# The Alan Turing Institute

**binder**

**Build a BinderHub for hosting Reproducible Software in the Cloud**

Sarah Gibson

https://doi.org/10.5281/zenodo.3404774

#TuringWay #ukrse19

# Turn code… into an environment!

# This Workshop

– What it is:

– What it's not:

– What we'll do:

# This Workshop

&mdash; What it is: <span style="color:#e6007e">Challenging!</span>

&mdash; What it's not:

&mdash; What we'll do:

# This Workshop

– What it is: Challenging!

– What it's not: A cloud/Azure workshop

– What we'll do:

#TuringWay #ukrse19

# This Workshop

– What it is: Challenging!

– What it's not: A cloud/Azure workshop

– What we'll do: Build a BinderHub!

#TuringWay #ukrse19

# Housekeeping

- Microsoft Azure: Please leave your email in #binderhub-workshop channel on RSE Slack

- Docker Hub: https://hub.docker.com/signup

- Code of Conduct: Be kind! https://rse.ac.uk/conf2019/code-of-conduct/

- HackMD: bit.ly/RSEConBinderHub

- post-its 🚦

#TuringWay #ukrse19

# Who?



**Sarah**
Research Data Scientist
Operator of mybinder.org

**Tania**
Microsoft Cloud
Developer Advocate

**Anna**
Research
Software Engineer

# (Rough) Agenda

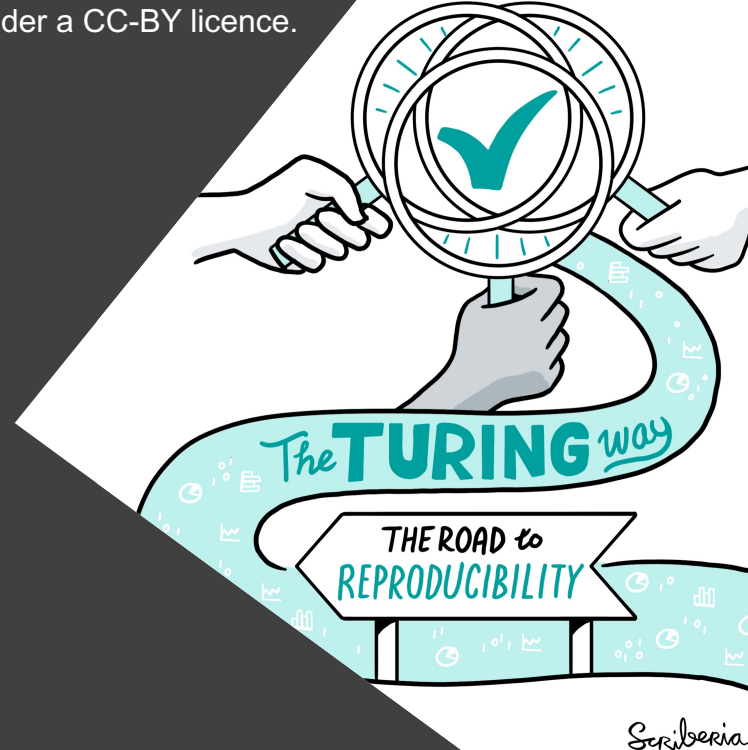| Time | Activity |
|------|----------|
| 09:00 – 09:30 | 👋 Introduction |
| 09:30 – 10:30 | 🚀 Deploy Kubernetes cluster* |
| 10:30 – 11:00 | ☕ Coffee break |
| 11:00 – 12:30 | 🧑‍💻 Install BinderHub |

*Don't worry if you don't know what this is yet, I'll explain!

# The Turing Way

A Handbook for Reproducible Data Science

*Making reproducibility too easy not to do!*

https://doi.org/10.5281/zenodo.3404774

#TuringWay #ukrse19

# Where does The Turing Way fit in?



Health

Digital twins: Cities

AI for science

Criminal justice system

Finance and economics

Digital twins: Complex systems engineering

Defence and Security

Public Policy

Tools, practices and systems for AI

#TuringWay #ukrse19

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

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

Kirstie Whitaker's talk at PyData LDN: https://youtu.be/IG3PcZ6EhiU

https://the-turing-way.netlify.com/reproducibility/03/definitions.html

#TuringWay #ukrse19

# 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, sofware 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, sofware 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, sofware 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 opinioned 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, sofware 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

# Built by a team… AND YOU!

– Please contribute!

github.com/alan-turing-institute/the-turing-way



**Contributors**

Thanks goes to these wonderful people (emoji key):

#TuringWay #ukrse19

Is not considered for promotion

Held to higher standards than others

Publication bias towards novel findings

**Barriers to reproducible research**

Requires additional skills

Plead the 5th

Support additional users

Takes time

# Market Research

Have you ever heard…?

*"Oh, it worked on my computer?"*

# Have you ever heard…?

## *"Oh, it worked yesterday?"*

#TuringWay #ukrse19

+ CI

"Oh, it worked on my computer?"

"Oh, it worked yesterday?"

Courtesy of Juliette Taka: https://twitter.com/mybinderteam/status/1082556317842264064

https://doi.org/10.5281/zenodo.3404774

#TuringWay #ukrse19

Courtesy of Juliette Taka: https://twitter.com/mybinderteam/status/1082556317842264064

#TuringWay #ukrse19

‹› Code     ⊙ Issues  2     ⌥ Pull requests  0     ⊞ Projects  0     ☰ Wiki     ⦿ Security     �III Insights

Branch: master ▾     **conda** / environment.yml          Find file    Copy path

betatim Update environment.yml          89dd429  on 11 Dec 2018

**4** contributors

14 lines (13 sloc) | 161 Bytes          Raw    Blame    History

```
 1   name: example-environment
 2   channels:
 3     - conda-forge
 4   dependencies:
 5     - numpy
 6     - psutil
 7     - toolz
 8     - matplotlib
 9     - dill
10     - pandas
11     - partd
12     - bokeh
13     - dask
```

binder-examples / **binder-r-description**

👁 Watch ▾ 1 | ★ Star 2 | ⑂ Fork 1

<> Code    ⓘ Issues 0    ⑈ Pull requests 0    ▥ Projects 0    ▤ Wiki    🛡 Security    ▥ Insights

Branch: master ▾    **binder-r-description** / **DESCRIPTION**    Find file   Copy path

👤 **gedankenstuecke** first commit      70f8b8e   on 18 Sep 2018

**1 contributor**

8 lines (7 sloc) | 282 Bytes      Raw | Blame | History   🖥 ✏ 🗑

```
1   Package: binderdescription
2   Version: 0.1
3   Date: 2018-09-18
4   Title: Binder R DESCRIPTION support
5   Description: Test that automatically building R packages works
6   Author: Bastian Greshake Tzovaras <bgreshake@googlemail.com>
7   Maintainer: Bastian Greshake Tzovaras <bgreshake@googlemail.com>
```

RESOURCES ✓

STEP 4

binder-examples / r

<> Code    ⚠ Issues 3    ⑂ Pull requests 1    ▦ Projects 0    ☰ Wiki    🛡 Security    �|�l| Insights

Branch: master ▾    r / install.R

Find file    Copy path

👤 betatim Add example Shiny app    8c01f0d    on 31 May 2018

4 contributors 👥👤👤👤

6 lines (5 sloc) | 148 Bytes

Raw    Blame    History    🖥 ✏ 🗑

```
1    install.packages("tidyverse")
2    install.packages("rmarkdown")
3    install.packages("httr")
4    install.packages("shinydashboard")
5    install.packages('leaflet')
```



BITBUCKET  GITLAB  SINGULARITY edge PARL  BINDER    NOTEBOOK ✓    EVERYONE CAN NOW RUN AND REPRODUCE HER COMPUTATIONS

{ } RESOURCES ✓    STEP ④

# What's the difference?

## mybinder.org

- Free to use service for everyone
- Must be public code
- Limited computational resources
- No GPU access

## Private BinderHub

- Service can be limited to teams or institutions
- Can be public or private code
- Set your own computational limits
- Deploy onto any type of machine you need

# The Vocab

– Binder → user interface/experience

– BinderHub → computational infrastructure

– mybinder.org → public BinderHub for everyone

~~Magic!~~ Technology

https://doi.org/10.5281/zenodo.3404774

#TuringWay #ukrse19

# BinderHub



Build and launch a repository

**GitHub repository name or URL**

https://github.com/alexmorley/binder-demo | GitHub ▾

**Git branch, tag, or commit** | **Path to a notebook file (optional)**

Git branch, tag, or commit | Path to a notebook file (optional) | File ▾ | launch

Clone GitHub Repo

1

# BinderHub

Build and launch a repository

**GitHub repository name or URL**

https://github.com/alexmorley/binder-demo | GitHub ▾

**Git branch, tag, or commit**

Git branch, tag, or commit

**Path to a notebook file (optional)**

Path to a notebook file (optional) | File ▾ | launch

**1** Clone GitHub Repo

repo2docker

**2** Build image according to instructions contained within the repo

# BinderHub

Build and launch a repository

**GitHub repository name or URL**

https://github.com/alexmorley/binder-demo                                    GitHub ▾

**Git branch, tag, or commit**                 **Path to a notebook file (optional)**

Git branch, tag, or commit                      Path to a notebook file (optional)      File ▾      launch

**1** Clone GitHub Repo

**2** Build image according to instructions contained within the repo

**3** Execute image

#TuringWay #ukrse19

# BinderHub

Build and launch a repository

**GitHub repository name or URL**

https://github.com/alexmorley/binder-demo                    GitHub ▾

**Git branch, tag, or commit**          **Path to a notebook file (optional)**

Git branch, tag, or commit  ⬆    Path to a notebook file (optional)    File ▾    launch

**1** Clone GitHub Repo

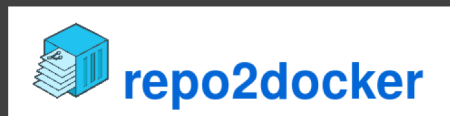**2** Build image according to instructions contained within the repo

**4** Allocate computational resources

**3** Execute image

# What is a JupyterHub?



**HTTP Proxy**

/hub/admin → **Admin**

**Config.py** ← **Hub** ♥ → **Database**

/hub/<url>

✓ ✗ **Authenticator**

**Spawners**  /user/<name>

JupyterHub is a way to help your humans use your computers. With notebooks!

All icons were obtained from Flaticon (https://www.flaticon.com/packs/essential-collection)

resources

Execute image

GitHub ▾

#TuringWay #ukrse19

# BinderHub

Build and launch a repository

**GitHub repository name or URL**

https://github.com/alexmorley/binder-demo       GitHub ▾

**Git branch, tag, or commit**          **Path to a notebook file (optional)**

Git branch, tag, or commit 🔼        Path to a notebook file (optional)   File ▾   launch

**1** Clone GitHub Repo

**2** Build image according to instructions contained within the repo

Allocate computational resources

**4**

**3** Execute image

# BinderHub

Build and launch a repository

**GitHub repository name or URL**

https://github.com/alexmorley/binder-demo          GitHub ▾

**Git branch, tag, or commit**          **Path to a notebook file (optional)**

Git branch, tag, or commit          Path to a notebook file (optional)          File ▾          launch

**1** Clone GitHub Repo

**2** Build image according to instructions contained within the repo

*jupyterhub*

Make image accessible at mybinder.org/some_url **5**

Allocate computational resources

**4**          **3**

Execute image

# BinderHub

Build and launch a repository

**GitHub repository name or URL**

https://github.com/alexmorley/binder-demo — GitHub ▾

**Git branch, tag, or commit**

Git branch, tag, or commit

**Path to a notebook file (optional)**

Path to a notebook file (optional) — File ▾ — launch

**6** Redirect User to mybinder.org/some_url

**1** Clone GitHub Repo

**2** Build image according to instructions contained within the repo

**5** Make image accessible at mybinder.org/some_url

binder

**3** Execute image

**4** Allocate computational resources

#TuringWay #ukrse19

# BinderHub

– Collection of tools
working in harmony
which BinderHub
orchestrates

# Scaling a BinderHub for multiple users

Problems if you run this on one computer:

– Resource intensive

– Resource control

– Security

#TuringWay #ukrse19

# Solution: Kubernetes!

– Resource intensive → Cluster management

– Resource control → Container management

– Security → Container isolation

# Solution: Kubernetes!

– Resource intensive → Cluster management

– Resource control → Container management

– Security → Container isolation

   Problem: Also Kubernetes… 😢


kubernetes

#TuringWay #ukrse19

# This Workshop

– What it is: Challenging!

– What it's not: A cloud/Azure workshop

– What we'll do: Build a BinderHub!

#TuringWay #ukrse19
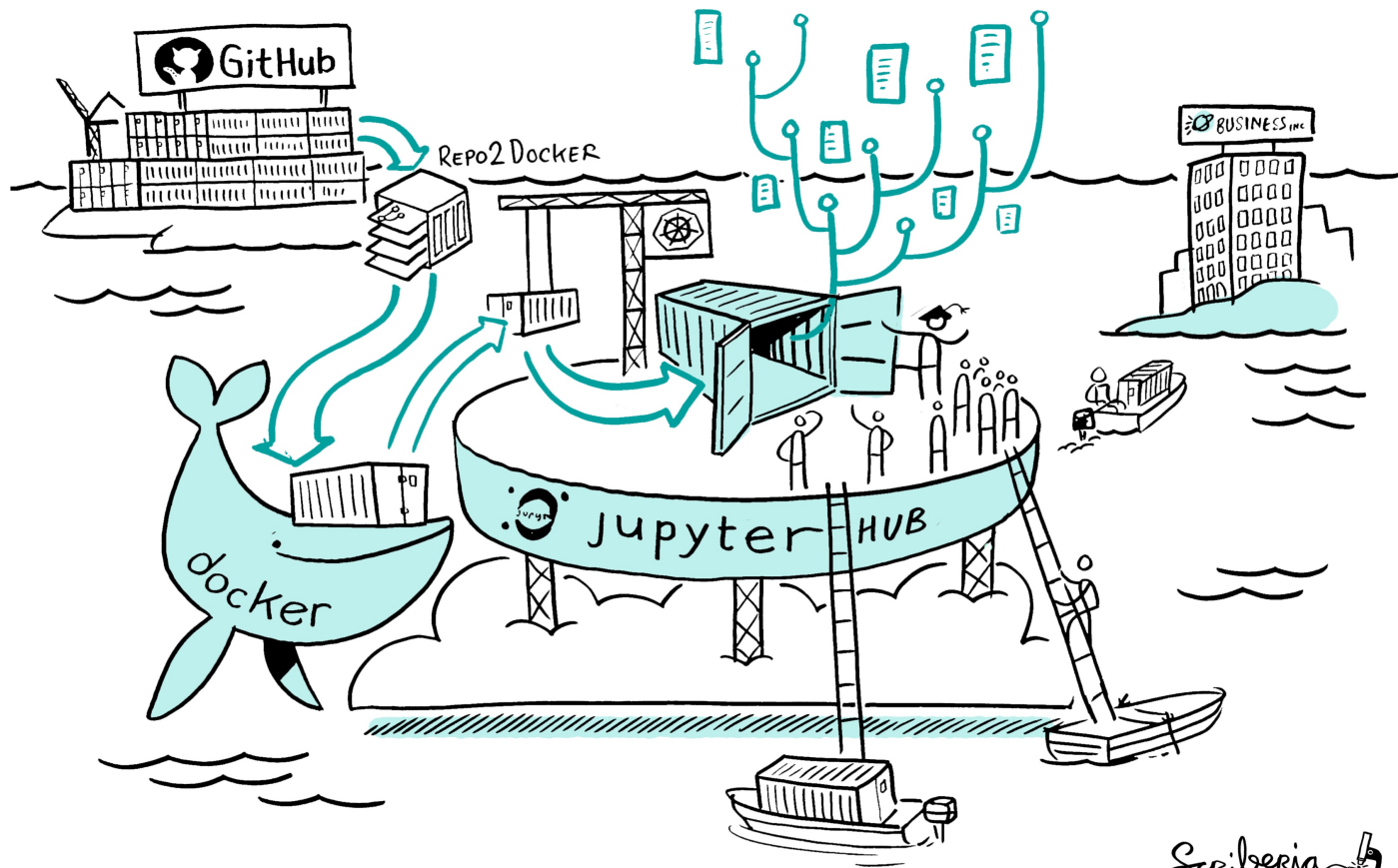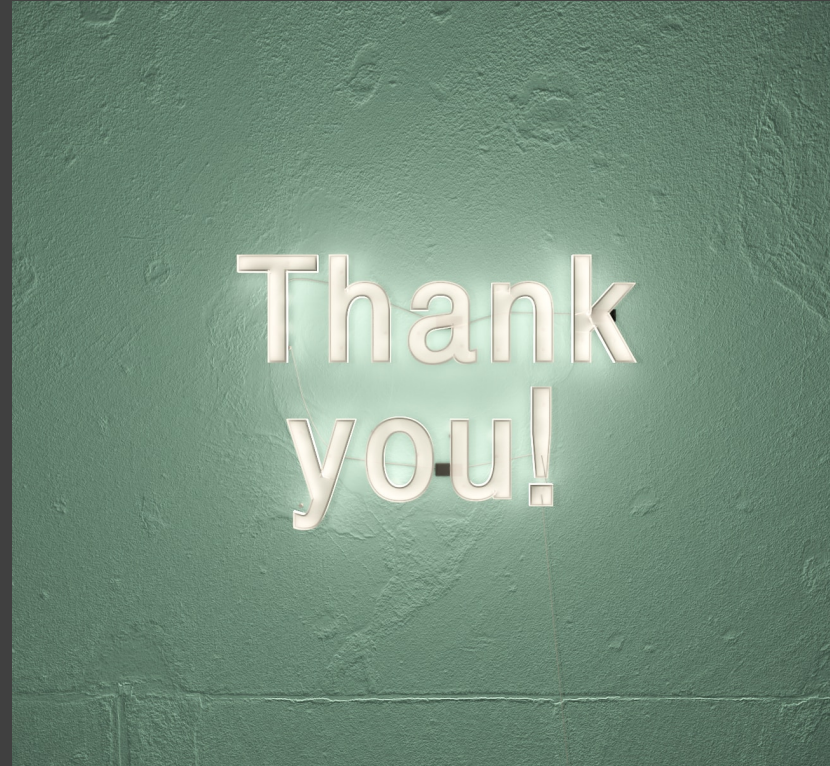
# This Workshop

## bit.ly/zero-to-binderhub-workshop

### HackMD: bit.ly/RSEConBinderHub

#TuringWay #ukrse19

– You have successfully built a BinderHub! 🤞

– Now check out this repo: github.com/alan-turing-institute/binderhub-deploy

– Please leave feedback in the HackMD: bit.ly/RSEConBinderHub

https://doi.org/10.5281/zenodo.3404774

# Thank You!
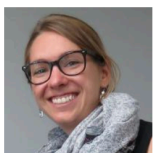
Jessica Forde — UC Berkeley — team red 📖

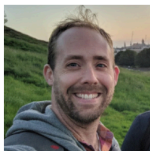Sarah Gibson — The Alan Turing Institute — team blue 💬 , 📖 , ✅
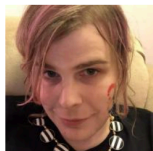
Tim Head — Wild Tree Tech — team red
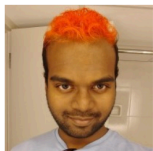
Lindsey Heagy — UC Berkeley — team blue 🤔 , 💡

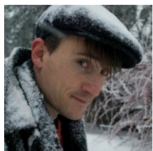Chris Holdgraf — Berkeley Institute for Data Science — team red 💻 , 🤔 , 📖 , 💬
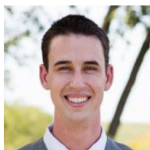
M Pacer — Netflix — team blue

Yuvi Panda — UC Berkeley — team blue 💻 , 💻

Min Ragan-Kelley — Simula — team lead — data, 💻

Zach Sailer — Project Jupyter — team blue 💻 , 🤔 , 💬

Erik Sundell — Sandvik CODE — team blue 💻 , 🚇

Carol Willing — Project Jupyter — team red

https://doi.org/10.5281/zenodo.3404774

#TuringWay #ukrse19