

Authoring tools to automate data sharing in scientific publishing

John Kitchin

March 14, 2016

Contents

1	Title slide	2
2	Introduction	2
3	Why is data sharing hard today?	2
4	Sharing is tedious, error prone	3
5	Our approach to this issue is integration of tools/data	3
6	Emacs + org-mode for reproducible, functional scientific documents	3
6.1	sub-heading	4
6.1.1	sub-sub-heading	4
7	We have used it extensively in scientific publishing	5
8	Reusing the data	6
9	Automating data embedding sharing	7
9.1	HTML export	7
9.2	PDF export	8
9.3	Vanilla export	8
10	What makes this integration possible?	8
11	Concluding thoughts	8

12 Extract the references	9
13 Getting started	10
14 Appendix	11
14.1 The custom export code	11

1 Title slide

2 Introduction

- Data sharing is not easy
- But it is increasingly required
- And increasingly desired
 - The big question is: "How do we make that happen?"
- Nobody wants more work
- We would all like to get more out of what we already do
- We are working on tools that make that happen that I will talk about today

3 Why is data sharing hard today?

- For many scientific authors data and analysis are not integrated into writing tools
 - Think about a manuscript in MS Word/L^AT_EX
 - Data in Excel, data files, . . .
 - Plots in Origin/SigmaPlot/etc.
 - Separate script files for analysis
 - all with limited interconnectivity
 - * e.g. where is the data in Fig 2 of the manuscript?
 - * Where is the file/script that did the data analysis?
 - * How do I reuse the data in Table 2 for a new purpose? Copy and Paste?

- These tools are not especially well-integrated and not easily adapted to new use cases

4 Sharing is tedious, error prone

- Sharing then becomes extra work to generate supporting information that reconstructs the effort, copies data you think you used, etc.
- Reconstruction of the work that went into the manuscript is error prone
- and tedious. . .
- Reuse is not much better
- Yet increasingly required, and desired

5 Our approach to this issue is integration of tools/data

- We developed an integrated set of tools that makes data/code/analysis part of the manuscript preparation/submission process could help with data sharing
- It would leverage the work we already do in scientific writing
- and provide access to reusable data/code
- This does require development of a relatively new tool chain for writing
- We have done this in Emacs + org-mode + code
- It is all open-source (<http://github.com/jkitchin/jmax>) and cross-platform

6 Emacs + org-mode for reproducible, functional scientific documents

- org-mode is basically a plain text markup language deeply integrated into Emacs (an editor)
- Outline mode at the core, and much much more

6.1 sub-heading

- Narrative text, equations $\int_0^1 e^x dx$, images



6.1.1 sub-sub-heading

Functional tables

a	b
2	4
0	1

3b2a6830d248f31580202ccdbadd5c49.csv:

- citations¹
- Integrated functional code

1 `date`

a9564ebc3289b7a14551baf8ad5ec60a:

```
RESULTS: Mon Mar 14 13:11:05 PDT 2016
```

- Functional links that can




- Open document locations
- Open mail, news, urls
- run user-defined code in almost any language

```
1 import time
2 print(time.asctime())
```

309ec9e61913846c23962de976f11728.py:


```
RESULTS: Mon Mar 14 13:11:36 2016
```

```
1 (current-time-string)
```

7f7bac1b1e471114e487e90437582fbd.elisp: 


```
RESULTS: Mon Mar 14 13:11:45 2016
```

```
1 Sys.time()
```

3f06a7efae851f9690a285a46ca320dc.R: 

```
RESULTS: [1] "2016-03-14 13:11:51 PDT"
```

```
1 #include <time.h>
2 #include <stdlib.h>
3 #include <stdio.h>
4
5 int main() {
6     time_t current_time;
7     char* c_time_string;
8     current_time = time(NULL);
9     /* Convert to local time format. */
10    c_time_string = ctime(&current_time);
11    printf("%s", c_time_string);
12    return 0;
13 }
```

812a422b1846832bae23932b69c07e4c.c: 

```
RESULTS: Mon Mar 14 13:11:56 2016
```

7 We have used it extensively in scientific publishing

- Over dozen papers in print by my group illustrating what can be done
- ² - all experimental

- [3](#) - mixed experiment/computation
- [4](#) - computational DFT on oxides
- [5](#) - computational coverage/site dependence
- [6](#) - full sql database described in SI
- [7](#) - exp/computation H2/D2 exchange on CuPd
- [8](#) - exp/computation XPS on CuPd alloys
- [9](#) - 1.8 GB dataset on Zenodo
- [10](#) - computation DFT+U
- [11](#) - examples of reusing data
- [1](#) - examples of reusing data
- [12](#) - Molecular simulation

Let's see a working example

- [13](#) The data is available in the SI
- The json database described in SI
- The source can be extracted from the PDF (goto line 336)

```
1 pdftk am4059149_si_001.pdf  unpack_files
```

3a4c9126c7b0a1737bd6a65d910ca25a: 

That SI file was generated here (goto line 336): [supporting-information.org](https://www.supporting-information.org)


8 Reusing the data

That data is human readable - and machine addressable

```
1 (remove-if-not (lambda (x) (string= "rutile" (nth 1 x))) data)
```


1f85d0194e997e7ea8f170bf787a8598.elisp: 

TiO ₂	rutile	LDA	-2801.64	30.58	259.47
TiO ₂	rutile	AM05	-2733.53	31.31	233.2
TiO ₂	rutile	PBEsol	-2759.29	31.22	239.76
TiO ₂	rutile	PBE	-2773.21	32.11	215.78

089d3777c582b5e548c243b5c9fa5229.csv: 

If you prefer Python, no problem. Here we get the anatase data:

```
1 return [x for x in data if x[1] == 'anatase']
```

32c09a6e325db19533e2e272caed35fd.py: 

TiO ₂	anatase	LDA	-2802.73	33.62	187.4
TiO ₂	anatase	AM05	-2741.12	34.33	178.26
TiO ₂	anatase	PBEsol	-2763.61	34.25	178.71
TiO ₂	anatase	PBE	-2781.16	35.13	171.42

6f02bbf48d7abde528dc614a9a7a84f1.csv: 

9 Automating data embedding sharing

- org-mode is great - If you use org-mode
- In an [Appendix](#) there is code that automatically embeds data and code in org-mode into HTML and PDF.
- One source to many outputs
- We can extract the source code and load it here

```
1 (org-babel-tangle)
2 (load-file "data-sharing.el")
```

feccf7cb9ed36d8b57790a1ed1e6933.elisp: 


9.1 HTML export

```
1 (custom-export-and-open-html)
```

94b98fc6d9a07d6041b0d18ad27acb59.elisp: 

9.2 PDF export

¹ `(custom-export-and-open-pdf)`

7ebdbc39a7bc97162b6a263e0e68f7f0.elisp: 

9.3 Vanilla export

`(org-open-file(org-latex-export-to-pdf))`

10 What makes this integration possible?

- An extensible editor
 - Extensible in a full programming language
 - This allows the tool to become what you want
 - Emacs is ideal for this
- A lightweight markup language
 - to differentiate text, code, data
 - Org-mode is also ideal for this
 - * Part structured markup, part api
 - * Very good compromise on function and utility with authoring ease
- `</code>` Since we use code to generate and analyse data, this solution works especially well

11 Concluding thoughts

- Emacs + org-mode + `</code>` enables a lot of very exciting capabilities in publishing and data sharing
 - Integrated narrative text, data, code
 - Export to a broad range of other formats
 - Interaction with the world (other computers, instruments) via APIs
 - * Materials Project, translation, Internet of Things, ...

- The future is very exciting
- We are not waiting for someone to figure out what we want
 - Anyway, by the time they deliver it we will need something else ;)

12 Extract the references

```

1 (save-window-excursion
2 (save-restriction
3 (widen)
4 (org-ref-bibliography)
5 (buffer-string)))

```

4e0979a1b8eb6bf328e800e4164ab2c7.elisp: 

1. cite:kitchin-2015-examp Kitchin, John R., "Examples of Effective Data Sharing in
2. cite:hallenbeck-2013-effec-o2 Hallenbeck, Alexander P. and Kitchin, John R., "Ef
3. cite:miller-2014-simul-temper Spencer D. Miller and Vladimir V. Pushkarev and An
4. cite:xu-2014-relat Zhongnan Xu and John R. Kitchin, "Relating the Electronic Str
5. cite:xu-2014-probin-cover Zhongnan Xu and John R. Kitchin, "Probing the Coverage
6. cite:curnan-2014-effec-concen Curnan, Matthew T. and Kitchin, John R., "Effects
7. cite:boes-2015-estim-bulk Jacob R. Boes and Gamze Gumuslu and James B. Miller an
8. cite:boes-2015-core-cu Jacob R. Boes and Peter Kondratyuk and Chunrong Yin and J
9. cite:xu-2015-linear-respon Xu, Zhongnan and Rossmeisl, Jan and Kitchin, John R.,
10. cite:xu-2015-accur-u Xu, Zhongnan and Joshi, Yogesh V. and Raman, Sumathy and K
11. cite:kitchin-2015-data-surfac-scienc John R. Kitchin, "Data Sharing in Surface Sc
12. cite:boes-2016-neural-networ Jacob R. Boes and Mitchell C. Groenenboom and John

- 13. cite:mehta-2015-ident-poten Prateek Mehta and Paul A. Salvador and John R. Kitchin
- 14. cite:pakin-attachfile Scott Pakin, "attachfile", , : ()

13 Getting started

Source code: <http://github.com/jkitchin/jmax>

Our starter-kit for Emacs + org-mode configured to do the things I showed you today Should work out of the box on Windows. Directions for using it on Mac/Linux.

Kitchingroup blog: <http://kitchingroup.cheme.cmu.edu>

Shon Feder @ShonFeder · Mar 11

The most exciting thing I've seen this year is a video of @johnkitchin
 [youtube.com/watch?v=2t925K...](https://www.youtube.com/watch?v=2t925K...) 

@johnkitchin

Check out our YouTube channel: <https://www.youtube.com/user/jrkitchin>



Introduction to org-ref

by John Kitchin

2 months ago · 1,333 views

This video shows new org-ref feat
 find the files used in this video he

1333 views (1800+ downloads of org-ref on MELPA!)



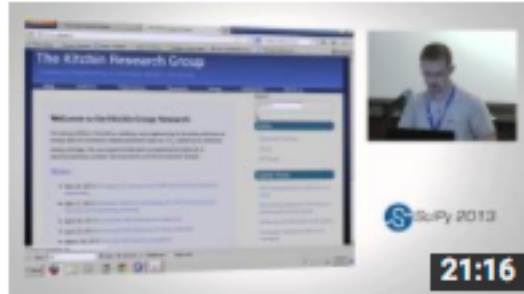
org mode is awesome

by John Kitchin

1 year ago · 20,984 views

Video showing many great

20,984 views



Emacs + org-mode + python in reproducibility

by Enthought

2 years ago • 23,094 views

Authors: Kitchin, John Carnegie Mellon University

...

23,094 views!

This talk: <https://github.com/jkitchin/ACS-2016-data-sharing>

14 Appendix

14.1 The custom export code

Here we define a custom table exporter. We use the regular table export mechanism, but save the contents of the table as a csv file. We define exports for two backends: \LaTeX and HTML. For \LaTeX , we use the `attachfile`¹⁴ package to embed the data file in the PDF. For HTML, we insert a link to the data file, and a data uri link to the HTML output. We store the filename of each generated table in a global variable named `*embedded-files*` so we can create a new Info metadata entry in the exported PDF.

```
1 (defvar *embedded-files* '())
2   "List of files embedded in the output."
3
4 (defun my-table-format (table contents info)
5   (let* ((tblstart (org-element-property
6                   :contents-begin table))
7         (tbl-data (save-excursion
8                    (goto-char tblstart)
9                    (org-babel-del-hlines
10                     (org-babel-read-table))))
11         (tblname (or (org-element-property :name table) (md5 (format "%s" tbl-data))))
12         (format (elt (plist-get info :back-end) 2))
13         (csv-file (concat tblname ".csv"))
14         (data-uri-data))
15
16   ;; Here we convert the table data to a csv file
17   (with-temp-file csv-file
18     (loop for row in tbl-data
19          do
20            (insert
21             (mapconcat
22              (lambda (x) (format "%s" x))
```

```

23         row
24         ", " )
25         (insert "\n"))
26     (setq data-uri-data
27         (base64-encode-string
28         (buffer-string)))
29
30     (add-to-list '*embedded-files* csv-file)
31
32     (cond
33     ;; HTML export
34     ((eq format 'html)
35      (concat
36      (org-html-table table contents info)
37      (format "<a href=\"%s\">%s</a>"
38             csv-file csv-file)
39      " "
40      (format (concat "<a href=\"data:text/csv;"
41                    "charset=US-ASCII;"
42                    "base64,%s\">data uri</a>"
43                    data-uri-data)))
44     ;; LaTeX/PDF export
45     ((eq format 'latex)
46      (concat
47      (org-latex-table table contents info)
48      "\n"
49      (format "%s: \\attachfile{%s}"
50             csv-file csv-file))))))

```

table-format.elisp: 

Next, we define an exporter for source blocks. We will write these to a file too, and put links to them in the exported files. We store the filename of each generated source file in a global variable named `*embedded-files*` so we can create a new Info metadata entry in the exported PDF.

```

1  (defun my-src-block-format (src-block contents info)
2    "Custom export for src-blocks.
3  Saves code in block for embedding. Provides backend-specific
4  output."
5    (let* ((srcname (org-element-property :name src-block))
6           (lang (org-element-property :language src-block))
7           (value (org-element-property :value src-block))
8           (format (elt (plist-get info :back-end) 2))
9           (exts `(("python" . ".py")
10                  ("emacs-lisp" . ".elisp")
11                  ("C" . ".c")
12                  ("R" . ".R"))))
13      (fname (concat
14              (or srcname (md5 value))
15              (cdr (assoc lang exts))))
16      (data-uri-data)
17
18      (with-temp-file fname

```

```

19     (insert value)
20     (setq data-uri-data (base64-encode-string
21                         (buffer-string))))
22
23     (add-to-list '*embedded-files* fname)
24
25     (cond
26     ;; HTML export
27     ((eq format 'html)
28      (concat
29       (org-html-src-block src-block contents info)
30       (format "<a href=\"%s\">%s</a>" fname fname)
31       " "
32       (format (concat "<a href=\"data:text/%s;"
33                     "charset=US-ASCII;base64,"
34                     "%s\">code uri</a>"
35                     lang data-uri-data))))
36     ;; LaTeX/PDF export
37     ((eq format 'latex)
38      (concat
39       (org-latex-src-block src-block contents info)
40       "\n"
41       (format "%s: \\attachfile{%s}" fname fname))))))

```


src-block-format.elisp: 

Finally, we also modify the results of a code block so they will appear in a gray box and stand out from the text more clearly.

```

1  (defun my-results (fixed-width contents info)
2    "Transform a results block to make it more visible."
3    (let ((results (org-element-property :results fixed-width))
4          (format (elt (plist-get info :back-end) 2))
5          (value (org-element-property :value fixed-width)))
6      (cond
7      ((eq 'latex format)
8       (format "\\begin{tcolorbox}
9       \\begin{verbatim}
10      RESULTS: %s
11      \\end{verbatim}
12      \\end{tcolorbox}"
13          value))
14      (t
15       (format "<pre>RESULTS: %s</pre>" value))))

```

dafeb6b72e57a11595885a79d0ce2cbe.elisp: 

RESULTS: my-results


An author may also choose to embed a file into their document, using the `attachfile` package for L^AT_EX. Here, we leverage the ability of org-mode

to create functional links that can be exported differently for L^AT_EX and HTML. We will create an attachfile link, and set it up to export as a L^AT_EX command or as a data URI for HTML.

```

1 (org-add-link-type
2   "attachfile"
3   (lambda (path) (org-open-file path))
4   ;; formatting
5   (lambda (path desc format)
6     (cond
7       ((eq format 'html)
8        ;; we want a data URI to the file name
9        (let* ((content
10              (with-temp-buffer
11                (insert-file-contents path)
12                (buffer-string)))
13              (data-uri
14                (base64-encode-string
15                  (encode-coding-string content 'utf-8))))
16              (add-to-list '*embedded-files* path)
17              (format (concat "<a href=\"data:;base64,\"
18                            \"%s\">%s</a>\"
19                            data-uri
20                            path))))
21              ((eq format 'latex)
22               ;; write out the latex command
23               (add-to-list '*embedded-files* path)
24               (format "\\attachfile{%s}" path))))))

```


attachfile-link.elisp: 

Here, we define a derived backend for HTML and L^AT_EX export. These are identical to the standard export backends, except for the modified behavior of the table and src-block elements.

```

1 (org-export-define-derived-backend 'my-html 'html
2   :translate-alist '((table . my-table-format)
3                     (src-block . my-src-block-format)
4                     (fixed-width . my-results)))
5
6 (org-export-define-derived-backend 'my-latex 'latex
7   :translate-alist '((table . my-table-format)
8                     (src-block . my-src-block-format)
9                     (fixed-width . my-results)))

```

2efb34a32a9c4653ff697c1d00fd294b.elisp: 

```


1 (defun custom-export-and-open-html ()
2   "Use my-html custom exporter and open the file."
3   (let* ((base (file-name-nondirectory

```

```

4         (file-name-sans-extension (buffer-file-name))))
5     (html (concat base ".html")))
6     (save-restriction (widen)
7         (browse-url (org-export-to-file 'my-html html))))))
8
9 (defun custom-export-and-open-pdf ()
10  "Use my-latex custom exporter and open pdf."
11  (save-restriction
12    (widen)
13    (let* ((org-latex-image-default-width "")
14           (*embedded-files* '())
15           (base (file-name-nondirectory
16                 (file-name-sans-extension (buffer-file-name))))
17           (tex (concat base ".tex"))
18           (pdf (concat base ".pdf"))
19           (org-latex-minted-options
20            (append
21             org-latex-minted-options
22             '(("xleftmargin" "\\parindent")))))
23      (org-export-to-file 'my-latex tex)
24      (ox-manuscript-latex-pdf-process tex)
25      (org-open-file pdf))))

```

209ed2ce9b5918b18e58d1c9eb91a1c2.elisp: 

References

- [1] John R. Kitchin. Examples of effective data sharing in scientific publishing. *ACS Catalysis*, 5(6):3894–3899, 2015. doi: 10.1021/acscatal.5b00538. URL <http://dx.doi.org/10.1021/acscatal.5b00538>.
- [2] Alexander P. Hallenbeck and John R. Kitchin. Effects of O₂ and SO₂ on the capture capacity of a primary-amine based polymeric CO₂ sorbent. *Industrial & Engineering Chemistry Research*, 52(31):10788–10794, 2013. doi: 10.1021/ie400582a. URL <http://pubs.acs.org/doi/abs/10.1021/ie400582a>.
- [3] Spencer D. Miller, Vladimir V. Pushkarev, Andrew J. Gellman, and John R. Kitchin. Simulating temperature programmed desorption of oxygen on Pt(111) using DFT derived coverage dependent desorption barriers. *Topics in Catalysis*, 57(1-4):106–117, 2014. doi: 10.1007/s11244-013-0166-3. URL <http://dx.doi.org/10.1007/s11244-013-0166-3>.
- [4] Zhongnan Xu and John R. Kitchin. Relating the electronic structure and reactivity of the 3d transition metal monoxide surfaces. *Catalysis Communications*, 52:60–64, 2014. ISSN 1566-7367. doi:

- 10.1016/j.catcom.2013.10.028. URL <http://dx.doi.org/10.1016/j.catcom.2013.10.028>.
- [5] Zhongnan Xu and John R. Kitchin. Probing the coverage dependence of site and adsorbate configurational correlations on (111) surfaces of late transition metals. *J. Phys. Chem. C*, 118(44):25597–25602, 2014. doi: 10.1021/jp508805h. URL <http://dx.doi.org/10.1021/jp508805h>.
- [6] Matthew T. Curnan and John R. Kitchin. Effects of concentration, crystal structure, magnetism, and electronic structure method on first-principles oxygen vacancy formation energy trends in perovskites. *The Journal of Physical Chemistry C*, 118(49):28776–28790, 2014. doi: 10.1021/jp507957n. URL <http://dx.doi.org/10.1021/jp507957n>.
- [7] Jacob R. Boes, Gamze Gumuslu, James B. Miller, Andrew J. Gellman, and John R. Kitchin. Estimating bulk-composition-dependent H₂ adsorption energies on Cu_xPd_{1-x} alloy (111) surfaces. *ACS Catalysis*, 5: 1020–1026, 2015. doi: 10.1021/cs501585k. URL <http://dx.doi.org/10.1021/cs501585k>.
- [8] Jacob R. Boes, Peter Kondratyuk, Chunrong Yin, James B. Miller, Andrew J. Gellman, and John R. Kitchin. Core level shifts in Cu-Pd alloys as a function of bulk composition and structure. *Surface Science*, 640:127–132, 2015. ISSN 0039-6028. doi: 10.1016/j.susc.2015.02.011. URL <http://dx.doi.org/10.1016/j.susc.2015.02.011>.
- [9] Zhongnan Xu, Jan Rossmeisl, and John R. Kitchin. A linear response DFT+U study of trends in the oxygen evolution activity of transition metal rutile dioxides. *The Journal of Physical Chemistry C*, 119(9): 4827–4833, 2015. doi: 10.1021/jp511426q. URL <http://dx.doi.org/10.1021/jp511426q>.
- [10] Zhongnan Xu, Yogesh V. Joshi, Sumathy Raman, and John R. Kitchin. Accurate electronic and chemical properties of 3d transition metal oxides using a calculated linear response U and a DFT + U(V) method. *The Journal of Chemical Physics*, 142(14):144701, 2015. doi: 10.1063/1.4916823. URL <http://scitation.aip.org/content/aip/journal/jcp/142/14/10.1063/1.4916823>.
- [11] John R. Kitchin. Data sharing in surface science. *Surface Science*, N/A:in press, 2015. ISSN 0039-6028. doi: 10.1016/j.susc.2015.05.007. URL <http://www.sciencedirect.com/science/article/pii/S0039602815001326>.

- [12] Jacob R. Boes, Mitchell C. Groenenboom, John A. Keith, and John R. Kitchin. Neural network and reaxff comparison for Au properties. *Accepted 1/2016, Int. J. Quantum Chemistry*, 2016.
- [13] Prateek Mehta, Paul A. Salvador, and John R. Kitchin. Identifying potential BO_2 oxide polymorphs for epitaxial growth candidates. *ACS Appl. Mater. Interfaces*, 6(5):3630–3639, 2015. doi: 10.1021/am4059149. URL <http://dx.doi.org/10.1021/am4059149>.
- [14] Scott Pakin. attachfile. <http://www.ctan.org/tex-archive/macros/latex/contrib/attachfile>. v1.5b.