# Agenda
## Boost your reproducibility with Binder

– 13:30    Registration and introductions

– 13:40    Introduction to the workshop and The Turing Way

– 13:50    Presentation: Why you need a reproducible computing
                     environment and how Binder can help

– 15:00    Coffee break

– 15:30    Code along demo: Zero to Binder, build a Binder resource

– 16:30    Build your own Binder

– 16:50    Feedback, group picture and close

# The Alan Turing Institute

---

## The Turing Way
## Boost your reproducibility with Binder workshop

## Kirstie Whitaker
Pronouns: she/her

# Agenda

## Boost your reproducibility with Binder

- 13:30    Registration and introductions

- 13:40    Introduction to the workshop and The Turing Way

- 13:50    Presentation: Why you need a reproducible computing environment and how Binder can help

- 15:00    Coffee break

- 15:30    Code along demo: Zero to Binder, build a Binder resource

- 16:30    Build your own Binder

- 16:50    Feedback, group picture and close

# The Turing Way is:

- a book
- a community
- a global collaboration
- a whole tonne of work



Rachael Ainsworth
Becky Arnold
Louise Bowler
Sarah Gibson
Patricia Herterich
James Hetherington
Rosie Higman
Anna Krystalli
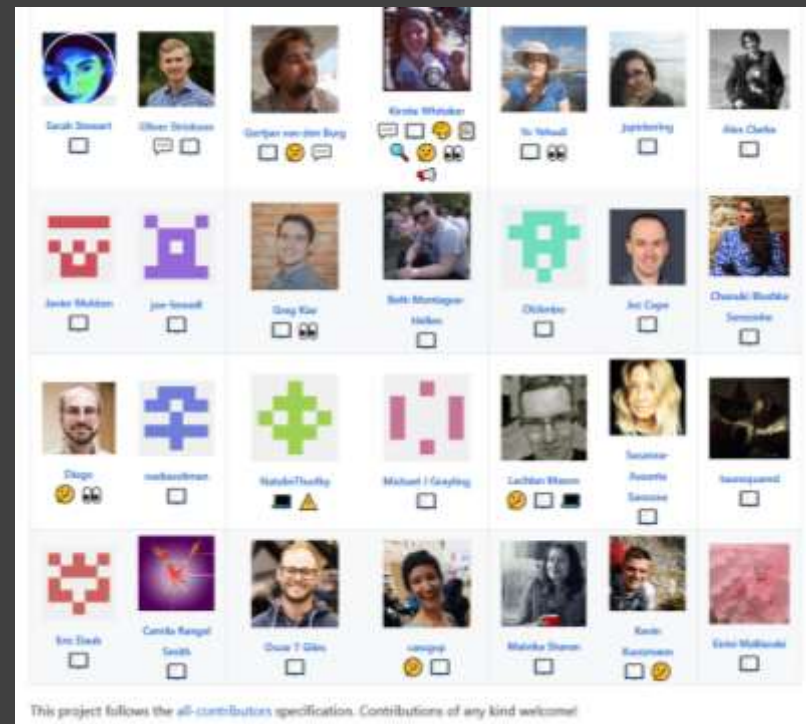Catherine Lawrence
Alex Morley
Martin O'Reilly
Malvika Sharan

# Thank you to all our contributors

# The Turing Institute

https://www.turing.ac.uk/news/enigma-machine-goes-display-alan-turing-institute
#TuringWay @kirstie_j @mybinderteam
https://doi.org/10.5281/zenodo.3632909

# University network

# The Institute's partners and collaborators

# Challenges

Advance data science and artificial intelligence to…


Revolutionise healthcare


Deliver safer, smarter engineering


Manage security in an insecure world


Shine a light on our economy


Make algorithmic systems fair, transparent, and ethical


Design computers for the next generation of algorithms


Supercharge research in science and humanities


Foster government innovation

The Alan Turing Institute

Home + News

# The Alan Turing Institute to spearhead new cutting-edge data science and AI research after £48 million government funding boost

Tuesday 18 Dec 2018

Learn more ↓

#TuringWay @kirstie_j @mybinderteam
https://doi.org/10.5281/zenodo.3632909

# The Turing Way

# Welcome to the Turing Way

The Turing Way is a lightly opinionated guide to reproducible data science.

Our goal is to provide all the information that researchers need at the start of their projects to ensure that they are easy to reproduce at the end.

This also means making sure PhD students, postdocs, PIs and funding teams know which parts of the "responsibility of reproducibility" they can affect, and what they should do to nudge data science to being more efficient, effective and understandable.

### A bit more background

Reproducible research is necessary to ensure that scientific work can be trusted. Funders and publishers are beginning to require that publications include access to the underlying data and the analysis code. The goal is to ensure that all results can be independently verified and built upon in future work. This is sometimes easier said than done. Sharing these research outputs means understanding data management, library sciences, software development, and continuous integration techniques: skills that are not widely taught or expected of academic researchers and data scientists.

The Turing Way is a handbook to support students, their supervisors, funders and journal editors

https://the-turing-way.netlify.com/introduction/introduction
#TuringWay @kirstie_j @mybinderteam
https://doi.org/10.5281/zenodo.3632909

## Welcome to the Turing Way

The Turing Way is a lightly opinionated guide to reproducible data science.

Our goal is to provide all the information that researchers need at the start of their projects to ensure that they are easy to reproduce at the end.

This also means making sure PhD students, postdocs, PIs and funding teams know which parts of the "responsibility of reproducibility" they can affect, and what they should do to nudge data science to being more efficient, effective and understandable.

### A bit more background

Reproducible research is necessary to ensure that scientific work can be trusted. Funders and publishers are beginning to require that publications include access to the underlying data and the analysis code. The goal is to ensure that all results can be independently verified and built upon in future work. This is sometimes easier said than done. Sharing these research outputs means understanding data management, library sciences, software development, and continuous integration techniques: skills that are not widely taught or expected of academic researchers and data scientists.

The Turing Way is a handbook to support students, their supervisors, funders and journal editors

# Welcome to the Turing Way

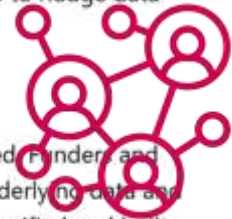The Turing Way is a lightly opinionated guide to reproducible data science.

Our goal is to provide all the information that researchers need at the start of their projects to ensure that they are easy to reproduce at the end.

This also means making sure PhD students, postdocs, PIs and funding teams know which parts of the "responsibility of reproducibility" they can affect, and what they should do to nudge data science to being more efficient, effective and understandable.

**A bit more background**

Reproducible research is necessary to ensure that scientific work can be trusted. Funders and publishers are beginning to require that publications include access to the underlying data and the analysis code. The goal is to ensure that all results can be independently verified and built upon in future work. This is sometimes easier said than done. Sharing these research outputs means understanding data management, library sciences, software development, and continuous integration techniques: skills that are not widely taught or expected of academic researchers and data scientists.

The Turing Way is a handbook to support students, their supervisors, funders and journal editors

https://the-turing-way.netlify.com/introduction/introduction
#TuringWay @kirstie_j @mybinderteam
https://doi.org/10.5281/zenodo.3632909

|  | Data | |
|---|---|---|
|  | Same | Different |
| **Analysis** Same | Reproducible | Replicable |
| Different | Robust | Generalisable |

Barriers to reproducible research

Held to higher standards than others

Is not considered for promotion

Publication bias towards novel findings

Requires additional skills

Plead the 5th

Support additional users

Takes time

https://doi.org/10.6084/m9.figshare.5537101
#TuringWay @kirstie_j @mybinderteam
https://doi.org/10.5281/zenodo.3632909

Barriers to reproducible research

Held to higher standards than others

Is not considered for promotion
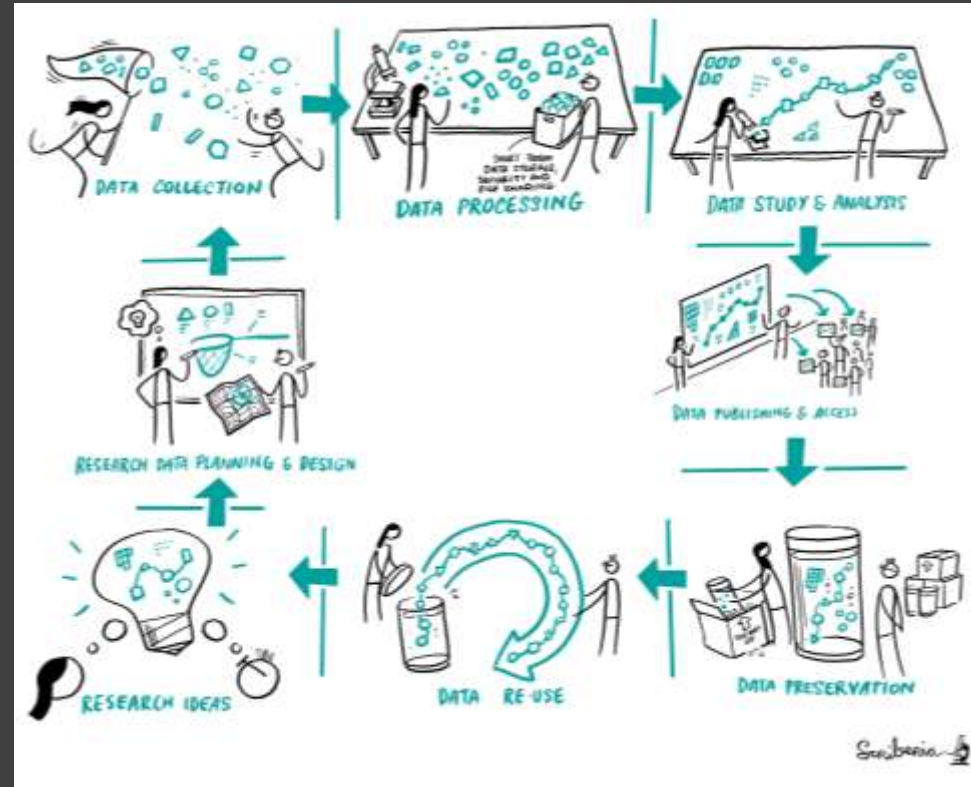
Publication bias towards novel findings

Requires additional skills

Plead the 5th

Support additional users

Takes time

To be fully reproducible we have to cover all the steps of the research cycle

And that is super overwhelming…but we're here to help

FAIR PRINCIPLES

PERSISTANT — FINDABLE
MEANINGFUL INTERACTION — ACCESSIBLE
INTEROPERABLE
FULL DISCLOSURE — REUSABLE
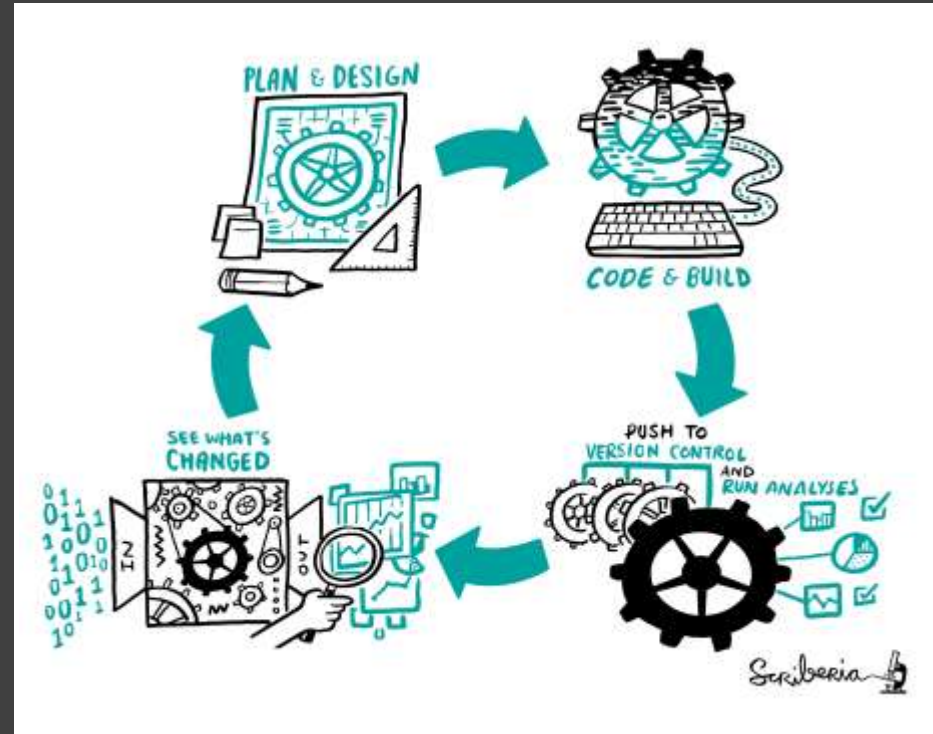
# Testing for research

```
Assert.AreEqual(

    GetTimeOfDay(),

    "Morning" )
```
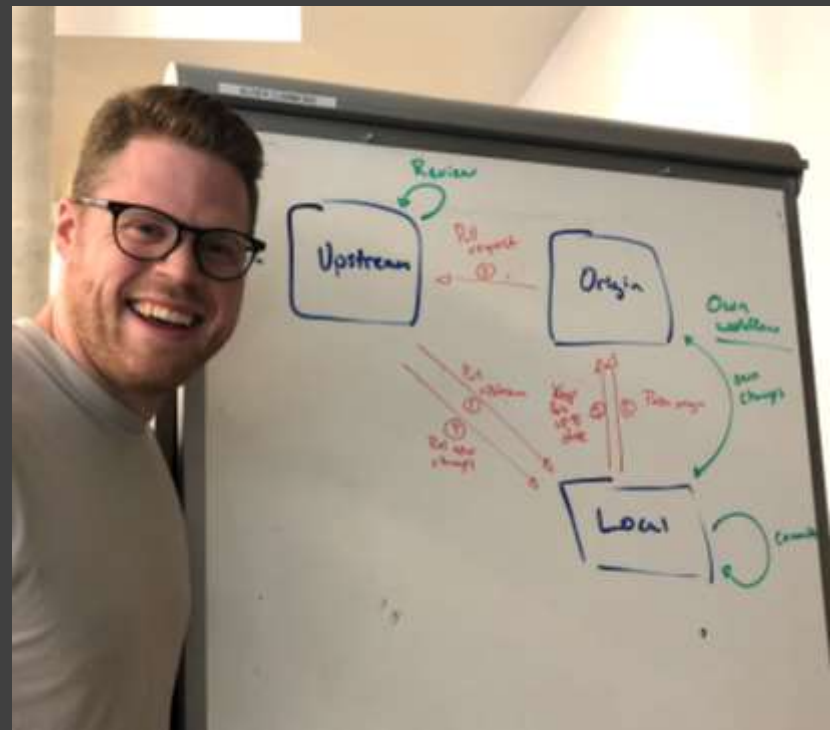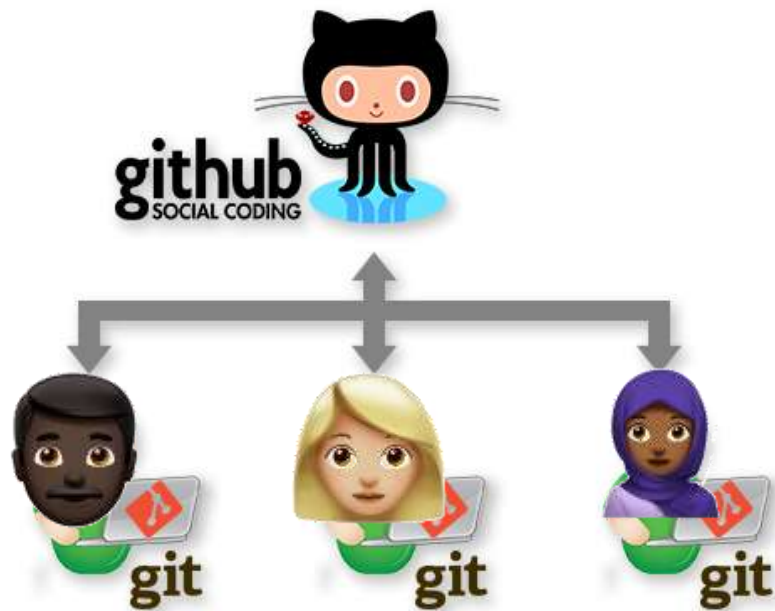
# Continuous integration

– Plan and design your experiment

– Write down the steps in code

– Push to version control and run the analyses

– Test to see what's changed

https://the-turing-way.netlify.com/collaborating_github/collaborating_github.html
https://the-turing-way.netlify.com/version_control/version_control.html    #TuringWay @kirstie_j @mybinderteam
https://neurohackademy.org                                                  https://doi.org/10.5281/zenodo.3632909

# Extension in 2020

– Expand scope to all data science practices

  – Reproducibility

  – Scoping and designing a data science project

  – Ethics

  – Communication and visualisation

  – Collaborative working



https://github.com/
    alan-turing-institute/the-turing-way/
    blob/master/project_management/
    tps-funding-application-20190429.md
        #TuringWay @kirstie_j @mybinderteam
        https://doi.org/10.5281/zenodo.3632909

# A global collaboration

# Patricia Herterich

"What really sets The Turing Way apart is HOW we're writing the book. The focus on community, the commitment to transparency and working open right from the beginning is an exciting (and terrifying) new way of working."

# Open Leadership Principles



**Understanding**
You make the work accessible and clear

**Sharing**
You make the work easy to adapt, reproduce, and spread

**Participation & Inclusion**
You build shared ownership and agency to make the work inviting and sustainable for all.

**Read more**
https://mozilla.github.io/olm-whitepaper

moz://a

#TuringWay @kirstie_j @mybinderteam
https://doi.org/10.5281/zenodo.3632909

# Agenda

## Boost your reproducibility with Binder

– 13:30    Registration and introductions

– 13:40    Introduction to the workshop and The Turing Way

– 13:50    Presentation: Why you need a reproducible computing
          environment and how Binder can help

– 15:00    Coffee break

– 15:30    Code along demo: Zero to Binder, build a Binder resource

– 16:30    Build your own Binder

– 16:50    Feedback, group picture and close

# Goals for the workshop

– Understand how your computational environment impacts reproducibility

– Learn what Binder is and how it can help make your research reproducible

– Build your own Binder!

# Our Code of Conduct

"The Turing Way team are dedicated to providing a welcoming and supportive environment for all people…we do not tolerate behaviour that is disrespectful to our community members or that excludes, intimidates, or causes discomfort to others."

– Be respectful of different viewpoints and experiences.

– Use welcoming and inclusive language.

– Do not harass people.

– Respect the privacy and safety of others.

   Please do not take pictures of anyone without their permission.

– Be considerate of others' participation.

– Don't be a bystander.

– Be respectful of different viewpoints and experiences.

– Use welcoming and inclusive language.

– Do not harass people.

– Respect the privacy and safety of others.

Please do not take pictures of anyone without their permission.

– Be considerate of others' participation.

– Don't be a bystander.

– **Anita**, **Felix** and **Jeremy** are here to help 👋

# Thank you to current (& future) contributors



https://github.com/alan-turing-institute/the-turing-way#contributors
https://allcontributors.org/docs/en/emoji-key

#TuringWay @kirstie_j @mybinderteam
https://doi.org/10.5281/zenodo.3632909

# Thank you

– Book: https://the-turing-way.netlify.com

– Newsletter: https://tinyletter.com/TuringWay

– GitHub: https://github.com/alan-turing-institute/the-turing-way

– Chat: https://gitter.im/alan-turing-institute/the-turing-way

– Unsplash photos by Adolfo Felix, James Pond, Jose Alejandro Cuffia, Kinson Leung, Mateo Vrbnjak, Mimi thian, Omar Albeik, Perry Grone, Toa Heftiba, Tomasz Frankows, Wilmer Martinez

– Noun Project icons by Aybige, Luis Prado, Edward Boatman, Becris, Rose Alice Design, Hyemm.work

– Original artwork by Scriberia: https://doi.org/10.5281/zenodo.3332807

#TuringWay @kirstie_j @mybinderteam

https://doi.org/10.5281/zenodo.3632909

# Agenda

## Boost your reproducibility with Binder

– 13:30    Registration and introductions

– 13:40    Introduction to the workshop and The Turing Way

– 13:50    Presentation: Why you need a reproducible computing
environment and how Binder can help

– 15:00    Coffee break

– 15:30    Code along demo: Zero to Binder, build a Binder resource

– 16:30    Build your own Binder

– 16:50    Feedback, group picture and close

# A note on the name

– I never thought the name would be approved!

– This is <u>not</u> a Turing project (although it has great support from the Institute)

– We are creating guidance together, <u>the way</u> is a journey not a set of rules

# Turing Way & Binder

Courtesy of Juliette Taka: https://twitter.com/mybinderteam/status/1082556317842264064
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

– Coordinate cloud computing resources with Kubernetes (k8s)

– Make it easy for users to access with a JupyterHub

– Set up the environment from your GitHub repository

repo2docker  Jupyter  kubernetes

Google Cloud

https://binderhub.readthedocs.io
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

# Sarah Gibson

"It took me a while to feel like I knew enough to contribute to Binder. But the team are always so excited to have my input. Its really motivating to be part of such a welcoming community."

– Check analysis on my phone

– Share the responsibility with busy PIs

– Requires version control, capturing environment and new build for each change

https://mybinder.readthedocs.io/en/latest/faq.html#how-much-does-running-mybinder-org-cost
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

https://mybinder.readthedocs.io/en/latest/faq.html#how-much-does-running-mybinder-org-cost
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909
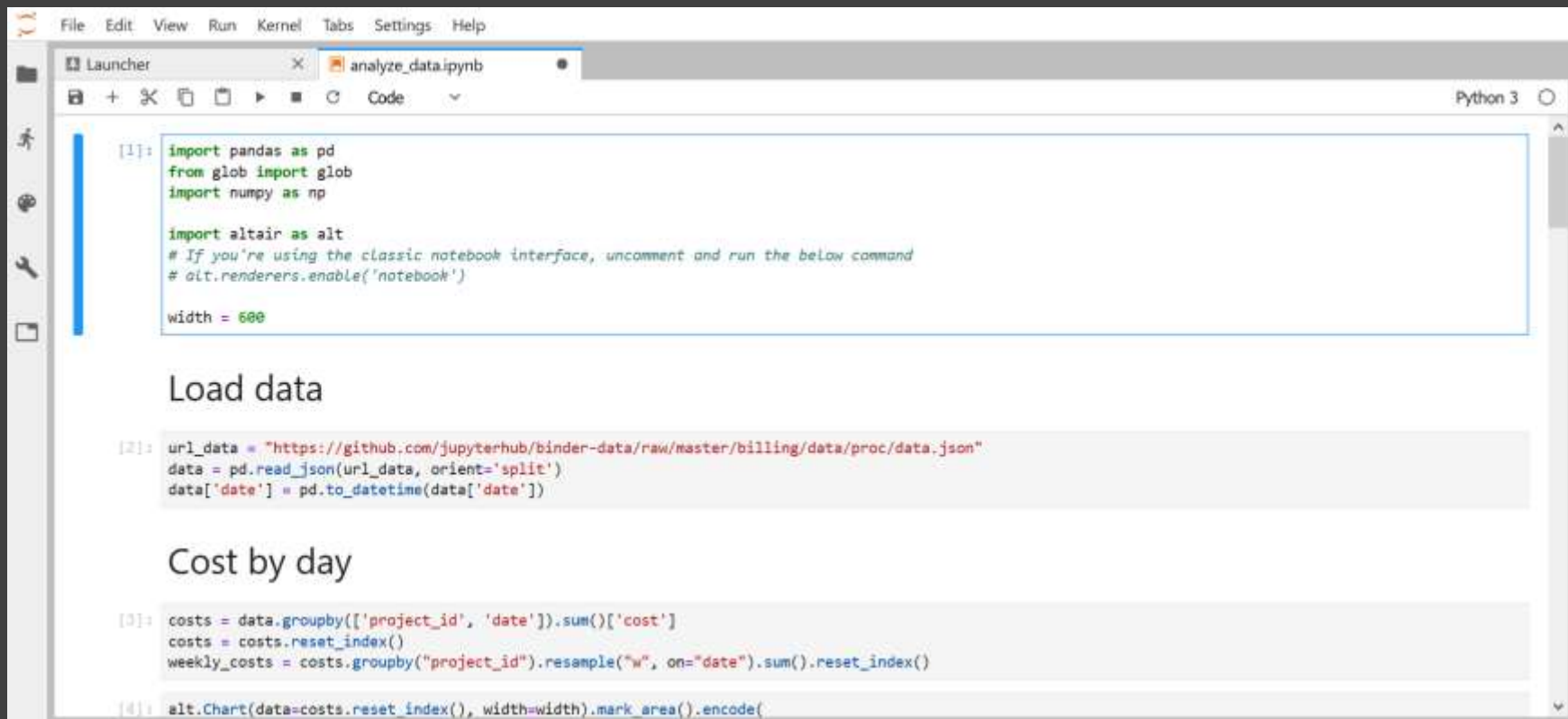
https://mybinder.readthedocs.io/en/latest/faq.html#how-much-does-running-mybinder-org-cost
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

Courtesy of Juliette Taka: https://twitter.com/mybinderteam/status/1082556317842264064
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

# Champion: Elena Kochkina

Turing at SemEval-2017 Task 8: Sequential Approach to Rumour Stance Classification with Branch-LSTM

Elena Kochkina, Maria Liakata, Isabelle Augenstein

**Source tweet**

**Support**
**Deny**
**Query**
**Comment**

... SDQC ... SDQC

France: 10 people dead after shooting at HQ of satirical weekly newspaper #CharlieHebdo, according to witnesses [link]

... SDQC

| Label \ Prediction | C | D | Q | S |
|---|---|---|---|---|
| Commenting | 760 | 0 | 12 | 6 |
| Denying | 68 | 0 | 1 | 2 |
| Querying | 69 | 0 | 36 | 1 |
| Supporting | 67 | 0 | 1 | 26 |

Table 5: Confusion matrix for testing set predictions

https://github.com/kochkinaelena/branchLSTM (on Turing Way Hub)
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

File   Edit   View   Run   Kernel   Tabs   Settings   Help

Console 1

```
| Time: | 1341.7773336343318|

[3]: %run depth_analysis.py

     trials.txt is not available


     --- Table 4 ---

     Number of tweets per depth and performance at each of the depths

     Depth      # tweets    # Support    # Deny    # Query    # Comment    Accuracy    MacroF    Support    Deny    Query
     Comment
     0          28          26           2         0          0            0.929       0.481     0.963      0.000   0.000
     0.000
     1          704         61           60        81         502          0.696       0.436     0.192      0.088   0.660
     0.806
     2          128         3            6         7          112          0.805       0.318     0.000      0.000   0.385
     0.887
     3          60          2            1         5          52           0.817       0.307     0.000      0.000   0.333
     0.895
     4          41          0            0         3          38           0.927       0.481     0.000      0.000   0.000
     0.962
     5          27          1            0         1          25           0.926       0.321     0.000      0.000   0.000
     0.962
     6+         61          1            2         9          49           0.803       0.223     0.000      0.000   0.000
     0.891


     --- Table 5 ---

     Confusion matrix

     Lab \ Pred   Comment    Deny    Query    Support
     Comment      667        5       62       44
     Deny         58         3       4        6
     Query        38         0       72       4
     Support      52         0       4        38
```

Name — Last Modified

| Name | Last Modified |
| --- | --- |
| dev_data | 15 days ago |
| downloaded_data | 16 days ago |
| output | 8 minutes ago |
| saved_data | 33 minutes ago |
| scorer | 15 days ago |
| src | 15 days ago |
| tokenizers | 35 minutes ago |
| badwords.txt | 15 days ago |
| bestparams_GN.txt | 15 days ago |
| depth_analysis.py | 15 days ago |
| environment.yml | 15 days ago |
| LICENSE | 16 days ago |
| outer.py | 15 days ago |
| postBuild | 15 days ago |
| predict.py | 15 days ago |
| preprocessing.py | 15 days ago |
| README.md | 15 days ago |
| requirements.txt | 15 days ago |
| subtaska.json | 15 days ago |
| subtaskb.json | 15 days ago |
| training.py | 15 days ago |

--- Table 5 ---

Confusion matrix

| Lab \ Pred | Comment | Deny | Query | Support |
|------------|---------|------|-------|---------|
| Comment | 667 | 5 | 62 | 44 |
| Deny | 58 | 3 | 4 | 6 |
| Query | 30 | 0 | 72 | 4 |
| Support | 52 | 0 | 4 | 38 |

--- Table 3 ---

Part 1: Results on testing set

Accuracy = 0.743565300286

Macro-average:
Precision  0.530
Recall     0.496
F-score    0.477
Support    —

Per-class:

|           | Comment | Deny  | Query | Support |
|-----------|---------|-------|-------|---------|
| Precision | 0.827   | 0.375 | 0.507 | 0.413   |
| Recall    | 0.857   | 0.042 | 0.679 | 0.404   |
| F-score   | 0.842   | 0.076 | 0.581 | 0.409   |
| Support   | 778     | 71    | 106   | 94      |

Part 2: Results on development set

As presented in the paper:

|         | Accuracy | Macro-F | Comment | Deny  | Query | Support |
|---------|----------|---------|---------|-------|-------|---------|
| Testing | 0.744    | 0.477   | 0.842   | 0.076 | 0.581 | 0.409   |

Could not find trials.txt; unable to generate results for development set in Table 3.

https://github.com/kochkinaelena/branchLSTM (on Turing Way Hub)
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

# Elena Kochkina

"How would I have known that it would be different on a different machine?! I only have access to the university HPC to run deep learning analyses."

# Gertjan van den Burg

"The fun part of data science is the modelling. Being able to read in information from a csv file should not be the hardest part."

– https://github.com/
alan-turing-institute/
CleverCSVDemo

– https://github.com/ alan-turing-institute/ CleverCSVDemo

– "Wrangling Messy CSV Files by Detecting Row and Type Patterns" arXiv:1811.11242

File　Edit　View　Insert　Cell　Kernel　Widgets　Help

Trusted　Python 3 ○

Markdown

# CSV dialect detection with CleverCSV

**Author**: [Gertjan van den Burg](#)

In this note we'll show some examples of using CleverCSV, a package for handling messy CSV files. We'll start with a motivating example and then show some other files where CleverCSV shines. CleverCSV was developed as part of a research project on automating data wrangling. It achieves an accuracy of 97% on over 9300 real-world CSV files and improves the accuracy on messy files by 21% over standard tools.

Handy links:

- [Paper on arXiv](#)
- [CleverCSV on GitHub](#)
- [CleverCSV on PyPI](#)
- [Reproducible Research Repo](#)

## IMDB Movie data

Alice is a data scientist who would like to analyse the movie ratings on IMDB for movies of different genres. She found [a dataset shared by a user on Kaggle](#) that contains information of over 14,000 movies. Great!

The data is stored in a CSV file, which is a very common data format for sharing tabular data. The first few lines of the file look like this:

🖫  ✚  ✂  ⎘  ⎗  ↑  ↓  ▶ Run  ■  C  ▶▶  | Markdown ▾ |  ⌨

# IMDB Movie data

Alice is a data scientist who would like to analyse the movie ratings on IMDB for movies of different genres. She found a dataset shared by a user on Kaggle that contains information of over 14,000 movies. Great!

The data is stored in a CSV file, which is a very common data format for sharing tabular data. The first few lines of the file look like this:

```
fn,tid,title,wordsInTitle,url,imdbRating,ratingCount,duration,year,type,nrOfWins,nrOfNominations,nrOfPhotos,nrOf
NewsArticles,nrOfUserReviews,nrOfGenre,Action,Adult,Adventure,Animation,Biography,Comedy,Crime,Documentary,Drama
,Family,Fantasy,FilmNoir,GameShow,History,Horror,Music,Musical,Mystery,News,RealityTV,Romance,SciFi,Short,Sport,
TalkShow,Thriller,War,Western
titles01/tt0012349,tt0012349,Der Vagabund und das Kind (1921),der vagabund und das kind,http://www.imdb.com/titl
e/tt0012349/,8.4,40550,3240,1921,video.movie,1,0,19,96,85,3,0,0,0,0,0,1,0,0,1,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
0,0
titles01/tt0015864,tt0015864,Goldrausch (1925),goldrausch,http://www.imdb.com/title/tt0015864/,8.3,45319,5700,19
25,video.movie,2,1,35,110,122,3,0,0,1,0,0,1,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0
titles01/tt0017136,tt0017136,Metropolis (1927),metropolis,http://www.imdb.com/title/tt0017136/,8.4,81007,9180,19
27,video.movie,3,4,67,428,376,2,0,0,0,0,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,1,0,0,0,0,0,0
titles01/tt0017925,tt0017925,Der General (1926),der general,http://www.imdb.com/title/tt0017925/,8.3,37521,6420,
1926,video.movie,1,1,53,123,219,3,1,0,1,0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0
titles01/tt0021749,tt0021749,Lichter der Großstadt (1931),lichter der gro stadt,http://www.imdb.com/title/tt0021
749/,8.7,70057,5220,1931,video.movie,2,0,38,187,186,3,0,0,0,0,0,1,0,0,1,0,0,0,0,0,0,0,0,0,0,0,1,0,0,0,0,0,0
```

Seems pretty standard, let's load it with Pandas!

In [1]: %xmode Minimal

In [1]:
```
%xmode Minimal
import pandas as pd
df = pd.read_csv('./data/imdb.csv')
```

Exception reporting mode: Minimal

ParserError: Error tokenizing data. C error: Expected 44 fields in line 66, saw 46

Oh, that doesn't work. Maybe there's something wrong with the file? Let's try opening it with the Python CSV reader:

In [2]:
```
import csv
with open('./data/imdb.csv', 'r', newline='') as fid:
    dialect = csv.Sniffer().sniff(fid.read())
    print("Detected delimiter = %r, quotechar = %r" % (dialect.delimiter, dialect.quotechar))
    fid.seek(0)
    reader = csv.reader(fid, dialect=dialect)
    rows = list(reader)

print("Loaded %i rows." % len(rows))
```

Detected delimiter = ' ', quotechar = '"'
Loaded 13928 rows.

Huh, that's strange, Python thinks the *space* is the delimiter and loads 13928 rows, but the file should contain 14,762 rows according to the documentation. What's going on here?

https://github.com/alan-turing-institute/CleverCSVDemo
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

Markdown

Huh, that's strange, Python thinks the *space* is the delimiter and loads 13928 rows, but the file should contain 14,762 rows according to the documentation. What's going on here?

It turns out that on the 65th line of the file, there's a movie with the title `Dr. Seltsam\, oder wie ich lernte\, die Bombe zu lieben (1964)` (the German version of Dr. Strangelove). The title has commas in it, that are escaped using the `\` character! Why are CSV files so hard? 😩

**CleverCSV to the rescue!**

CleverCSV detects the dialect of CSV files much more accurately than existing approaches, and it is therefore robust against these kinds of format variations. It even has a wrapper that works with DataFrames!

```python
In [3]: from ccsv.wrappers import csv2df

df = csv2df('./data/imdb.csv')
df
```

Out[3]:

| | fn | tid | title | wordsInTitle | url | imdbRating | ratingCount | duration | year | type | ... | News |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | titles01/tt0012349 | tt0012349 | Der Vagabund und das Kind (1921) | der vagabund und das kind | http://www.imdb.com/title/tt0012349/ | 8.4 | 40550.0 | 3240.0 | 1921.0 | video.movie | ... | 0 |
| 1 | titles01/tt0015864 | tt0015864 | Goldrausch (1925) | goldrausch | http://www.imdb.com/title/tt0015864/ | 8.3 | 45319.0 | 5700.0 | 1925.0 | video.movie | ... | 0 |
| 2 | titles01/tt0017136 | tt0017136 | Metropolis (1927) | metropolis | http://www.imdb.com/title/tt0017136/ | 8.4 | 81007.0 | 9180.0 | 1927.0 | video.movie | ... | 0 |
| 3 | titles01/tt0017925 | tt0017925 | Der General (1926) | der general | http://www.imdb.com/title/tt0017925/ | 8.3 | 37521.0 | 6420.0 | 1926.0 | video.movie | ... | 0 |
| | | | Lichter der | lichter der gro | http://www.imdb.com | | | | | | | |

Markdown

Huh, that's strange, Python thinks the *space* is the delimiter and loads 13928 rows, but the file should contain 14,762 rows according to the documentation. What's going on here?

It turns out that on the 65th line of the file, there's a movie with the title `Dr. Seltsam\, oder wie ich lernte\, die Bombe zu lieben` German version of Dr. Strangelove). The title has commas in it, that are escaped using the `\` character! Why are CSV files so hard? 😩

**CleverCSV to the rescue!**

CleverCSV detects the dialect of CSV files much more accurately than existing approaches, and it is therefore robust against these kinds of f even has a wrapper that works with DataFrames!

```python
In [3]: from ccsv.wrappers import csv2df

df = csv2df('./data/imdb.csv')
df
```

Out[3]:

| | fn | tid | title | wordsInTitle | url | imdbRating | ratingCount | duration | year | type | ... | News |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | titles01/tt0012349 | tt0012349 | Der Vagabund und das Kind (1921) | der vagabund und das kind | http://www.imdb.com/title/tt0012349/ | 8.4 | 40550.0 | 3240.0 | 1921.0 | video.movie | ... | 0 |
| 1 | titles01/tt0015864 | tt0015864 | Goldrausch (1925) | goldrausch | http://www.imdb.com/title/tt0015864/ | 8.3 | 45319.0 | 5700.0 | 1925.0 | video.movie | ... | 0 |
| 2 | titles01/tt0017136 | tt0017136 | Metropolis (1927) | metropolis | http://www.imdb.com/title/tt0017136/ | 8.4 | 81007.0 | 9180.0 | 1927.0 | video.movie | ... | 0 |
| 3 | titles01/tt0017925 | tt0017925 | Der General (1926) | der general | http://www.imdb.com/title/tt0017925/ | 8.3 | 37521.0 | 6420.0 | 1926.0 | video.movie | ... | 0 |
| | | | Lichter der | lichter der gro | http://www.imdb.com | | | | | | | |

Markdown

```
df
```

|  |  |  | Episode 2005) | episode | /title/tt0672466/ |  |  |  |  |  |  |  |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 14757 | titles04/index.html.9992 | tt0675644 | "Playhouse 90" The Miracle Worker (TV Episode ... | playhouse the miracle worker tv episode | http://www.imdb.com /title/tt0675644/ | 7.3 | 8.0 | 5400.0 | 1957.0 | video.episode | ... | 0 |
| 14758 | titles04/index.html.9994 | tt0679222 | "Private Screenings" Robert Mitchum and Jane R... | private screenings robert mitchum and jane rus... | http://www.imdb.com /title/tt0679222/ | 7.0 | 20.0 | 3600.0 | 1996.0 | video.episode | ... | 0 |
| 14759 | titles04/index.html.9995 | tt0680064 | "Providence" All the King's Men (TV Episode 2002) | providence all the king s men tv episode | http://www.imdb.com /title/tt0680064/ | NaN | NaN | 3600.0 | 2002.0 | video.episode | ... | 0 |
| 14760 | titles04/index.html.9997 | tt0681024 | "QI" Adam (TV Episode 2003) | qi adam tv episode | http://www.imdb.com /title/tt0681024/ | 7.6 | 89.0 | 1800.0 | 2003.0 | video.episode | ... | 0 |

14761 rows × 44 columns

Hooray! 🎉

How does it work? CleverCSV searches the space of all possible dialects of a file, and computes a *data consistency measure* that quantifies how much the resulting table "looks like real data". The consistency measure combines patterns of row lengths in the parsing result and the data type of the resulting cells. This mimicks how a human would identify the dialect. If you're wondering why this problem is hard, it's because every dialect will give you *some* table, but not necessarily the correct one. More details can be found in the paper.

https://github.com/alan-turing-institute/CleverCSVDemo
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

# Gertjan van den Burg

"The fun part of data science is the modelling. Being able to read in information from a csv file should not be the hardest part.

There is no AI. I am the AI."

File    Edit    View    Insert    Cell    Kernel    Widgets    Help

Trusted    Python 3 ○

Markdown

# CSV dialect detection with CleverCSV

**Author**: Gertjan van den Burg

In this note we'll show some examples of using CleverCSV, a package for handling messy CSV files. We'll start with a motivating example and then show some other files where CleverCSV shines. CleverCSV was developed as part of a research project on automating data wrangling. It achieves an accuracy of 97% on over 9300 real-world CSV files and improves the accuracy on messy files by 21% over standard tools.

Handy links:

- Paper on arXiv
- CleverCSV on GitHub
- CleverCSV on PyPI
- Reproducible Research Repo

## IMDB Movie data

Alice is a data scientist who would like to analyse the movie ratings on IMDB for movies of different genres. She found a dataset shared by a user on Kaggle that contains information of over 14,000 movies. Great!

The data is stored in a CSV file, which is a very common data format for sharing tabular data. The first few lines of the file look like this:

– https://github.com/ alan-turing-institute/ CSV_Wrangling

– "Wrangling Messy CSV Files by Detecting Row and Type Patterns" arXiv:1811.11242

– https://github.com/ alan-turing-institute/ CSV_Wrangling

– "Wrangling Messy CSV Files by Detecting Row and Type Patterns" arXiv:1811.11242

#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

# The Turing Way

**What is Make**

Make is a build automation tool. It uses a configuration file called a Makefile that contains the *rules* for what to build. Make builds *targets* using *recipes*. Targets can optionally have *prerequisites*. Prerequisites can be files on your computer or other targets. Make determines what to build based on the dependency tree of the targets and prerequisites (technically, this is a directed acyclic graph). It uses the *modification time* of prerequisites to update targets only when needed.

**Why use Make for Reproducible Research?**

There are several reasons why Make is a good tool to use for reproducible research:

1. Make is available on many platforms
2. Make is easy to learn
3. Makefiles are text files, which makes them easy share and keep in version control.
4. Many people are already familiar with Make
5. Using Make doesn't exclude using other tools such as Travis, Docker, etc.

## Learn Make by Example

One of the things that might scare people off from using Make is that existing Makefiles can seem daunting and it may seem difficult to tailor to your own needs. In this hands-on tutorial we will

# Case studies

– Show that it can be done

– Provide templates and starting points

– Inspire

# A global collaboration

## Contributors

Thanks goes to these wonderful people (emoji key):

This project follows the all-contributors specification. Contributions of any kind welcome!

- **Sam**, who has no GitHub experience
- **Alex**, who has a lot of GitHub experience
- **Amal**, who knows they want to contribute, and does
- **Noor**, who doesn't know they want to contribute, but does
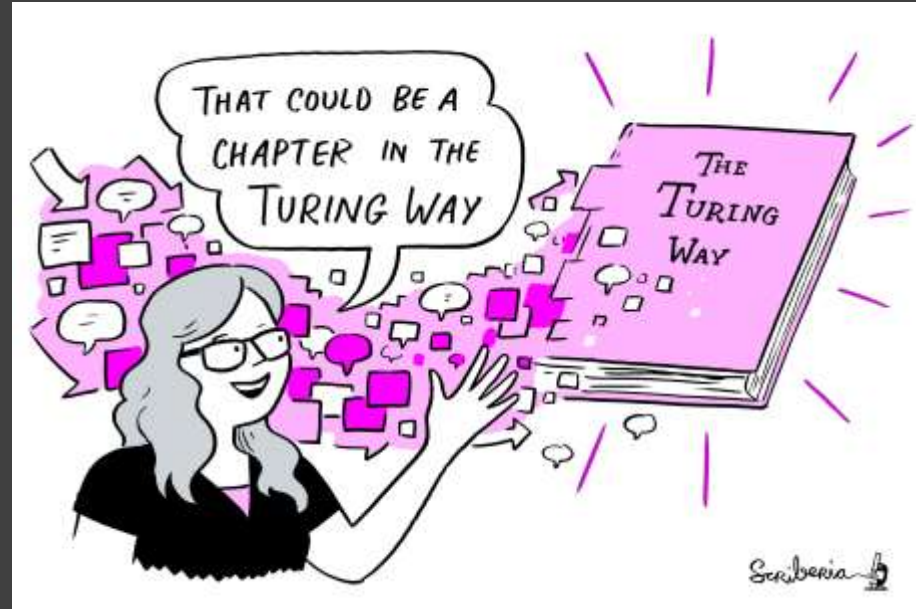
https://github.com/alan-turing-institute/the-turing-way/pull/421
#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

# The future

# **Funding extension**

– Expand scope to all data science practices

    – Ethics, model selection, project management, collaborative working

– Full time community manager, contributions from Turing & beyond



https://github.com/
alan-turing-institute/the-turing-way/
blob/master/project_management/
tps-funding-application-20190429.md

#PyDataLDN #TuringWay @kirstie_j
https://doi.org/10.5281/zenodo.3632909

# Metrics for success

- 20 new chapters

- 100 authors

- 200 contributors

- 1000 mailing list subscribers

- 50 first pull requests

- 20 new contributors to other open source projects