

File and Data Management

Jing Su

20180223

What You Will Learn

- Why file management of your research data is important
- Specific techniques for organizing your research data,
 - File structures
 - File naming
 - Version control
 - Storage & Backup
- Including:
 - Small group discussion
 - Exercise for organizing your own data
- Focuses on research data, also applies to other types of files

Small Group Discussion

- What kind of data do you work with?
- What organizational challenges have you faced?
- What tools or techniques work for you?

Research Data Lifecycle



Data Management Checklist 1/2

- Are you using standardised and consistent procedures to collect, process, check, validate and verify data?
- Are your structured data self-explanatory in terms of variable names, codes and abbreviations used?
- Which descriptions and contextual documentation can explain what your data mean, how they were collected and the methods used to create them?
- How will you label and organise data, records and files?
- Will you apply consistency in how data are catalogued, transcribed and organised, e.g. standard templates or input forms?
- Which data formats will you use? Do formats and software enable sharing and long-term validity of data, such as non-proprietary software and software based on open standards?
- When converting data across formats, do you check that no data or internal metadata have been lost or changed?
- Are your digital and non-digital data, and any copies, held in a safe and secure location?
- Do you need to securely store personal or sensitive data?

Data Management Checklist 2/2

- If data are collected with mobile devices, how will you transfer and store the data?
- If data are held in various places, how will you keep track of versions?
- Are your files backed up sufficiently and regularly and are back-ups stored safely?
- Do you know what the master version of your data files is?
- Do your data contain confidential or sensitive information? If so, have you discussed data sharing with the respondents from whom you collected the data?
- Are you gaining (written) consent from respondents to share data beyond your research?
- Do you need to anonymise data, e.g. to remove identifying information or personal data, during research or in preparation for sharing?
- Have you established who owns the copyright of your data? Might there be joint copyright?
- Who has access to which data during and after research? Are various access regulations needed?
- Who is responsible for which part of data management?
- Do you need extra resources to manage data, such as people, time or hardware?

Data Management Checklist

- **What** types of data and for how long?

Five steps to decide what data to keep

- Step 1. Identify purposes that the data could fulfill
- Step 2. Identify data that must be kept
- Step 3. Identify data that should be kept
- Step 4. Weigh up the costs
- Step 5. Complete the data appraisal

- **Who** will be responsible to collect and document the data?

Roles and responsibilities. Legal and ethical obligations and right. Plan and consent to share.

- **How** to document different types of data?

Study-level, Data-level, and Metadata

Wet lab: [Electronic Lab Notebook \(ELN\)](#)

Computational: large size sequencing data, consortium data (TCGA, ICGC)

Some current ELN products

	Suitability	Platform	Storage	Comments
Benchling	Individual, Group	Browser	Vendor cloud	Free, user-friendly, self-contained, Molecular Biology bias.
Biovia	Group, Department	Macintosh, Windows	Vendor cloud or local server	Basic but robust feature set and workflow, strong in compliance, deployed campus-wide at some institutions.
Docollab	Individual, Group	Browser	Vendor cloud	Basic feature set with simple, modern interface.
e-Notebook	Group, Department	Windows	Local server	Complex, Chemistry/Pharma bias.
e-Workbook	Group, Department	Browser	Vendor cloud	Strong inventory management, Chemistry/Pharma bias.
eLabFTW	Group, Department	Browser	Local server	Free, Open Source, requires local server (Docker containers recommended). Community-driven development, sponsored by Institut Curie.
eLABJournal	Individual, Group, Department	Browser	Vendor cloud, private cloud, or local server	Comprehensive product with strong inventory management integration
Findings	Individual	Macintosh, iOS	Local HD, Dropbox sync	Simple, attractive interface, good synchronisation with Apple devices.
Hivebench	Individual, Group, Department	Browser, Macintosh, iOS	Vendor cloud or local server	New product - pilot deployment in one Cambridge research group has been a positive experience.
LabArchives	Individual, Group	Browser, iOS	Vendor cloud	Comprehensive features including Graphpad Prism integration. <small>[Trial comments (Cambridge only)]</small>

Formats: Data type and sources

File formats currently recommended by UK Data Archive for long term preservation for research data

FILE FORMATS CURRENTLY RECOMMENDED BY THE UK DATA ARCHIVE FOR LONG-TERM PRESERVATION
OF RESEARCH DATA

TYPE OF DATA	RECOMMENDED FILE FORMATS FOR SHARING, RE-USE AND PRESERVATION
Quantitative tabular data with extensive metadata a dataset with variable labels, code labels, and defined missing values, in addition to the matrix of data	SPSS portable format (.por) delimited text and command ('setup') file (SPSS, Stata, SAS, etc.) containing metadata information some structured text or mark-up file containing metadata information, e.g. DDI XML file
Quantitative tabular data with minimal metadata a matrix of data with or without column headings or variable names, but no other metadata or labelling	comma-separated values (CSV) file (.csv) tab-delimited file (.tab) including delimited text of given character set with SQL data definition statements where appropriate
Geospatial data vector and raster data	ESRI Shapefile (essential: .shp, .shx, .dbf ; optional: .prj, .sbx, .sbn) geo-referenced TIFF (.tif, .tfw) CAD data (.dwg) tabular GIS attribute data
Qualitative data textual	eXtensible Mark-up Language (XML) text according to an appropriate Document Type Definition (DTD) or schema (.xml) Rich Text Format (.rtf) plain text data, ASCII (.txt)
Digital image data	TIFF version 6 uncompressed (.tif)
Digital audio data	Free Lossless Audio Codec (FLAC) (.flac)
Digital video data	MPEG-4 (.mp4) motion JPEG 2000 (.jp2)
Documentation	Rich Text Format (.rtf) PDF/A or PDF (.pdf) OpenDocument Text (.odt)

File Naming Conventions

- Make file names unique
- Include most important identifying information of the project:
 - ✓ project name
 - ✓ acronym, or research data name
 - ✓ study title
 - ✓ location information
 - ✓ researcher initials
 - ✓ date (consistently formatted, e.g. YYYYMMDD)
 - ✓ version
- Use underscores to separate elements; avoid special characters, spaces and periods.
- Use leading zeros when incorporating numbers to enable sorting (a sequence of 1-100 should be numbered 001-100).
- File names should be short enough to be readable, while still conveying enough pertinent information (limits 255 chars)

File Naming Conventions Examples

- **The Good:** DryValleySoil_ICPOES_20101115_JDSv2.dat
 - DryValleySoil, project name
 - ICPOES, instrument name
 - 20101115 date of sample created
 - JDS, initials of the scientist
 - V2, second version
- **The Bad:** my Data @DryValley November 15 2010.v2.dat
- **The Ugly:**

Can you understand/use these data files? Would anyone 5 years from now?

 - SrvMthdDraft.doc
 - SrvMthdFinal.doc
 - SrvMthdLastOne.doc
 - SrvMthdRealVersion.doc

Use content-or descriptive information

Batching Renaming Tools

- Windows:
 - Adobe Bridge (via any Creative Cloud products): <http://ist.mit.edu/adobe-creative-cloud>
 - Ant Renamer: <http://www.antp.be/software/renamer>
 - Bulk Rename Utility: <http://www.bulkrenameutility.co.uk/>
 - ImageMagick: <http://www.imagemagick.org/>
 - PSRenamer: <http://www.powersurgepub.com/products/psrenamer.html>
 - RenameIT: <http://sourceforge.net/projects/renameit>
- Mac:
 - Adobe Bridge (via any Creative Cloud products): <http://ist.mit.edu/adobe-creative-cloud>
 - ImageMagick: <http://www.imagemagick.org/>
 - Name Changer: http://web.mac.com/mickeyroberson/MRR_Software/NameChanger.html
 - PSRenamer: <http://www.powersurgepub.com/products/psrenamer.html>
 - Renamer4Mac : <http://renamer4mac.com/>
 - Name Mangler: <http://manytricks.com/namemangler/>
- Linux:
 - GNOME Commander: <http://www.nongnu.org/gcmd/>
 - GPRename: <http://gprename.sourceforge.net/>
 - ImageMagick: <http://www.imagemagick.org/>
 - PSRenamer: <http://www.powersurgepub.com/products/psrenamer.html>
- Unix
 - The use of the grep command to search for regular expressions

Version Control

Aim: Keep raw data untouched and reverse to earlier version

- Save an untouched copy of the raw data, work on save untouched copy
- Use a file naming convention (like v001, v002 or v1_0, v1_2, v2_0)
- Use a directory structure naming convention that includes version information
- Date can be part of the file name, e.g.
2012-02-27_Template_soil_testing.xlsx
- Append the author's name to the file name, e.g.
Template_soil_testing_modified_by_AH.xlsx
- Add a version number after each major edit, e.g.
Template_soil_testing_v03.xlsx
- Directory top-level folders should include the project title, unique identifier, and date (year), but the files themselves should be well-described independent of the directory structure.
- Version control tools:
 - Wet lab: Electronic Lab Notebooks/Box/LIMS
 - Dry lab: SVN/GitHub

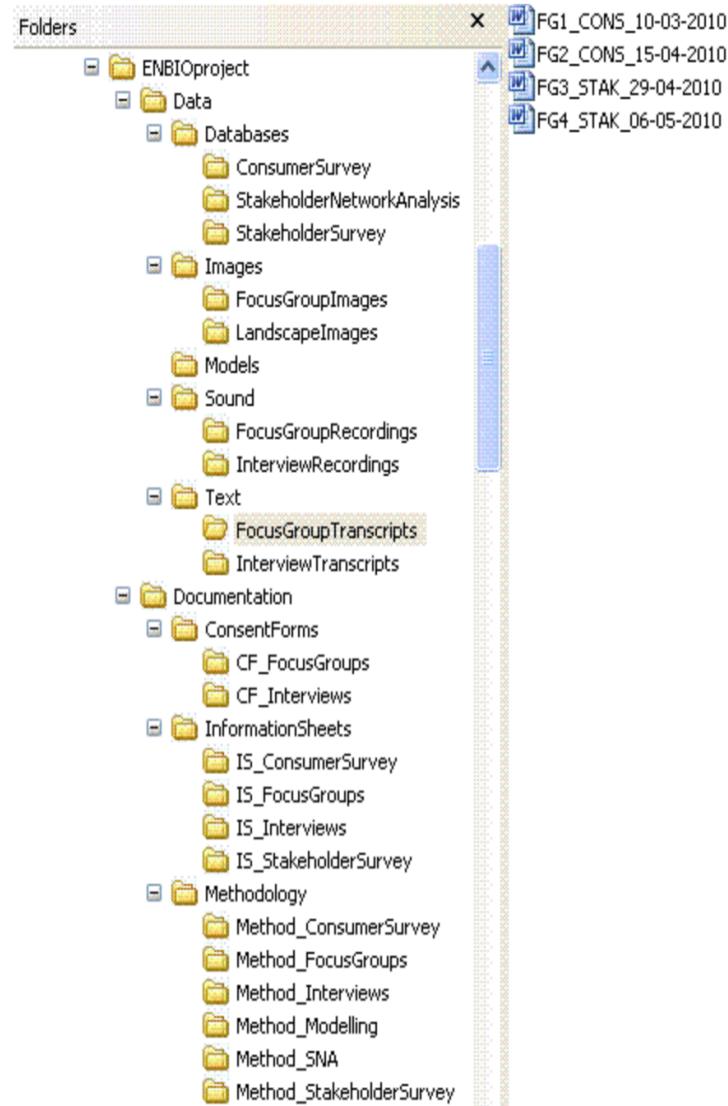
Version Control Example

VERSION CONTROL TABLE FOR A DATA FILE			
Title:	Vision screening tests in Essex nurseries		
File Name:	VisionScreenResults_00_05		
Description:	Results data of 120 Vision Screen Tests carried out in 5 nurseries in Essex during June 2007		
Created By:	Chris Wilkinson		
Maintained By:	Sally Watsley		
Created:	04/07/ 2007		
Last Modified:	25/11/ 2007		
Based on:	VisionScreenDatabaseDesign_02_00		
VERSION	RESPONSIBLE	NOTES	LAST AMENDED
00_05	Sally Watsley	Version 00_03 and 00_04 compared and merged by SW	25/11/2007
00_04	Vani Yussu	Entries checked by VY, independent from SK	17/10/2007
00_03	Steve Knight	Entries checked by SK	29/07/2007
00_02	Karin Mills	Test results 81-120 entered	05/07/2007
00_01	Karin Mills	Test results 1-80 entered	04/07/2007

Folder Structure

- Methods of organising electronic material
 - **Hierarchical**: Items organised in folders and sub-folders
 - **Tag-based**: Each item assigned one or more tags
 - **Hybrid** combination of hierarchical and tag-based

Folder Structure Examples – Hierarchical



Folder Structure Examples – Tag-based

Demo i

+ New 42 files found Search Q

File Ext.	Title ↑	Tags	Size	Date Modified
TXT		20160930-161208	125 B	2016.10.15 - 11:54:30
HTML		20161110-201627 medium	2.1 kB	2016.11.10 - 20:28:52
JPG	034-IMG_29263	Sstar waiting 48.6764537+21.9122314 20161214	558.1 kB	2016.07.21 - 21:50:15
JPG	458f25ac-1692- 9b52-591fd280d 8f8	20170126 car restaurant	419.5 kB	2015.12.29 - 11:19:18
WAV	asd	test	557.1 kB	2016.09.22 - 13:33:06
PDF	bitmessage	paper	198.9 kB	2014.04.02 - 10:09:42
PDF	Cafe Wedekind	restaurant location linux low tan1	140.7 kB	2014.04.02 - 10:08:08

HTML i Meeting Notes x

audi 20130629 eo-34

Search Q

File ≡ ↶ ↷ ↶ ↷ ?

UBUNTU

Participants:

- Max Mustermann
- John Smith
- Beatrice Karl
- John Gallius

Agenda:

- Budget discussion
- Vacation plan
- Misceleneous

1. Budget discussion

2. Vacation plan

3. Misceleneous

4. Miscellaneous

5. Miscellaneous

6. Miscellaneous

7. Miscellaneous

8. Miscellaneous

9. Miscellaneous

10. Miscellaneous

11. Miscellaneous

12. Miscellaneous

13. Miscellaneous

14. Miscellaneous

15. Miscellaneous

16. Miscellaneous

17. Miscellaneous

18. Miscellaneous

19. Miscellaneous

20. Miscellaneous

21. Miscellaneous

22. Miscellaneous

23. Miscellaneous

24. Miscellaneous

25. Miscellaneous

26. Miscellaneous

27. Miscellaneous

28. Miscellaneous

29. Miscellaneous

30. Miscellaneous

31. Miscellaneous

32. Miscellaneous

33. Miscellaneous

34. Miscellaneous

35. Miscellaneous

36. Miscellaneous

37. Miscellaneous

38. Miscellaneous

39. Miscellaneous

40. Miscellaneous

41. Miscellaneous

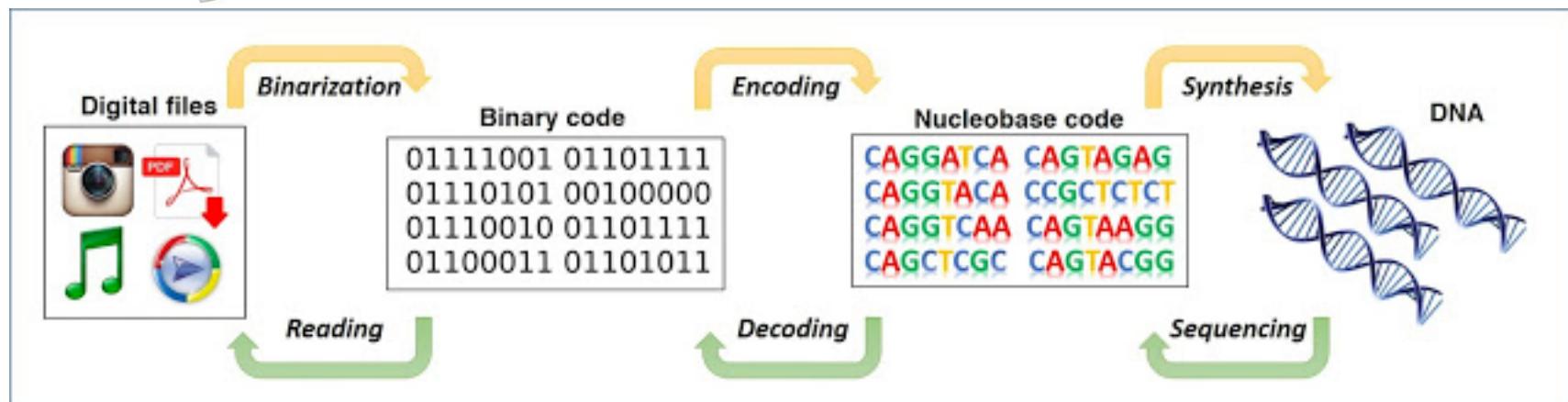
42. Miscellaneous

Small Group Discussion

- What sort of structure(s) do you currently use?
- What do you see as the key advantages and disadvantages of the different types of system?
- Are there specific tasks one sort of system seems particularly suitable for? How does this apply to your research project?

Data Storage

The everlasting external disks



Are they really permanent? What if...

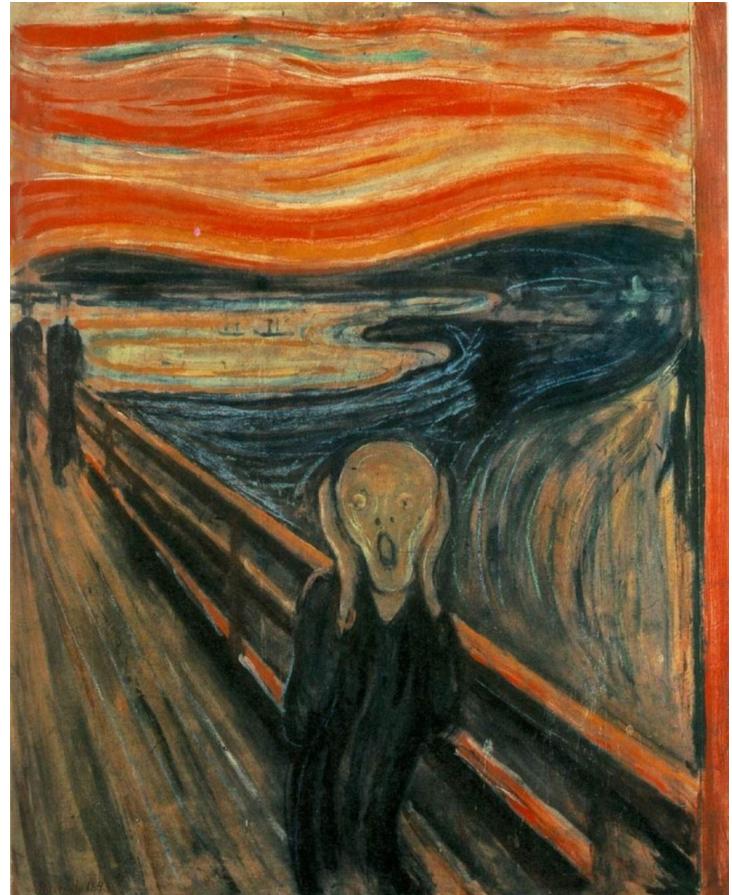
What if your data is lost



Cancer Research UK – University of Manchester – 27 April 2017

What if your data is lost

- Your laptop got stolen
- Your office/house burnt
- Your USB stick is lost
- Your portable hard disk is damaged
- Data copied to Dropbox disappeared



Storage + Security + Encryption + Backup + Sharing





Backup ↕

Ctrl

Storage + Security + Encryption + Backup + Sharing

- University Storage Service

- SharePoint Online
- Dropbox
- OneDrive
- Research Data Store
- Research Cold Store
- Research File Share



Save 50% on Dropbox Business subscriptions when you buy through the University

- CRUK CI IT
- Lab
- Individual (Time machine)
- CLOUD?



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

Data backup and file sharing

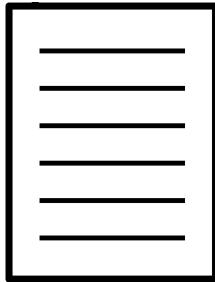


Space/price	2 GB (free) Unlimited (£55/year)	15 GB (free) 1 TB (~£80/year)	1 TB (free)
File history and recovery	Yes, unlimited	Yes	Last 90 days
File size limit	None	5 GB	15 GB
Support	UIS	Unsupported	UIS
OS	Windows, Mac, Linux, Android, iOS	Windows, Mac, Android, iOS	Windows, Mac, Android, iOS
Accessibility	Sync anywhere on any devices	Live editing	Integration with Microsoft Office

Data Backup

- Q: If manual ... how often?
A: How much would you be willing to lose?
- *Software allows you to set up **backup time automatically***

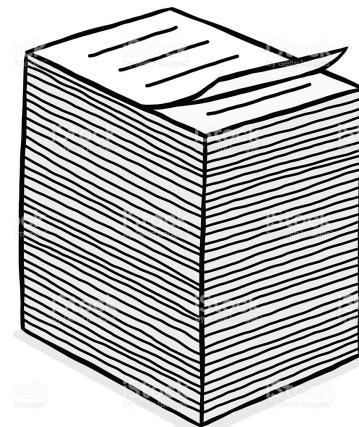
1 day



1 week



1 month-year



More ... file sharing



Email



Website



FTP