

We will be using an electronic whiteboard (Etherpad).

Please log onto your computer and access:

<https://public.etherpad-mozilla.org/p/2019-2-26-cruk-ci-myrd>

Data sharing and Backup

Anne Pajon

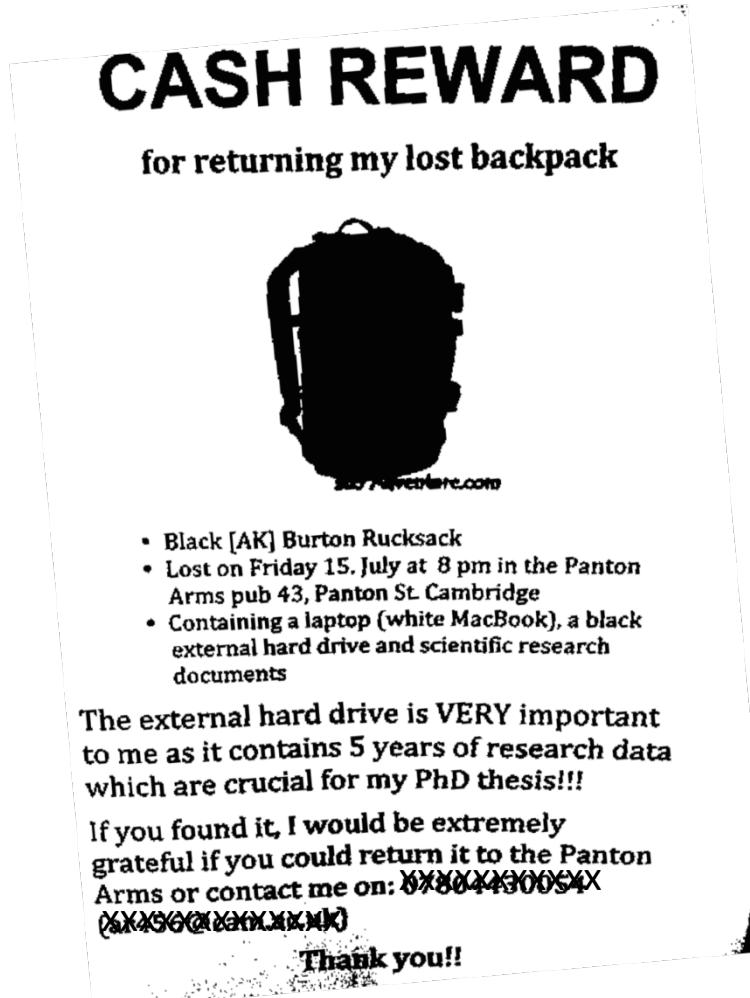
26th February 2019



This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

To avoid data *disasters* ...

What would you do if you'd lose your data tomorrow?



<https://blogs.ch.cam.ac.uk/pmr/2011/08/01/why-you-need-a-data-management-plan/>

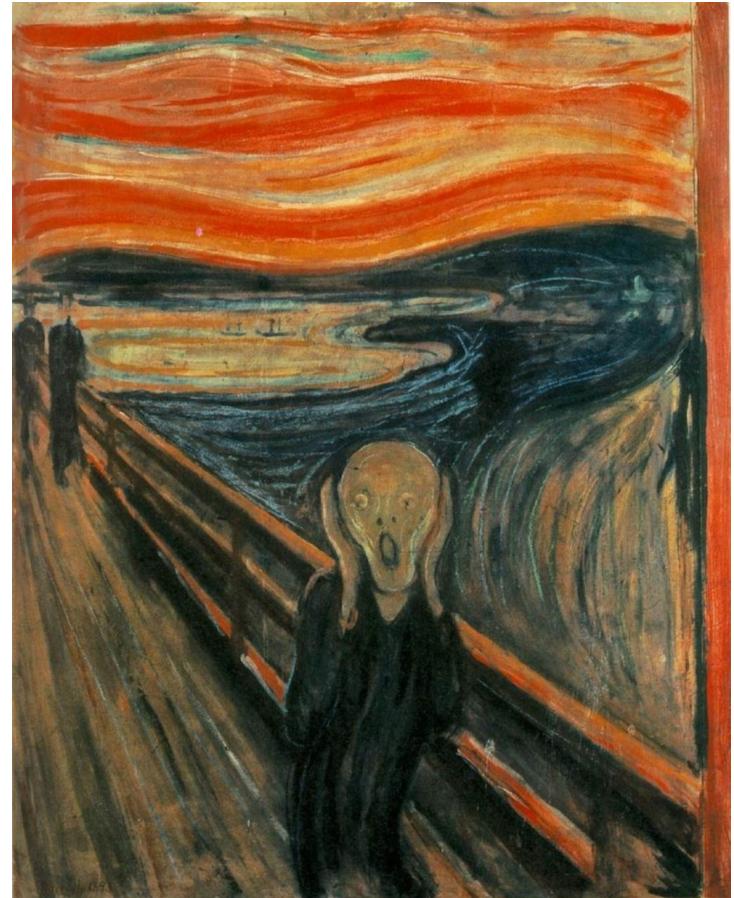
What if?



Cancer Research UK – University of Manchester – 27 April 2017

What would you do if you'd lose your data tomorrow?

- Your laptop got stolen
- Your office/house burnt
- Your USB stick is lost
- Your portable hard disk is damaged
- Your data in Dropbox disappeared



https://en.wikipedia.org/wiki/The_Scream

Never work directly on the raw data

Leave it intact

Always **make a copy**, and work on the copy

Data backup

At least 2 backups at 2 different locations

External disks



Online backup



Servers

Department
College
IT



Cheap
£10-15 / TB (1024GB)



Failure rate
1.5%/year

Accessibility
Free (limit)

Personal data
Hacking

Managed by
experts

Moving between
institutions

Data backup



Manual

Copying files to relevant folders



Automated

- Install software
e.g. Time machine
(Mac users)

- RAID technology
- Checksums



Copying files to relevant folders

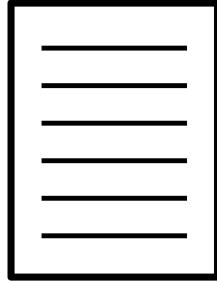
Automatically upload files to the cloud when any changes are saved

If manual ... how often?



How much would you be willing to lose?

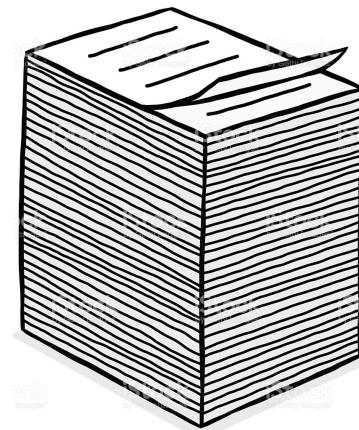
1 day



1 week



1 month-year



*Software allows you to set up **backup time** automatically*

Data backup and file sharing



Space/price	2 GB (free) 1TB (£96/year)	15 GB (free) 1 TB (£80/year)	5 GB (free) 1 TB (£56)
File history and recovery	Yes, unlimited	Yes	Last 90 days
File size limit	None	5 GB	15 GB
Support	UIS	UIS	UIS
OS	Windows, Mac, Linux, Android, iOS	Windows, Mac, Android, iOS	Windows, Mac, Android, iOS
Accessibility	Sync anywhere on any devices	Live editing	Integration with Microsoft Office

More ... file sharing



Email



Website



FTP

Why data sharing is important?



CC-BY Danny Kingsley & Sarah Brown

Data should be shared to move our knowledge forward

What data to share?

- Raw data and associated metadata
- Software and scripts
- Methods used
- Processed data
- Papers
 - results and figures
- What about non-positive results?
 - RIO - Research Ideas and Outcomes



<https://riojournal.com/>

- Would you consider Pre-Print for your draft manuscripts?



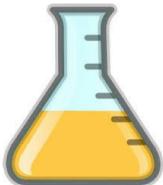
<https://www.biorxiv.org/>

DOAJ DIRECTORY OF
OPEN ACCESS
JOURNALS

<https://doaj.org/>

When to share?

Close data



Initial experiments

Method optimisation

Answering biological questions

Write up story

Ideas



Paper



Tests

First scripts

Robust analysis pipeline

Figures
Tables

Document computational method

Private code

Public code

Gather information early

- Think of data submission from the start
- Collect metadata before submission
- Keep a **ReadMe** file about your project
- Consider using an electronic notebook

Making your publication Open Access

The Open Access Team will check your funder and journal policies and advise on how to comply with Open Access requirements.



Accepted for publication?
Upload manuscript

<https://www.openaccess.cam.ac.uk/>

How to share your data?

Store, describe and deposit your data in suitable and trusted public data repositories and add a link to your data in your publication.

Repositories for datasets

- Discipline specific
 - Registry of Research Data Repositories
<http://www.re3data.org/>
 - EMBL-EBI services
<https://www.ebi.ac.uk/services>
- General purpose
 - Zenodo <https://zenodo.org/>

Repositories for software

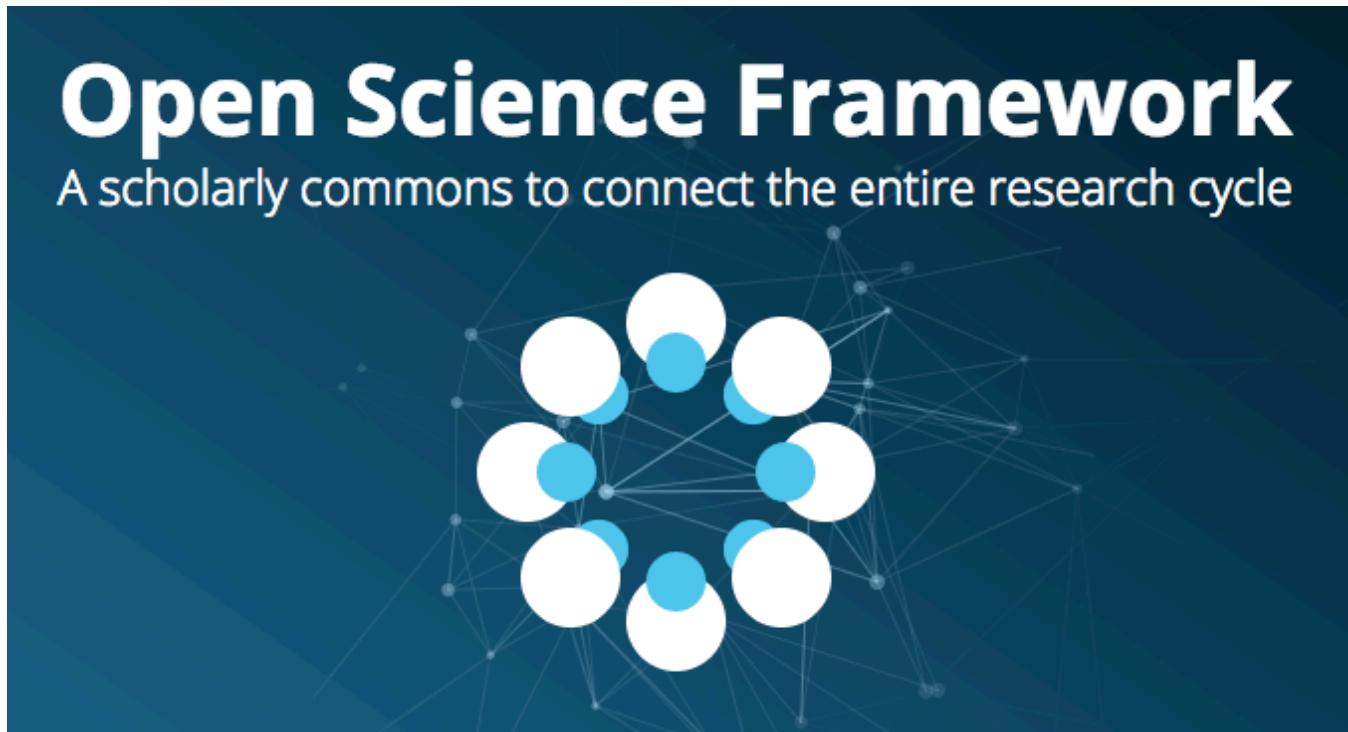
- GitHub <https://github.com>
- GitLab <https://gitlab.com>
- Bitbucket
<https://bitbucket.org>



Zenodo assigns a Digital Object Identifier (DOI) to make the upload easily and uniquely citeable, with GitHub integration to enable tracking of each release.

Open Science Framework

Cloud-based management for your projects @ <https://osf.io/>



Under which license?

Share your work with one of the Creative Commons licenses
<https://creativecommons.org/>

The screenshot shows a teal-colored web page titled "Choose a license". It features three main steps: "Choose Features" (with a hand icon), "Optional Info" (with a document icon), and "Get License" (with a CC logo). A large yellow "Get Started" button is at the bottom. The text below the first step says: "This chooser helps you determine which Creative Commons License is right for you in a few easy steps. If you are new to Creative Commons, you may also want to read [Licensing Considerations](#) before you get started."

Software licenses

<https://choosealicense.com>

Choose an open source license

{ Which of the following best describes your situation? }



I want it simple and permissive.

The [MIT License](#) is a permissive license that is short and to the point. It lets people do anything they want with your code as long as they provide attribution back to you and don't hold you liable.

[jQuery](#), [.NET Core](#), and [Rails](#) use the MIT License.



I'm concerned about patents.

The [Apache License 2.0](#) is a permissive license similar to the MIT License, but also provides an express grant of patent rights from contributors to users.

[Android](#), [Apache](#), and [Swift](#) use the Apache License 2.0.



I care about sharing improvements.

The [GNU GPLv3](#) is a copyleft license that requires anyone who distributes your code or a derivative work to make the source available under the same terms, and also provides an express grant of patent rights from contributors to users.

[Bash](#), [GIMP](#), and [Privacy Badger](#) use the GNU GPLv3.

{ What if none of these work for me? }

My project isn't software.

There are licenses for that.

I want more choices.

More licenses are available.

I don't want to choose a license.

You don't have to.

Conclusion



- Always **make a copy** of your raw data
- **Backup** your data at least **twice** at two different locations
- Document your process using a **ReadMe** file
- Ideally most data should be shared
 - Sharing is essential for all publicly funded research
 - Share as early as possible
 - Using suitable repositories and DOI
- Under Creative Commons or Open Source license