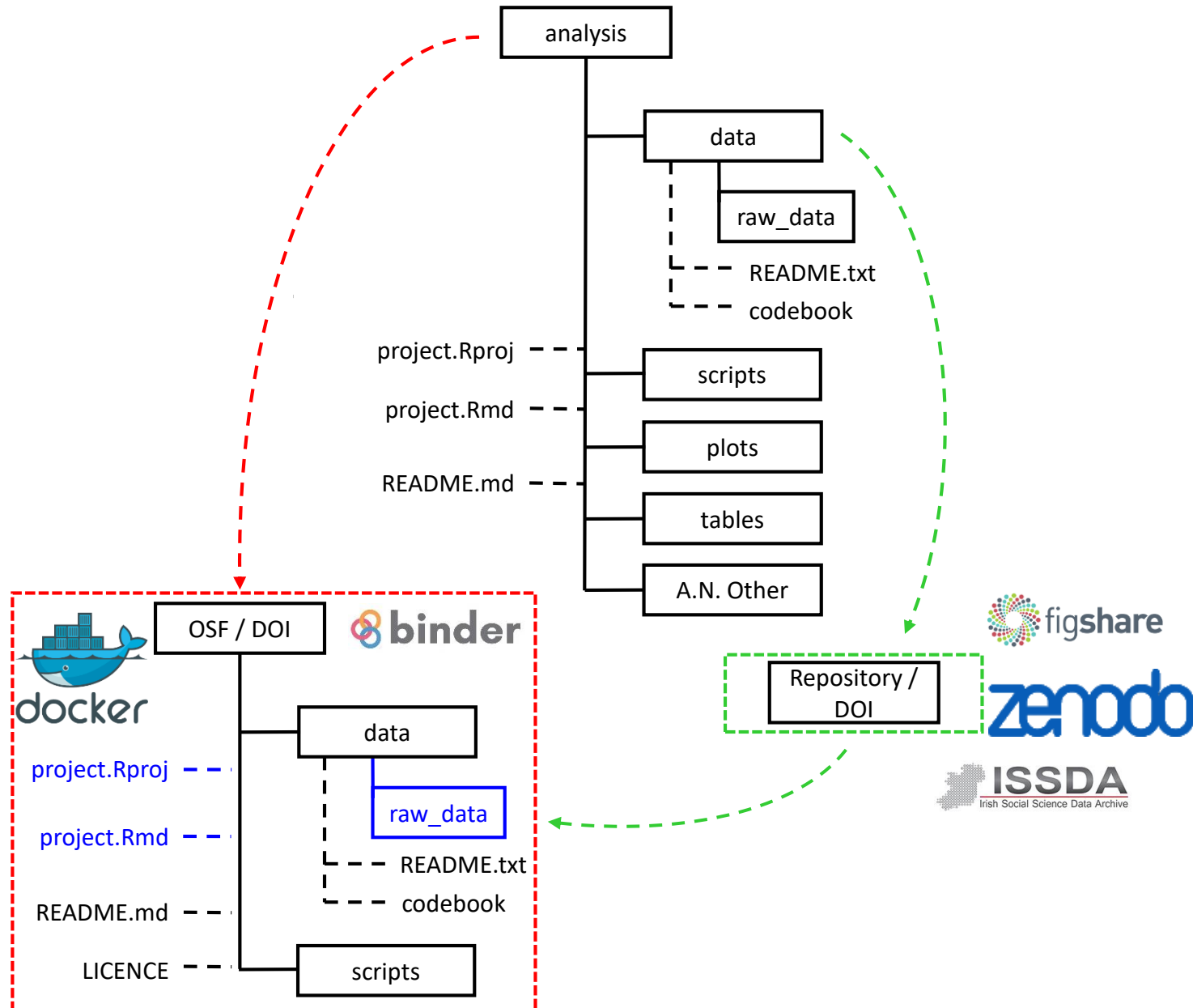
OPEN
SUBJECT CATEGORIES

» Research data
» Publication
characteristics

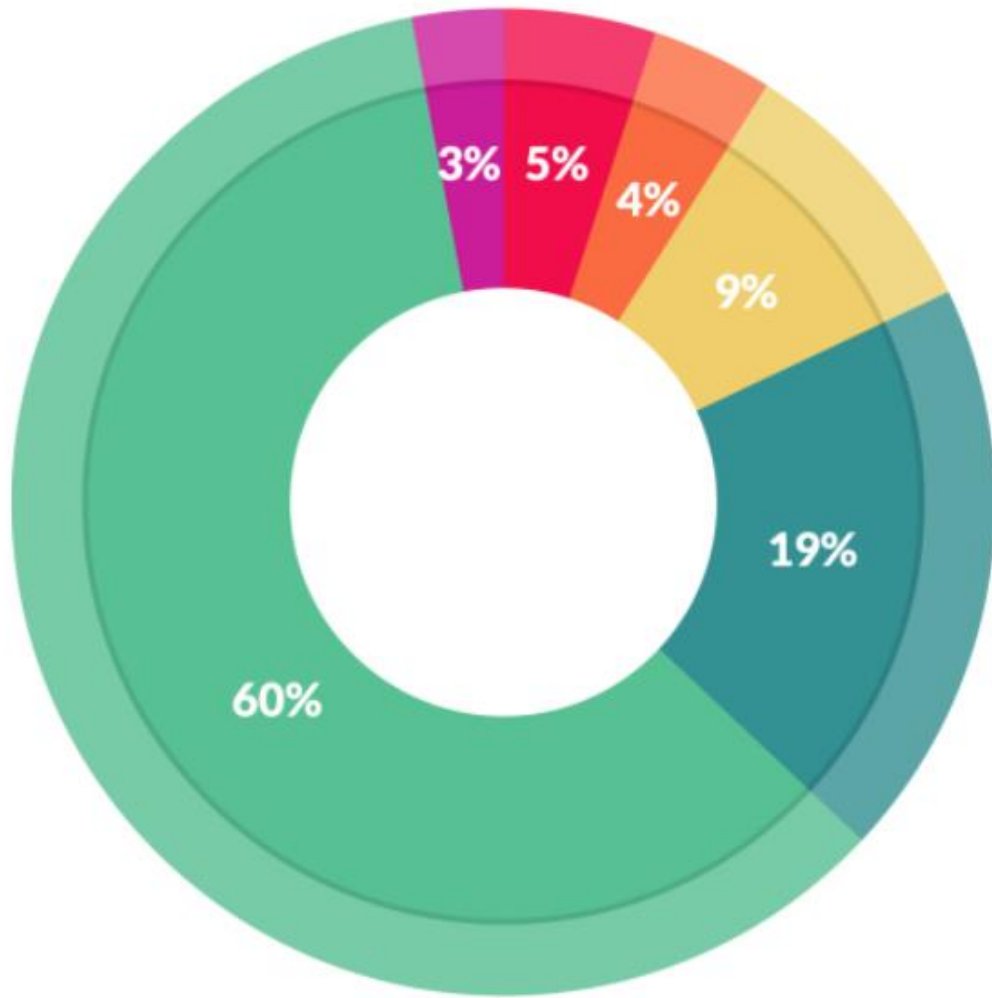
**Comment: The FAIR Guiding
Principles for scientific data
management and stewardship**

Mark D. Wilkinson *et al.*[#]

What does this allow us to do?



Resources are being wasted by not doing this



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%