

# Reproducible Document Stack: Towards a Scalable Solution for Reproducible Articles

Emmy Tsang, Innovation Community Manager @ eLife  
Twitter: @eLifeInnovation | Labs: [elifesci.org/labs](https://elifesci.org/labs)



# What is eLife?

hhmi | Howard Hughes  
Medical Institute



Knut and Alice  
Wallenberg  
Foundation

- A non-profit **run by researchers**, backed by research funders to drive reform in research communication
- £21m investment to date
- We invest heavily in open-source technology development and innovation **on behalf of the community**

Helping scientists **accelerate discovery** by  
operating a platform for research **communication**  
that encourages and recognises **the most**  
**responsible behaviours in science.**

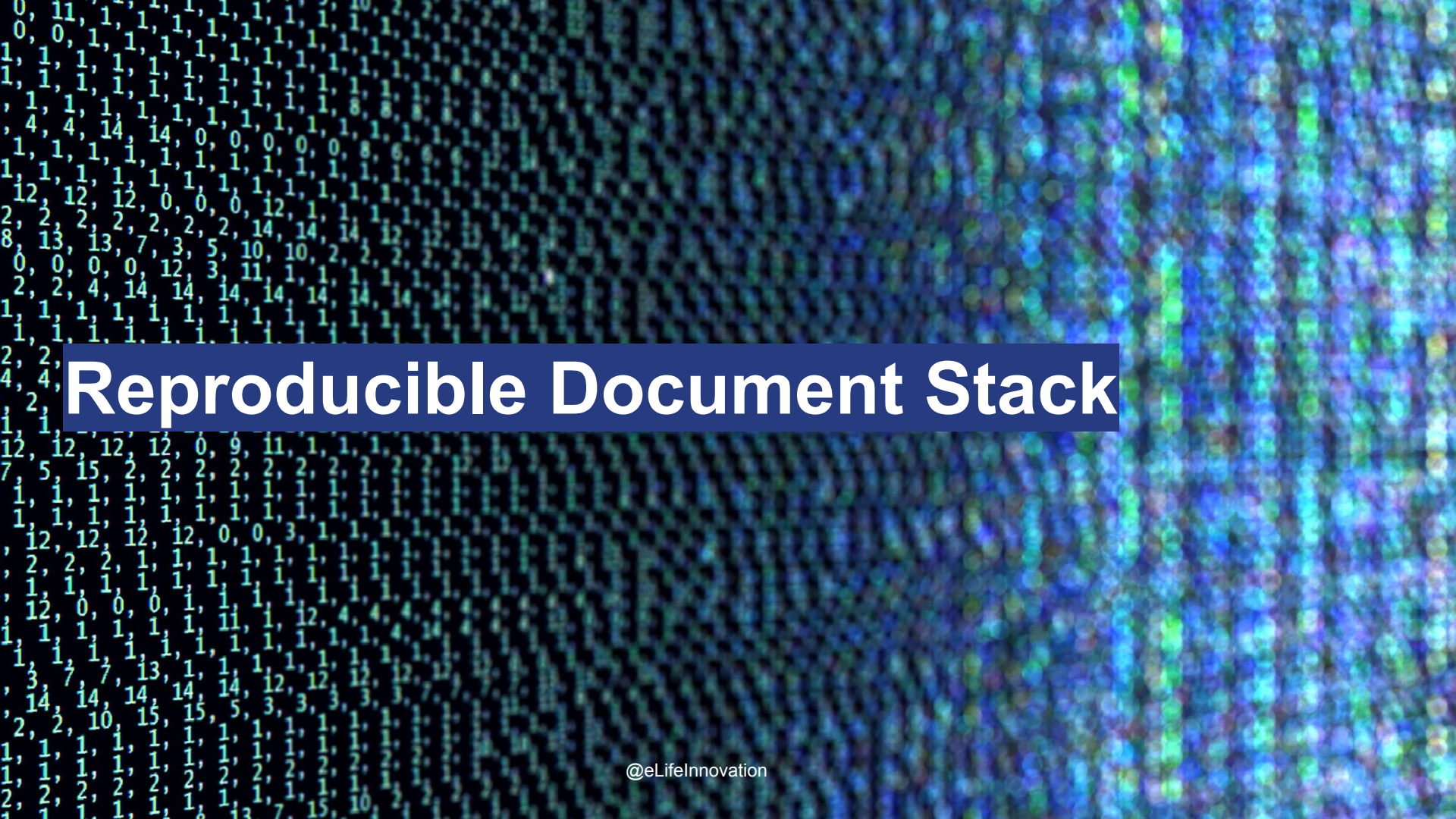
# eLife Innovation

# eLife Innovation's mission

Drive open-source innovation for open science, through:

- Building **open-source tools, platforms and technologies** to improve the ways research is discovered, consumed, shared and evaluated
- Supporting a **community of open-source innovators** to develop these tools





# Reproducible Document Stack

# Is this you?

The screenshot displays the RStudio interface. On the left, the R Markdown source file '1-example.Rmd' is open. A 'Knit' menu is visible, with options: 'Knit to HTML', 'Knit to PDF' (selected), 'Knit to Word', and 'Knit with Parameters...'. The R code in the document includes a library call for 'viridis' and a text block explaining the purpose of the code. The right pane shows the rendered PDF, titled 'Viridis Demo'. The PDF content includes a title, a descriptive paragraph, and two sections: 'Viridis colors' and 'Magma colors', each featuring a contour map of the Maunga Whau volcano. The 'Viridis colors' section shows a plot using the 'viridis(200)' color palette, and the 'Magma colors' section shows a plot using the 'viridis(200, option = "A")' color palette. The PDF viewer shows 'Page: 1 of 2' and 'Automatic Zoom'.

```
1- example.Rmd
1 - 
2 t 
3 o 
4 
5 
6 ```{r include = FALSE}
7 library(viridis)
8 ```
9 
10 The code below demonstrates two color palettes in the
11 [viridis](https://github.com/sjmgarnier/viridis) package.
12 Each plot displays a contour map of the Maunga Whau volcano
13 in Auckland, New Zealand.
14 
15 ## Viridis colors
16 
17 ```{r}
18 image(volcano, col = viridis(200))
19 ```
20 
21 ## Magma colors
22 
23 ```{r}
24 image(volcano, col = viridis(200, option = "A"))
25 ```
```

Viridis Demo

The code below demonstrates two color palettes in the [viridis](https://github.com/sjmgarnier/viridis) package. Each plot displays a contour map of the Maunga Whau volcano in Auckland, New Zealand.

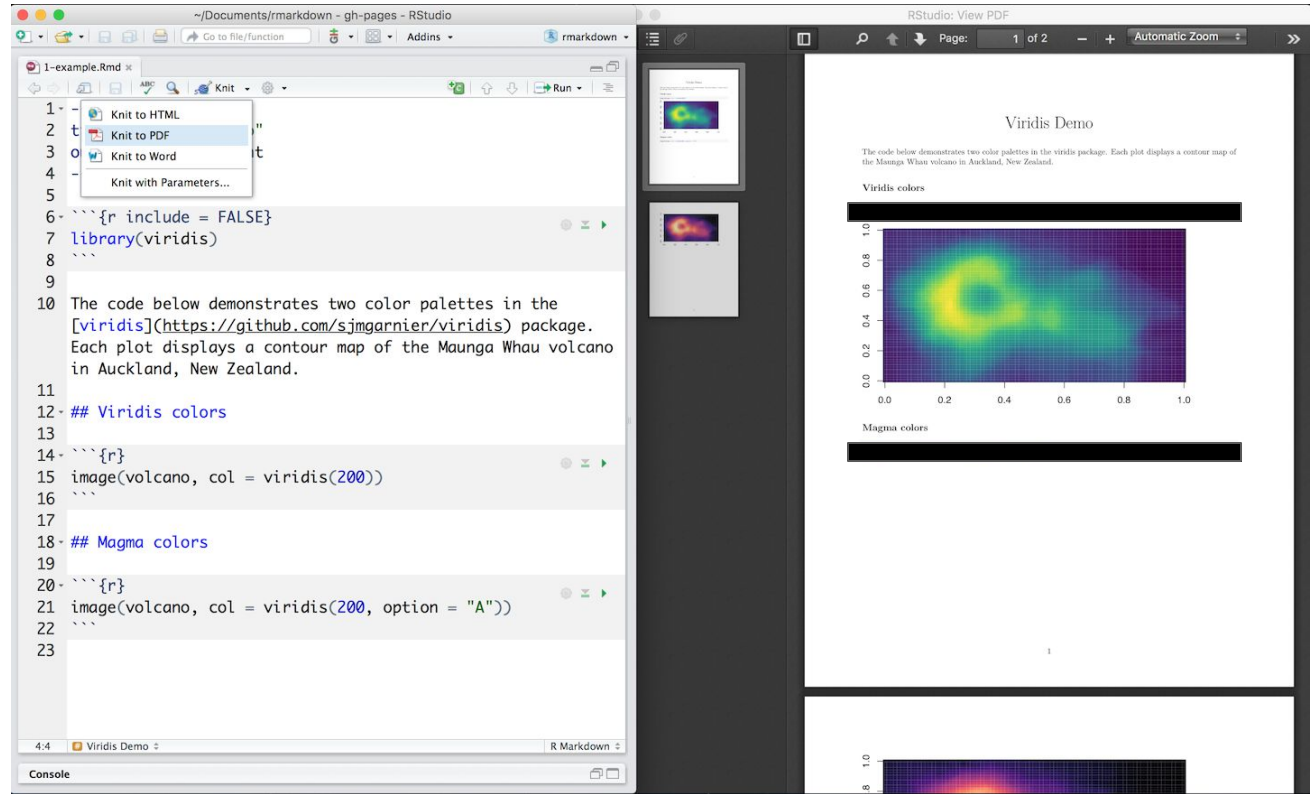
Viridis colors

```
image(volcano, col = viridis(200))
```

Magma colors

```
image(volcano, col = viridis(200, option = "A"))
```

# Is this you?



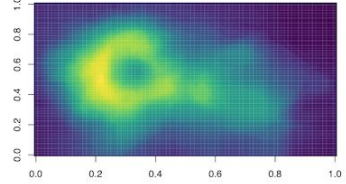
The screenshot displays the RStudio interface. On the left, the R Markdown source file '1-example.Rmd' is open. A context menu is visible over the first few lines of code, with 'Knit to PDF' selected. The code in the editor includes a header, a title, and two R code blocks using the 'viridis' package to generate contour maps of the Maunga Whau volcano. The first block uses the default 'viridis' palette, and the second block uses the 'magma' option. The right pane shows the rendered PDF, titled 'Viridis Demo'. It contains the same header and title, followed by explanatory text and two contour plots. The top plot is labeled 'Viridis colors' and the bottom plot is labeled 'Magma colors'. Both plots show a contour map of the volcano with axes ranging from 0.0 to 1.0. The 'Viridis' plot uses a sequential color palette, while the 'Magma' plot uses a diverging color palette. The PDF viewer shows 'Page: 1 of 2' and 'Automatic Zoom'.

```
1- example.Rmd
1 - 
2 t 
3 O 
4 
5 
6 ```{r include = FALSE}
7 library(viridis)
8 ```
9 
10 The code below demonstrates two color palettes in the
11 [viridis](https://github.com/sjmgarnier/viridis) package.
12 Each plot displays a contour map of the Maunga Whau volcano
13 in Auckland, New Zealand.
14 
15 ## Viridis colors
16 
17 ```{r}
18 image(volcano, col = viridis(200))
19 ```
20 
21 ## Magma colors
22 
23 ```{r}
24 image(volcano, col = viridis(200, option = "A"))
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
```


Viridis Demo

The code below demonstrates two color palettes in the [viridis](https://github.com/sjmgarnier/viridis) package. Each plot displays a contour map of the Maunga Whau volcano in Auckland, New Zealand.

Viridis colors



Magma colors





# Our vision: Reproducible Documents

- Encapsulates usable code and data within the flow of a manuscript.
- Delivers **progressive enhancement** from static research article, to full data and code interaction
- **Future-proof**: Platform, tool, language agnostic
- **Accessible**: Easy and accessible for everyone
- Encourage reuse of published research

# Reproducible Document Stack



Demo: [elifesci.org/reproducible-example](https://elifesci.org/reproducible-example)

# Towards a scalable infrastructure for reproducible document publishing

1. Interoperable authoring and conversion tools
2. Portable reproducible documents
3. Reliable and performant reproducible execution environments
4. Publisher tools



# Authoring and Conversion tools

- **Authoring with Stencila Desktop:** an intuitive, clean text editor built on top of Texture, with code cells and reproducible figures
  - “Mini” formula language for Excel-like graphing
- **Conversion with Encoda:** allow conversion from commonly used formats (e.g. Jupyter notebooks, R Markdown, Google Doc, LaTeX, PDF) to DAR



# Pandoc: pandoc.org

## Pandoc a universal document converter

[Donate](#)[About](#)[Installing](#)[Getting started](#)[Demos ▾](#)[Documentation ▾](#)[Help](#)[Extras](#)[Releases](#)

### About pandoc

If you need to convert files from one markup format into another, pandoc is your swiss-army knife. Pandoc can convert between the following formats:

(← = conversion from; → = conversion to; ↔ = conversion from and to)

#### Lightweight markup formats

- ↔ [Markdown](#) (including [CommonMark](#) and [GitHub-flavored Markdown](#))
- ↔ [reStructuredText](#)
- [AsciiDoc](#)
- ↔ Emacs [Org-Mode](#)
- ↔ Emacs [Muse](#)
- [Textile](#)
- ← [txt2tags](#)

#### HTML formats

- ↔ (X)HTML 4
- ↔ HTML5

#### Ebooks

- ↔ [EPUB](#) version 2 or 3
- ↔ [FictionBook2](#)

#### Documentation formats

- [GNU TexInfo](#)

#### Word processor formats

- ↔ Microsoft Word [docx](#)
- ↔ OpenOffice/LibreOffice [ODT](#)
- [OpenDocument XML](#)
- Microsoft [PowerPoint](#)

#### Page layout formats

- [InDesign ICML](#)

#### Outline formats

- ↔ [OPML](#)

#### Wiki markup formats

- ↔ [MediaWiki markup](#)
- ↔ [DokuWiki markup](#)
- ← [TikiWiki markup](#)
- ← [TWiki markup](#)
- [Vimwiki markup](#)
- [XWiki markup](#)
- [ZimWiki markup](#)



## Pandoc

pandoc sunspots.ipynb -o sunspots.docx

### Getting started

OK, let's just dive right in and fill in details as we go. I'll be using Python for this exploration but will focus on the story and not the code.

First things first, let's load the sunspots data, which is easy to find (e.g. from NOAA) and conveniently included in a popular Python package for doing statistical work...

```
import statsmodels.api as sm
import pandas as pd
data_loader = sm.datasets.sunspots.load_pandas()
df = data_loader.data
```

`df` is shorthand for "dataframe", which we can think of as an Excel-like table of values. Dataframes have various methods that can be called to easily learn about the data contained in them, and we'll step through calling some of these methods. Below, we see that we have 309 pairs of (year, activity) to examine...

```
df
<class 'pandas.core.frame.DataFrame'>
Int64Index: 309 entries, 0 to 308
Data columns:
YEAR      309 non-null values
SUNACTIVITY  309 non-null values
dtypes: float64(2)
```

We can quickly inspect the first and last handful of values to get an idea of what the data look like...

```
df.head()
```

```
YEAR  SUNACTIVITY
0 1700         5
1 1701        11
2 1702        16
3 1703        23
4 1704        36
```

```
df.tail()
```

```
YEAR  SUNACTIVITY
304 2004         40.4
305 2005         29.8
306 2006         15.2
307 2007          7.5
308 2008          2.9
```

## Stencila + Pandoc

stencila convert sunspots.ipynb sunspots.docx

### Getting started

OK, let's just dive right in and fill in details as we go. I'll be using Python for this exploration but will focus on the story and not the code.

First things first, let's load the sunspots data, which is easy to find (e.g. from NOAA) and conveniently included in a popular Python package for doing statistical work...

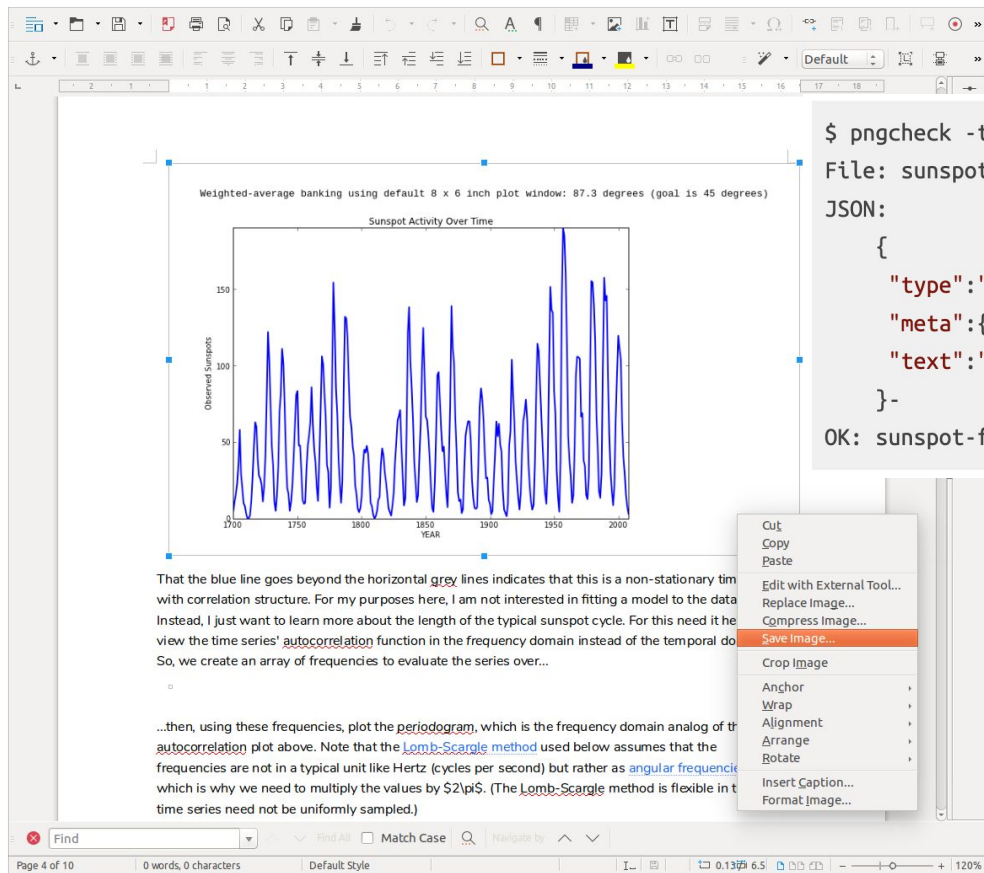
`df` is shorthand for "dataframe", which we can think of as an Excel-like table of values. Dataframes have various methods that can be called to easily learn about the data contained in them, and we'll step through calling some of these methods. Below, we see that we have 309 pairs of (year, activity) to examine...

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 309 entries, 0 to 308
Data columns:
YEAR      309 non-null values
SUNACTIVITY  309 non-null values
dtypes: float64(2)
```

We can quickly inspect the first and last handful of values to get an idea of what the data look like...

```
0 1700 5
1 1701 11
2 1702 16
3 1703 23
4 1704 36
```

# Encode



```
$ pngcheck -tc sunspot-figure.png
```

```
File: sunspot-figure.png (50854 bytes)
```

```
JSON:
```

```
{
  "type": "CodeChunk",
  "meta": {"execution_count": 18},
  "text": "import scipy.optimize as spo\n\ntarget = np.radians(45) \\\n\n"}-
```

```
OK: sunspot-figure.png (832x518, 32-bit RGB+alpha, non-interlaced, 97.1%)
```

# Encoda

## Documents, markup, and notebook formats

Format	Status
Markdown	status alpha
RMarkdown	status alpha
Latex	status alpha
HTML	status alpha
PDF	-
Google Doc	status alpha

## Tabular data and spreadsheet formats

Format	Status
CSV	status alpha
Yaml front matter for CSV <a href="#">CSVY</a>	<a href="#">#25</a>
Excel (.xlsx)	status alpha
OpenDocument Spreadsheet	status alpha
<a href="#">Tabular Data Package</a>	status alpha

# Portable reproducible documents

- [DAR \(Document ARchive\)](#) Container for text, code, data and media assets
- Standard JATS-XML based data formats
- Open format
- Extension to support R Markdown inline code cells to enhance interoperability

# Reproducible execution environments

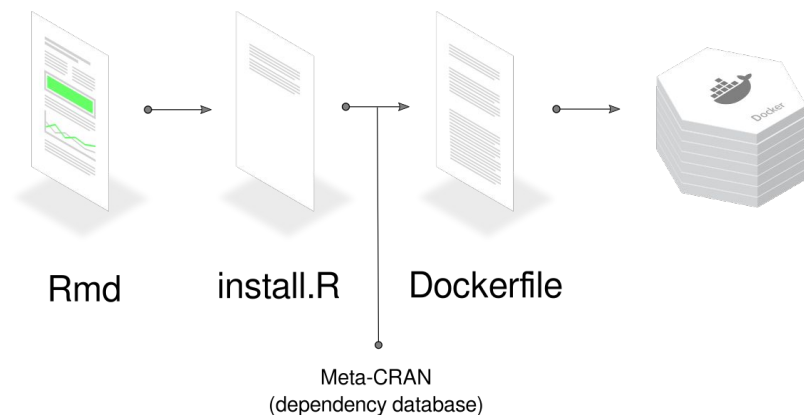
- [Stencila Hub](#) to provide reliant and performant execution environments to run live-code elements
- Building on top of existing technologies: Jupyter Kernels, Binder Hub, Docker



# Dockta: smaller Docker images optimised for reproducible articles

- Performs static code analysis to determine package requirements.
- Uses package databases to determine package system dependencies and generate linked metadata
- Quicker installation of package dependencies

<https://stencila.github.io/dockta/>



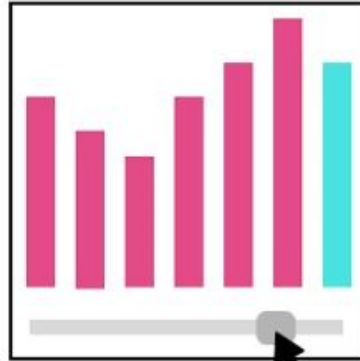
# Publisher tools: Progressive enhancement via multi-level output

1.



Plain .PNG  
(e.g., mobile)

2.



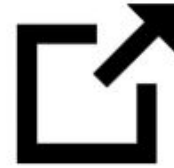
Interactive

3.

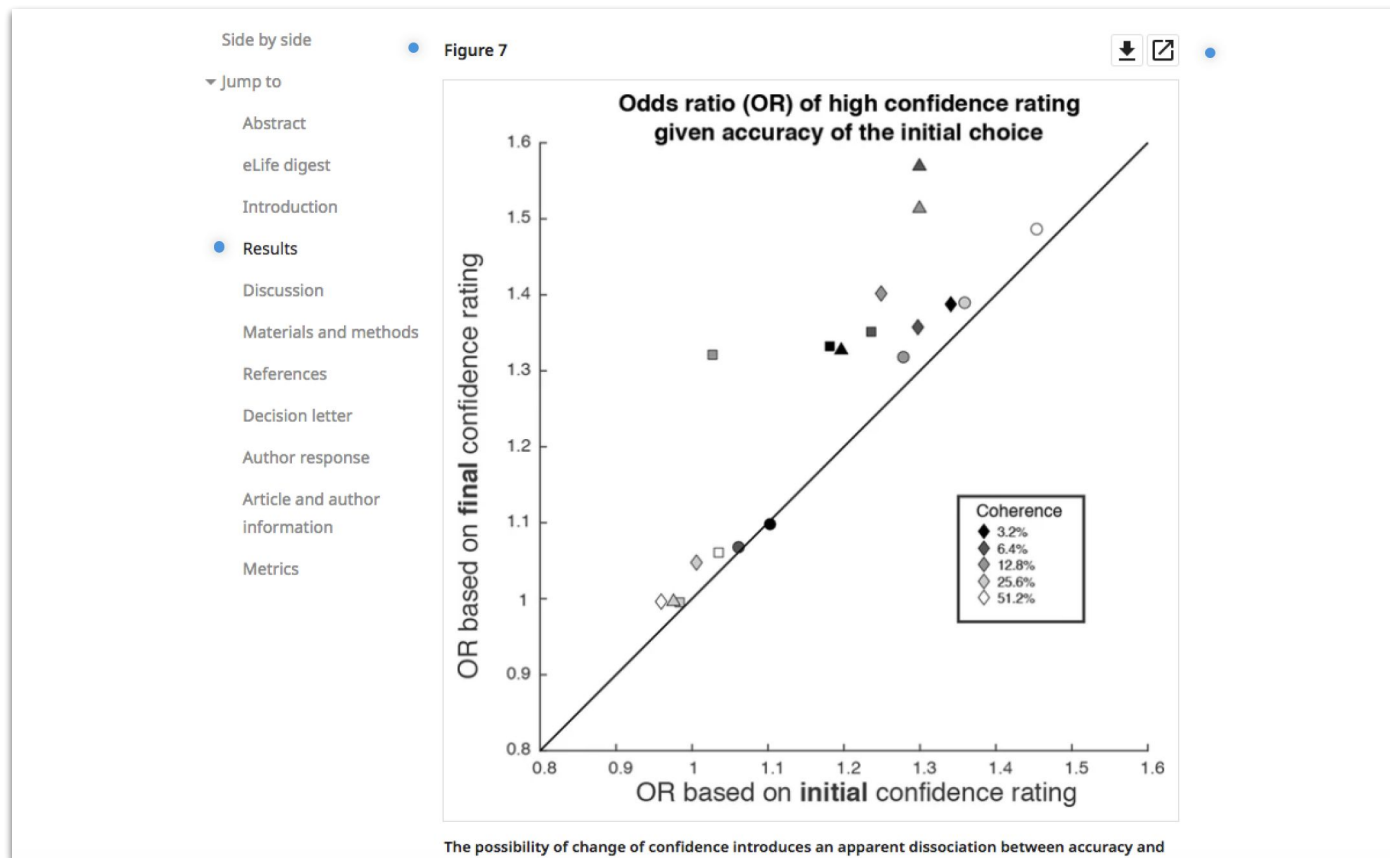


Code / Data view  
(edit, execute)

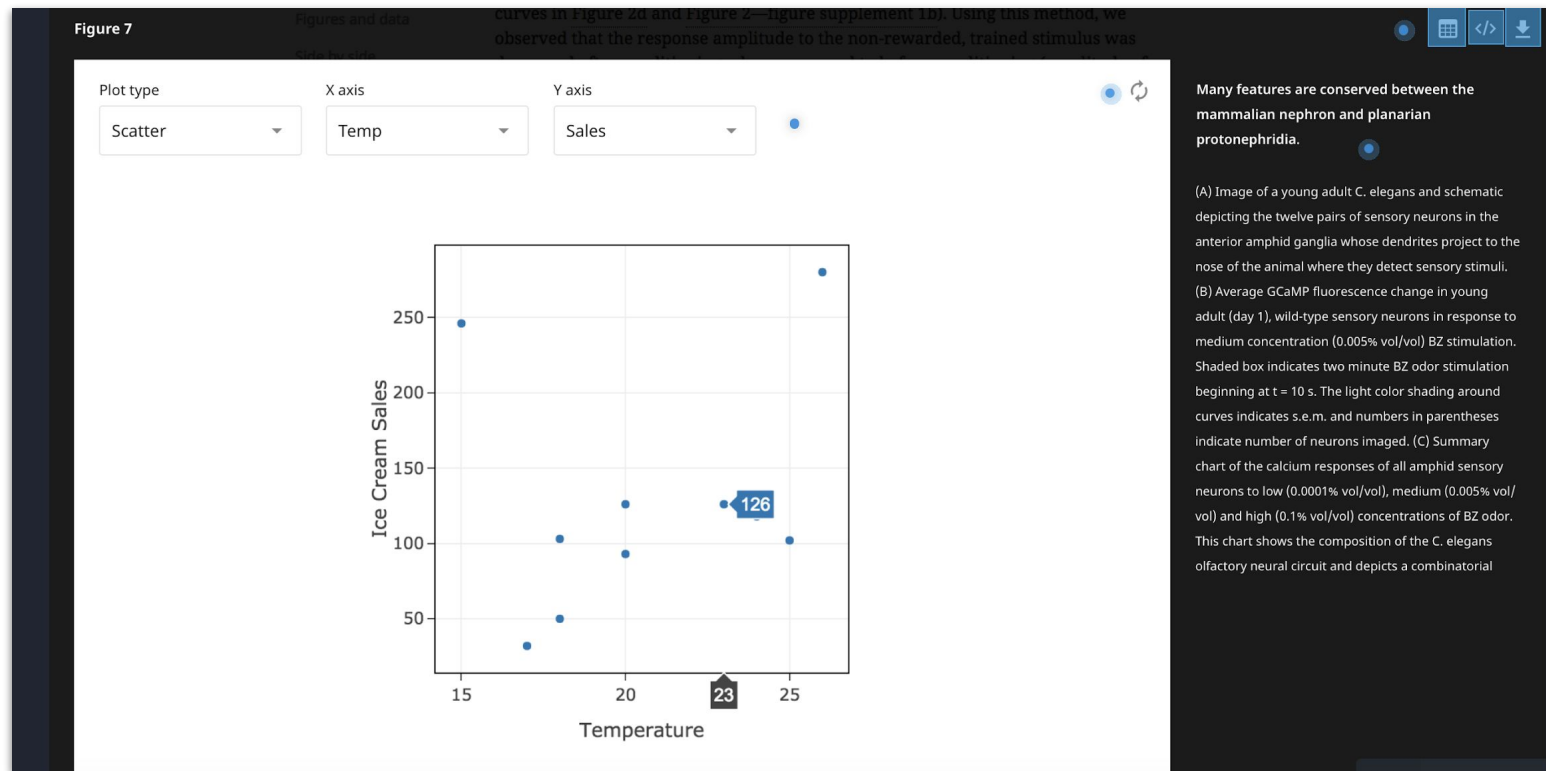
4.



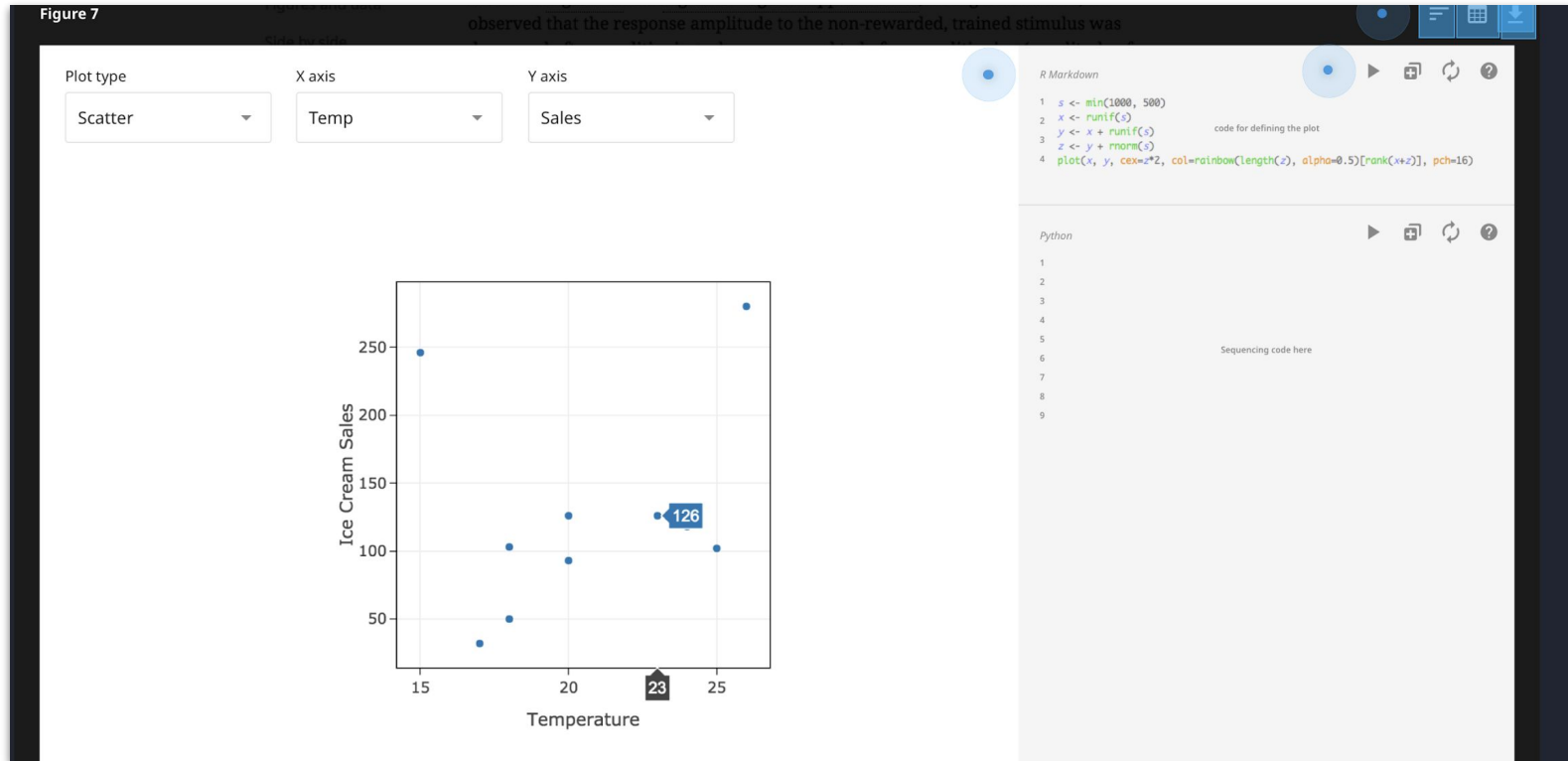
# Casual reader



# Interactive figures



# Interacting with code and data





# Publisher tools

- Web-based publishing of fully reproducible documents with in-browser code interaction and execution
  - Converted to HTML and served from Stencila Hub; or
  - Rendered by Javascript in browser using a Texture Reader interface
- Quick export to PDF for legacy systems
- Journal submission infrastructure integration

# Core development principles

- Open
  - Not trying to “win” a tools race
- Interoperable
  - Easy for scientists to create / publishers to publish reproducible documents from multiple starting points
- Modular
  - Tools within the stack can be taken out and integrated into other pipelines
  - Minimise dependencies for reuse

# You can help

- Share your use case
- Provide feedback
- Learn about progress and opportunities to help

Sign up: [elifesci.org/RDSupdates](https://elifesci.org/RDSupdates)

This will take you to a form asking for your consent to be added to a mailing list for ~bimonthly emails with updates about this project, including calls for contributions and feedback.

“Alice saw this deep learning model online and would like to try it on her data. To first ensure that she can run the model she decided to download the original dataset and apply the model to check if she can get the same results. When she tried to run the code cloned from Github, error messages were appearing one after another, suggesting that she has the wrong software dependencies and system paths.”

Are you Alice? Join our [Birds of a Feather at BOSC](#)  
Wednesday July 24, 4:40pm

# Questions?

Email: [e.tsang@elifesciences.org](mailto:e.tsang@elifesciences.org)

Twitter: [@eLifeInnovation](https://twitter.com/eLifeInnovation) / [@emmy\\_ft](https://twitter.com/emmy_ft)

Labs: [elifesci.org/labs](http://elifesci.org/labs)

Stencila: [stenci.la](http://stenci.la)

Substance: [substance.io](http://substance.io)



[doi.org/10.6084/m9.figshare.8983640](https://doi.org/10.6084/m9.figshare.8983640)



These slides are licensed for reuse under  
Creative Commons Attribution License with  
attribution to eLife

External content retained under individual  
licenses as indicated on slides