

# Bash scripting

Lecture 8  
BIOS 6660, Spring 2019  
Instructor: Pam Russell



# Bash

The shell language we use on the command line

A full-featured programming language

# Bash scripts

Save Bash commands in a text file

Script runs in shell environment

Common for genomic workflows using command line tools

# Bash script in HW4: check\_files.sh

```
#!/bin/bash

# This script checks for the existence of required files for HW4,
# and prints some basic information about each.
# You need to complete the section marked "TODO".

# "-e" means "process escape characters" e.g. render "\n" as a
# newline instead of literal "\n"
echo -e "\nChecking for required files...\n"

# Declare some string variables
# Note the "$" to refer to a variable that has already been declared
# Note how you combine strings into a concatenated string
dir_repo="..../.."
dir_proj=..
dir_data="$dir_proj/data"
dir_output="$dir_proj/output"
dir_src="$dir_proj/src"

# Create an array of names of the required files for the homework
# Note "\" means continue the command onto the next line, so this
# array declaration is a single command written across multiple lines
required_files=\
($dir_data/README.txt \
$dir_repo/.gitignore \
$dir_src/src.Rproj \
$dir_src/workflow.html \
$dir_output/avg_wkday_traffic.png)
"check_files.sh" 86L, 2291C
```

# Running a bash script

```
MacBook-Pro-9:src Pamela$ chmod +x check_files.sh
MacBook-Pro-9:src Pamela$ ./check_files.sh
Checking for required files...
Missing required file: ../data/README.txt
Missing required file: ../src/src.Rproj
Missing required file: ../src/workflow.html
Missing required file: ../output/avg_wkday_traffic.png
Some required files are missing. Exiting.
MacBook-Pro-9:src Pamela$
```

Make file executable by adding **+x** permissions (only need to do once)

Run script.  
Note you need to supply the directory (".")

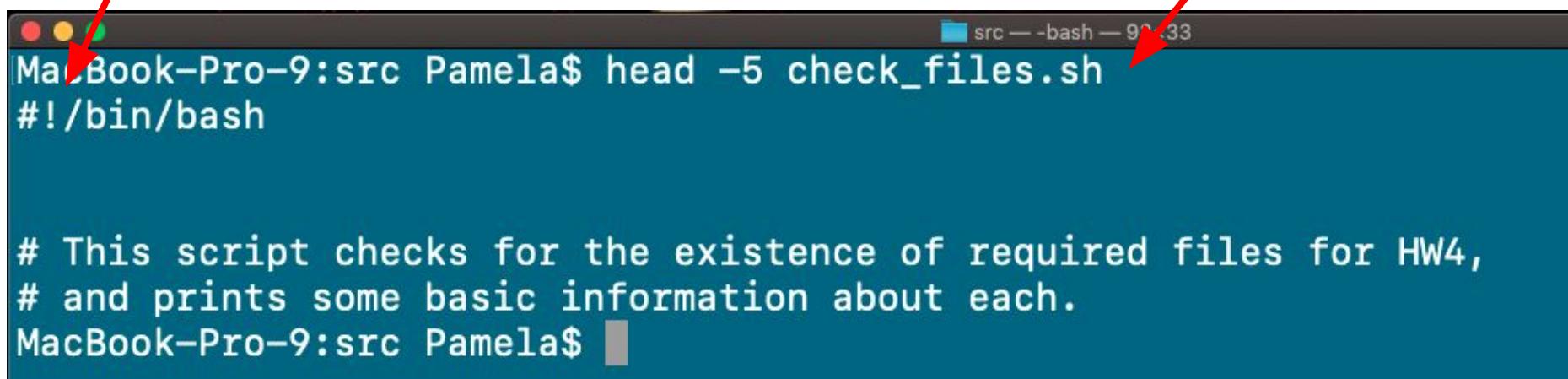
# File extension and shebang

Shebang: **#!** ("hash bang")

Instructs shell which interpreter to use.

Use **#!/bin/bash** in all Bash scripts.

File extension actually doesn't matter except to remind us. Shebang determines how file will be interpreted.

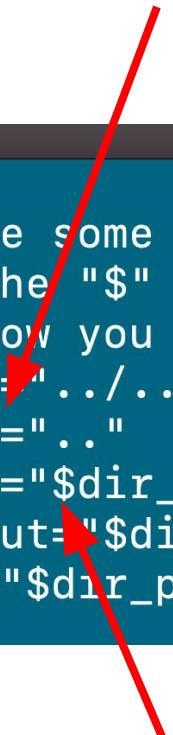


```
MacBook-Pro-9:src Pamela$ head -5 check_files.sh
#!/bin/bash

# This script checks for the existence of required files for HW4,
# and prints some basic information about each.
MacBook-Pro-9:src Pamela$
```

# Variable declarations and references

Declare variable **dir\_proj**



```
# Declare some string variables
# Note the "$" to refer to a variable that has already been declared
# Note how you combine strings into a concatenated string
dir_repo=".../..."
dir_proj=...
dir_data="$dir_proj/data"
dir_output="$dir_proj/output"
dir_src="$dir_proj/src"
```

Reference **dir\_proj** with \$

# Control structures

For loop



```
# Check if all the files exist
# Use a Bash for loop to loop through the required files
all_ok=true
for file in "${required_files[@]}"
do
    # "-f $file" means file exists and is a regular file
    # "!" -f $file" means the negation
    if [ ! -f $file ]; then
        all_ok=false
        echo "Missing required file: $file"
    fi
done
```

If/else  
statement



```
# Print a message about whether all files were found
# If any files are missing, terminate the program here
if [ $all_ok = true ]; then
    echo -e "\nFound all required files."
else
    echo -e "\nSome required files are missing. Exiting.\n"
    exit -1
fi
```

# Line continuation

```
# Create an array of names of the required files for the homework
# Note "\" means continue the command onto the next line, so this
# array declaration is a single command written across multiple lines
required_files=\
($dir_data/README.txt \
$dir_repo/.gitignore \
$dir_src/src.Rproj \
$dir_src/workflow.html \
$dir_output/avg_wkday_traffic.png)
```



Backslash lets you continue a  
single command onto multiple  
lines

# Capturing output of a command as a variable

Capture output of the **whoami** command in a variable called **me**

```
MacBook-Pro-9:src Pamela$ me=$(whoami)
MacBook-Pro-9:src Pamela$ echo $me
Pamela
MacBook-Pro-9:src Pamela$ files=$(ls)
MacBook-Pro-9:src Pamela$ echo $files
check_files.sh functions.R tests.R workflow.Rmd
MacBook-Pro-9:src Pamela$
```

\$ for references to existing variables

# Pipes and awk

Don't want to include  
the file name

Try to capture number of lines in  
**tests.R** as a variable

```
MacBook-Pro-9:src Pamela$ nlines_tests=$(wc -l tests.R)
MacBook-Pro-9:src Pamela$ echo $nlines_tests
21 tests.R
MacBook-Pro-9:src Pamela$ nlines_tests=$(wc -l tests.R | awk '{print $1}')
MacBook-Pro-9:src Pamela$ echo $nlines_tests
21
MacBook-Pro-9:src Pamela$
```

Pipe the original command to an **awk** command that prints the first field of each line of input

# Bash scripts for genomics

Get arguments from command line



```
#!/bin/bash

module load samtools
module load bedtools2/2.26.0

if [ "$#" -ne 2 ]; then
    echo $#
    echo -e "usage: compute_coverage.sh [bam directory] [\"[region1 (e.g. 11:230000-250000)] [region2] ...\"]"
    exit -1
fi

bsub_dir="$outdir/bsub"
mkdir $bsub_dir
```

Do some text processing and write several output files



```
IFS=' ' read -r -a regions <<< "$2"
for region in "${regions[@]}"
do
    region_printable=$(echo $region | perl -p -e 's/[:-]/_/g')
    region_bed="$outdir/${region_printable}.bed"
    region_tabs=$(echo $region | perl -p -e 's/[:-]/\t/g')
    echo -e "${region_tabs}\t${region_printable}\t$0\t+" > ${region_bed}
done

for bam in $1/*bam; do
```

Loop through all .bam files in a directory and run a genomic program for each



```
    sample=$(basename $bam ".bam")

    for region in "${regions[@]}"
    do
        region_printable=$(echo $region | perl -p -e 's/[:-]/_/g')
        outfile="$outdir/${sample}.coverage_${region_printable}.txt.gz"
        if [ ! -f $outfile ]; then
            tmp_bed="$outdir/${region_printable}.bed"
            bsub_output="$bsub_dir/bsub_${sample}_${region_printable}.out"
            rm -f $bsub_output
            eight=$8'
            bsub -o $bsub_output -R "rusage[mem=29000]" -R "select[mem>29000]" \
                "bedtools coverage -a ${tmp_bed} -b $bam -d | awk '{print $eight}' | gzip -f > $outdir/${sample}.coverage_${reg
```

```
done
compute_coverage.sh lines 1-43/45 100%
```

# Looping through files in a directory

- `$1` is a special variable holding the first command line argument, a directory in this case
- `*bam` means all files ending with "bam"

```
for bam in $1/*bam; do
    sample=$(basename $bam ".bam")
    for region in "${regions[@]}"
    do
        region_printable=$(echo $region | perl -p -e 's/[:-]/_/g')
        outfile="$outdir/${sample}.coverage_${region_printable}.txt"
        if [ ! -f $outfile ]; then
            tmp_bed="$outdir/${region_printable}.bed"
            bsub_output="$bsub_dir/bsub_${sample}_${region_printable}"
            rm -f $bsub_output
            eight='\$8'
            bsub -o $bsub_output -R "rusage[mem=29000]" -R "select[8]"
                "bedtools coverage -a ${tmp_bed} -b $bam -d | awk"
        else
            echo "File exists $outfile"
        fi
    done
done
```

# Why use Bash instead of an R script?

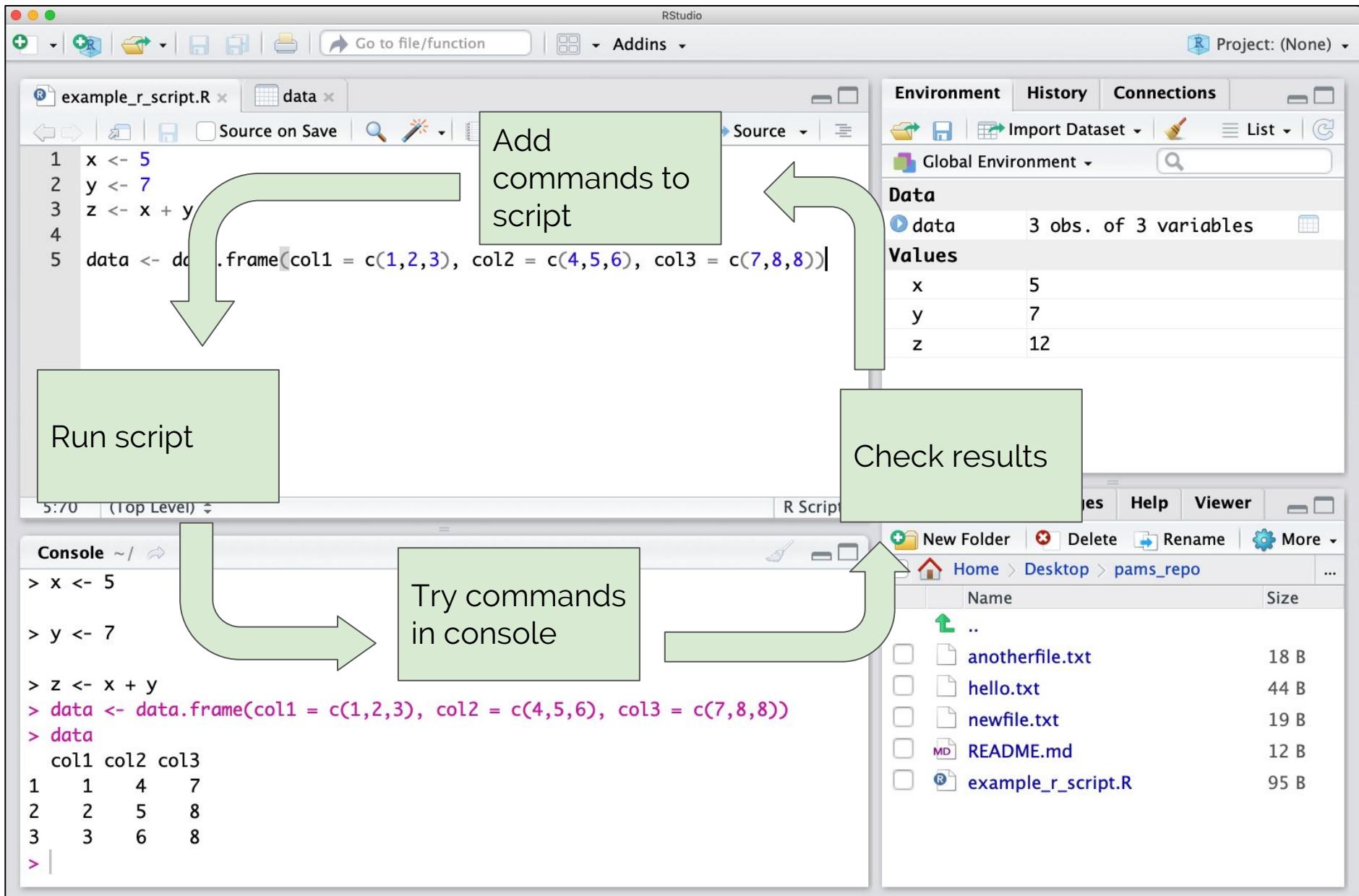
R:

- Has analogs of many Bash commands e.g.  
**list.files()**
- Theoretically can run any Linux command from R (see  
[here](#))

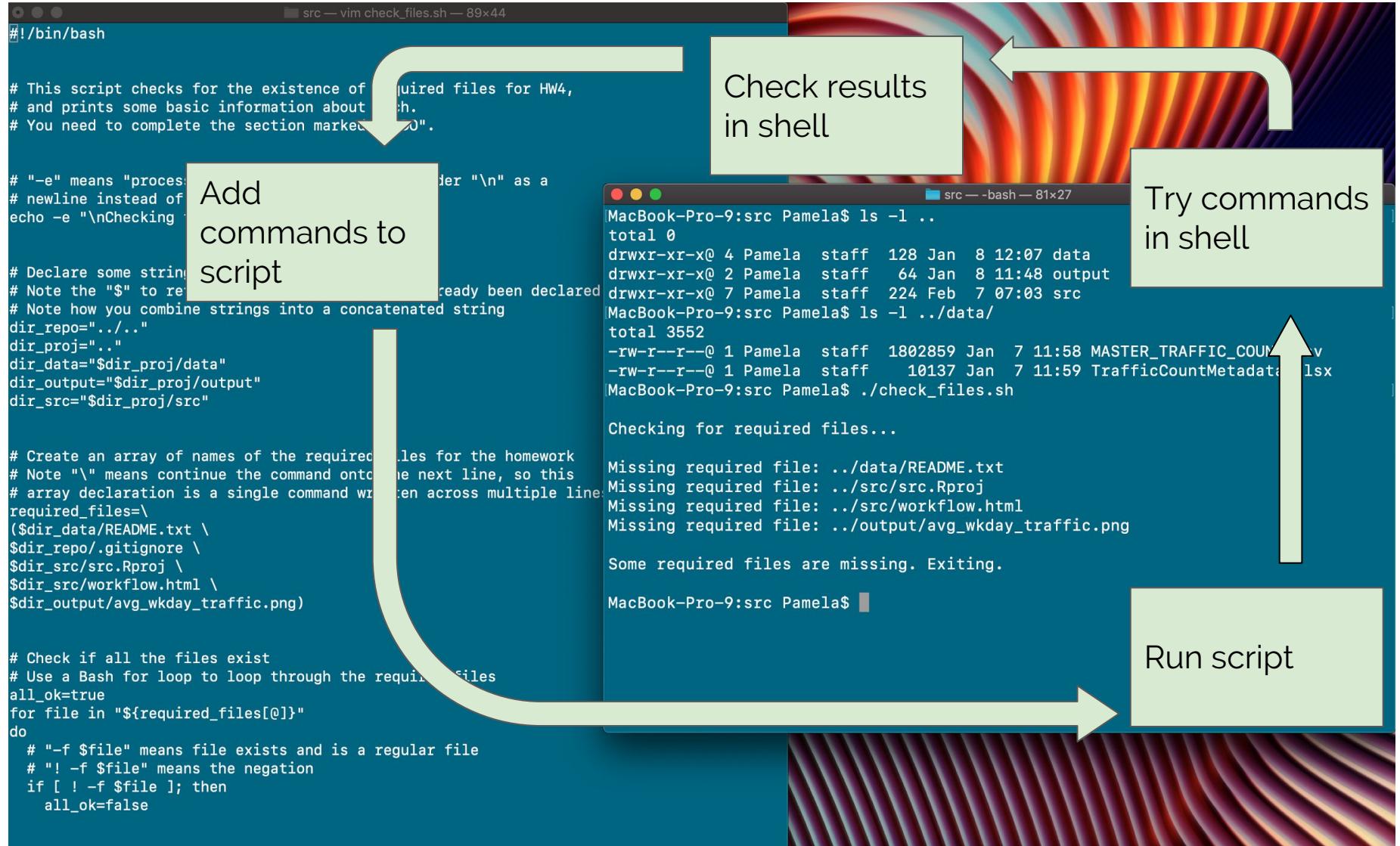
Bash script:

- Direct access to shell environment  
    => Efficient development cycle
- Readability
- Integration of command line programs
- Job/process control for complex pipelines

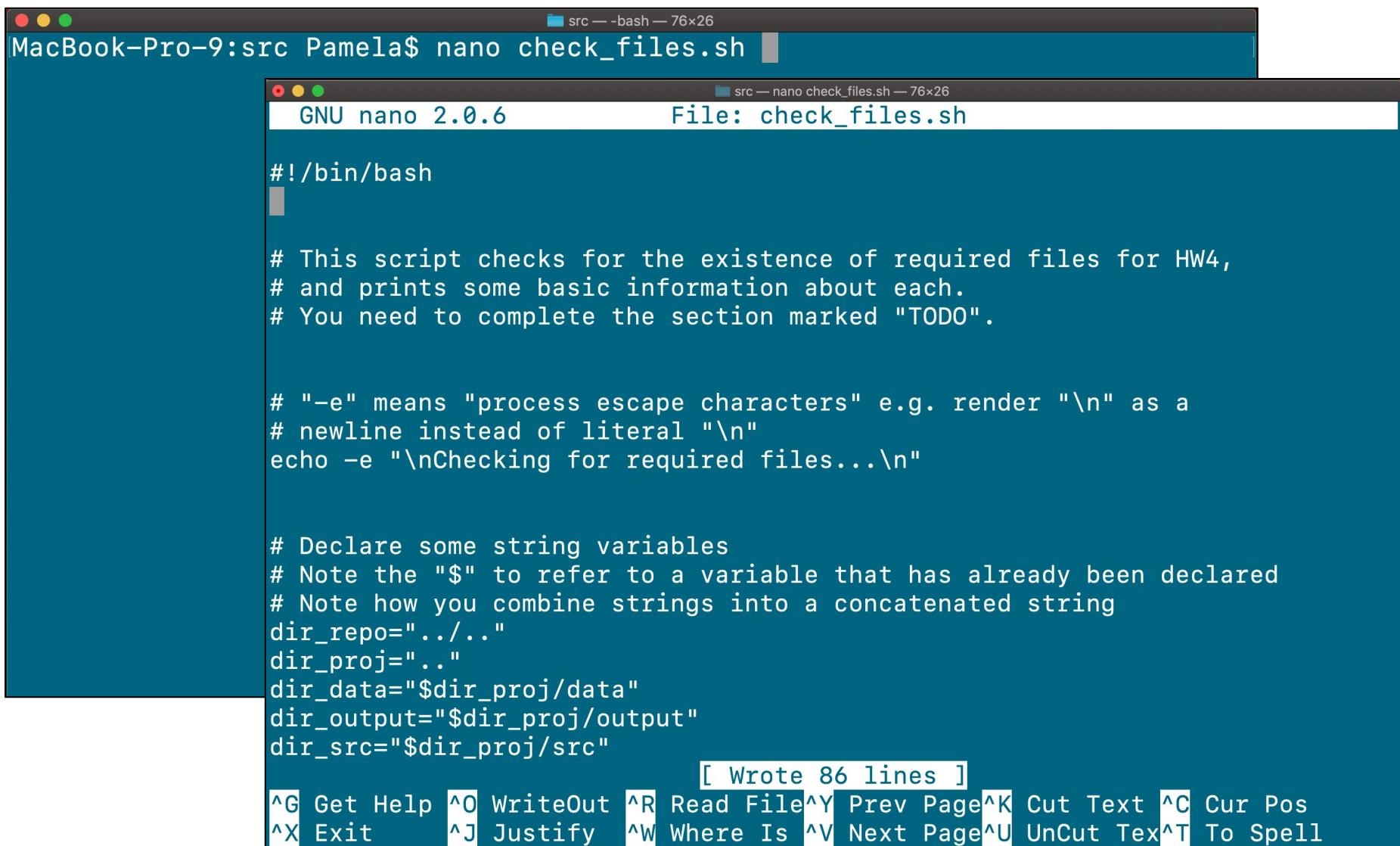
# R development cycle (Lecture 3)



# Bash development cycle



# Editing a text file from the shell: Nano editor



```
MacBook-Pro-9:src Pamela$ nano check_files.sh

GNU nano 2.0.6          File: check_files.sh

#!/bin/bash

# This script checks for the existence of required files for HW4,
# and prints some basic information about each.
# You need to complete the section marked "TODO".

# "-e" means "process escape characters" e.g. render "\n" as a
# newline instead of literal "\n"
echo -e "\nChecking for required files...\n"

# Declare some string variables
# Note the "$" to refer to a variable that has already been declared
# Note how you combine strings into a concatenated string
dir_repo=".../..."
dir_proj=...
dir_data="$dir_proj/data"
dir_output="$dir_proj/output"
dir_src="$dir_proj/src"
[ Wrote 86 lines ]
^G Get Help ^O WriteOut ^R Read File ^Y Prev Page ^K Cut Text ^C Cur Pos
^X Exit      ^J Justify   ^W Where Is ^V Next Page ^U UnCut Tex ^T To Spell
```

# Nano controls

The screenshot shows a terminal window titled "src — nano check\_files.sh — 76x26". The title bar also displays "GNU nano 2.0.6" and "File: check\_files.sh". The main area contains a shell script:

```
#!/bin/bash

# This script checks for the existence of required files for HW4,
# and prints some basic information about each.
# You need to complete the section marked "TODO".

# "-e" means "process escape characters" e.g. render "\n" as a
# newline instead of literal "\n"
echo -e "\nChecking for required files...\n"

# Declare some string variables
# Note the "$" to refer to a variable that has already been declared
# Note how you combine strings into a concatenated string
dir_repo="..../.."
dir_proj=..
dir_data="$dir_proj/data"
dir_output=''$dir_data'/'$dir_proj'_output"
[ Wrote 86 lines ]
```

At the bottom of the screen, there is a menu bar with the following options:

- Get Help
- ^O WriteOut
- ^R Read File
- ^Y Prev Page
- ^K Cut Text
- ^C Cur Pos
- ^X Exit
- ^J Justify
- ^W Where Is
- ^V Next Page
- ^U UnCut Tex
- ^T To Spell

Red boxes highlight several key commands: "Save" (over "WriteOut"), "Exit" (over "Get Help" and "Exit"), and "^X Exit".

# Graphical editors: convenient if you're working on your laptop

The screenshot shows the RStudio IDE interface. On the left, a terminal window displays a Bash script named `check_files.sh`. The script checks for required files and prints messages based on their presence. On the right, the RStudio environment is visible, showing the Global Environment pane which is currently empty, and the Files pane which lists the contents of the `src` directory of a project named `homework`.

```
~/.Dropbox/Documents/Teaching/bios6660_spring2019/homework - master - RStudio
+ Go to file/function G Addins

check_files.sh x
31 $dir_src/workflow.html \
32 $dir_output/avg_wkday_traffic.png)
33
34
35 # Check if all the files exist
36 # Use a Bash for loop to loop through the required files
37 all_ok=true
38 for file in "${required_files[@]}"
39 do
40     # "-f $file" means file exists and is a regular file
41     # "!" -f $file" means the negation
42     if [ ! -f $file ]; then
43         all_ok=false
44         echo "Missing required file: $file"
45     fi
46 done
47
48
49 # Print a message about whether all files were found
50 # If any files are missing, terminate the program here
51 if [ $all_ok = true ]; then
52     echo -e "\nFound all required files."
53 else
54     echo -e "\nSome required files are missing. Exiting.\n"
55     exit -1
56 fi
57
```

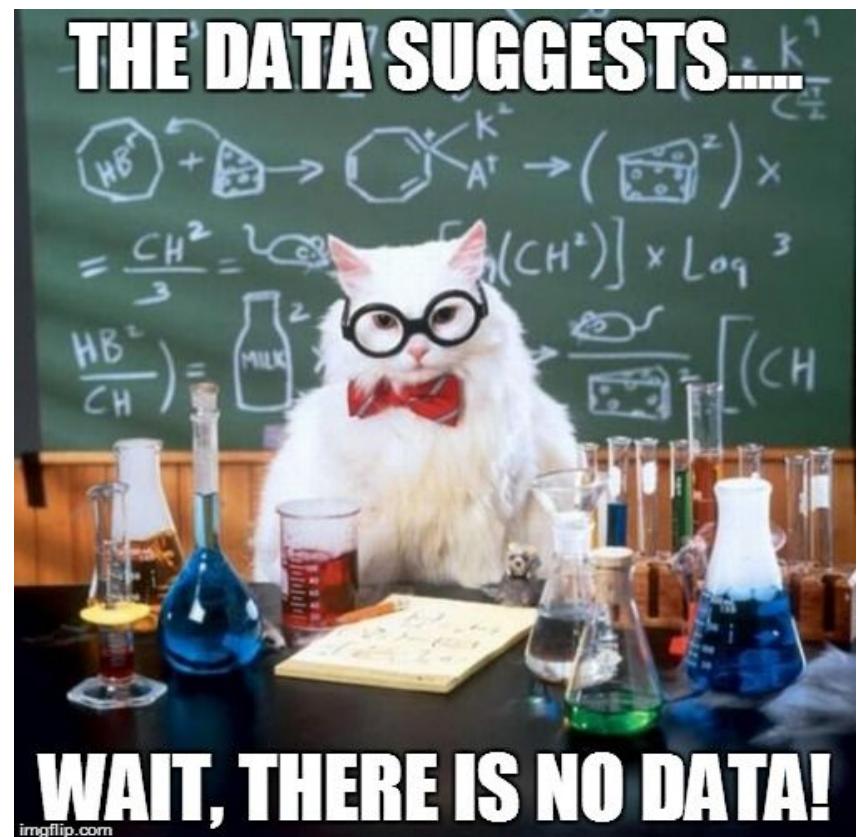
Environment is empty

| Name           | Size   |
|----------------|--------|
| check_files.sh | 2.2 KB |
| functions.R    | 1 KB   |
| tests.R        | 598 B  |
| workflow.Rmd   | 1.5 KB |

Console

# Data management

Lecture 8  
BIOS 6660, Spring 2019  
Instructor: Pam Russell



# Data: singular or plural?



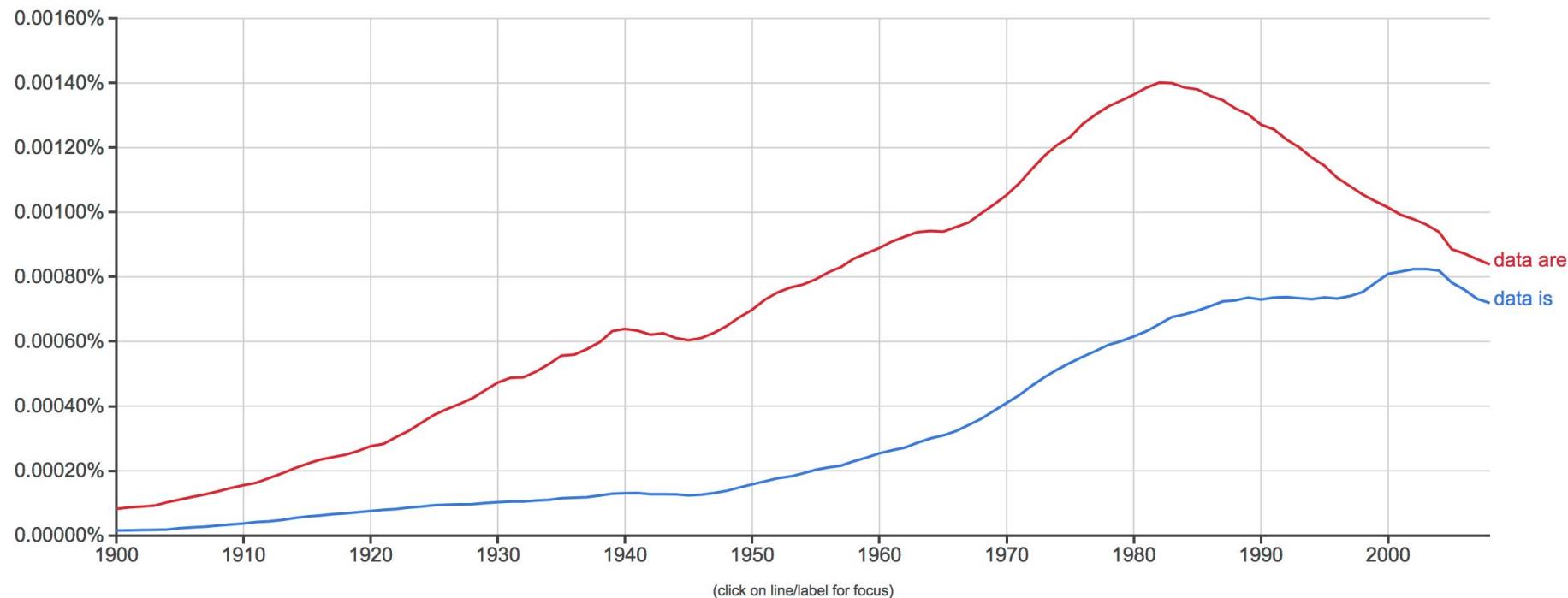
Guardian style guide:

"Data takes a singular verb (like agenda), though strictly a plural; you come across datum, the singular of data, about as often as you hear about an agendum"

# Google Books Ngram Viewer

Graph these comma-separated phrases:   case-insensitive

between  and  from the corpus English



# Google Books Ngram Viewer

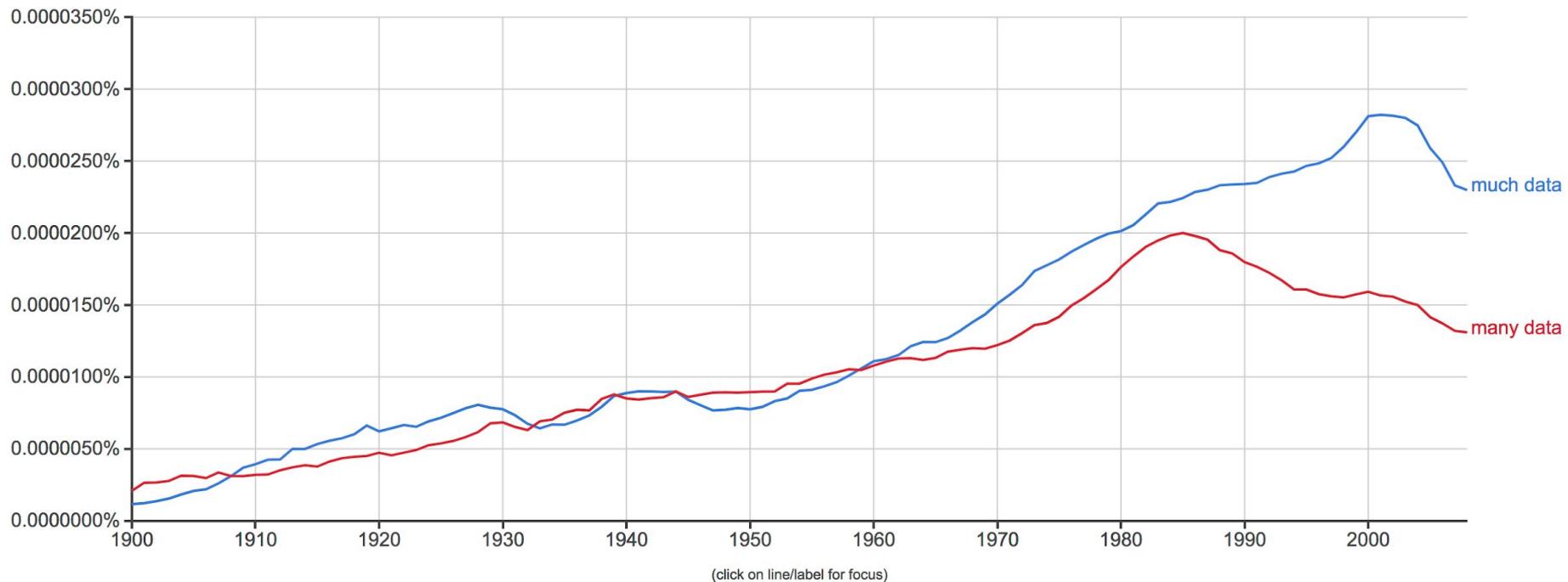
Graph these comma-separated phrases: much data,many data

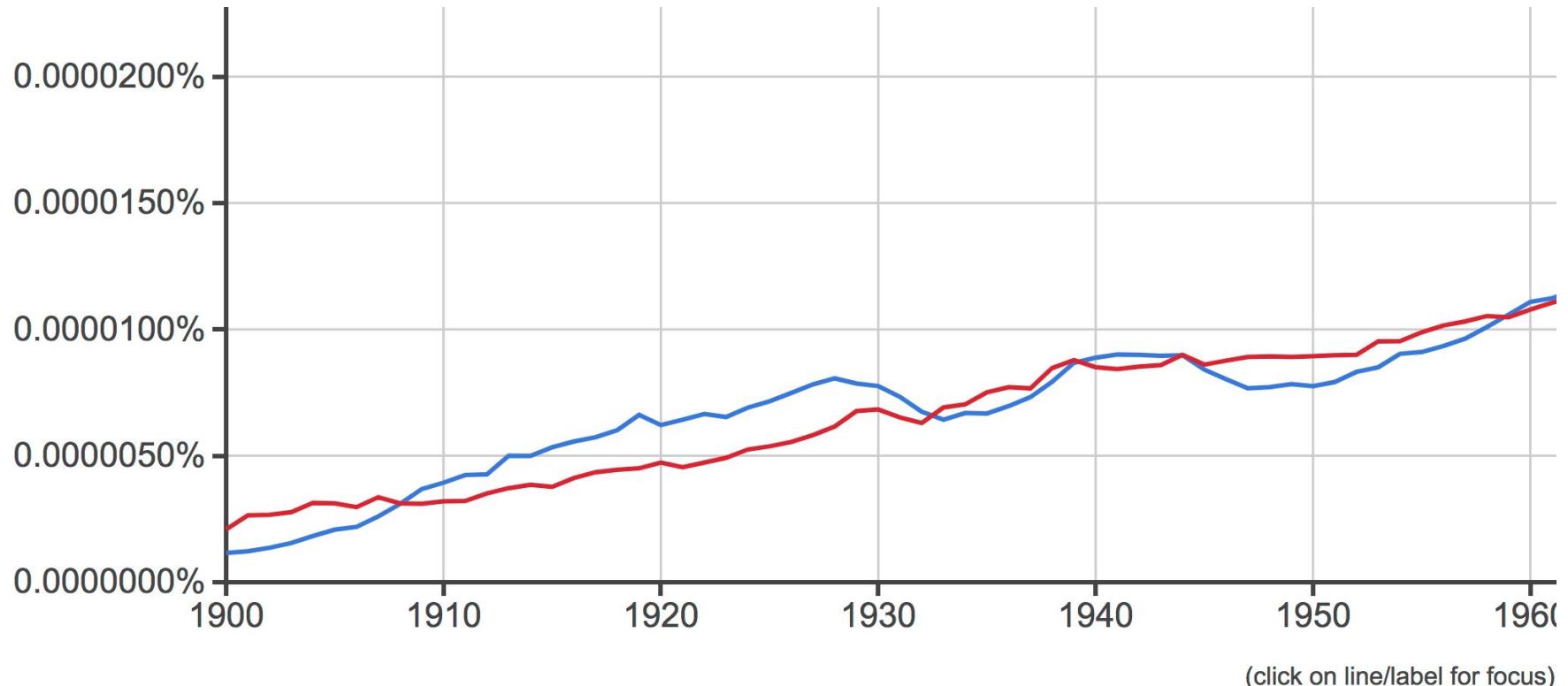
case-insensitive

between 1900 and 2008 from the corpus English

with smoothing of 3

**Search lots of books**





Search in Google Books:

[1900 - 1928](#)

[1900 - 1928](#)

[1929 - 1991](#)

[1929 - 1987](#)

[1992 - 1997](#)

[1988 - 1993](#)

[1998 - 2002](#)

[1994 - 2000](#)

[2003 - 2008](#)

[2001 - 2008](#)

Run your own experiment! Raw data is available for download [here](#).

# What is data management?

Handling, documentation, sharing, and preservation of research data

# Why data management?

## For you:

- Data easier to find, use, analyze
- Get credit (citations) for your data
- Minimal effort to share data

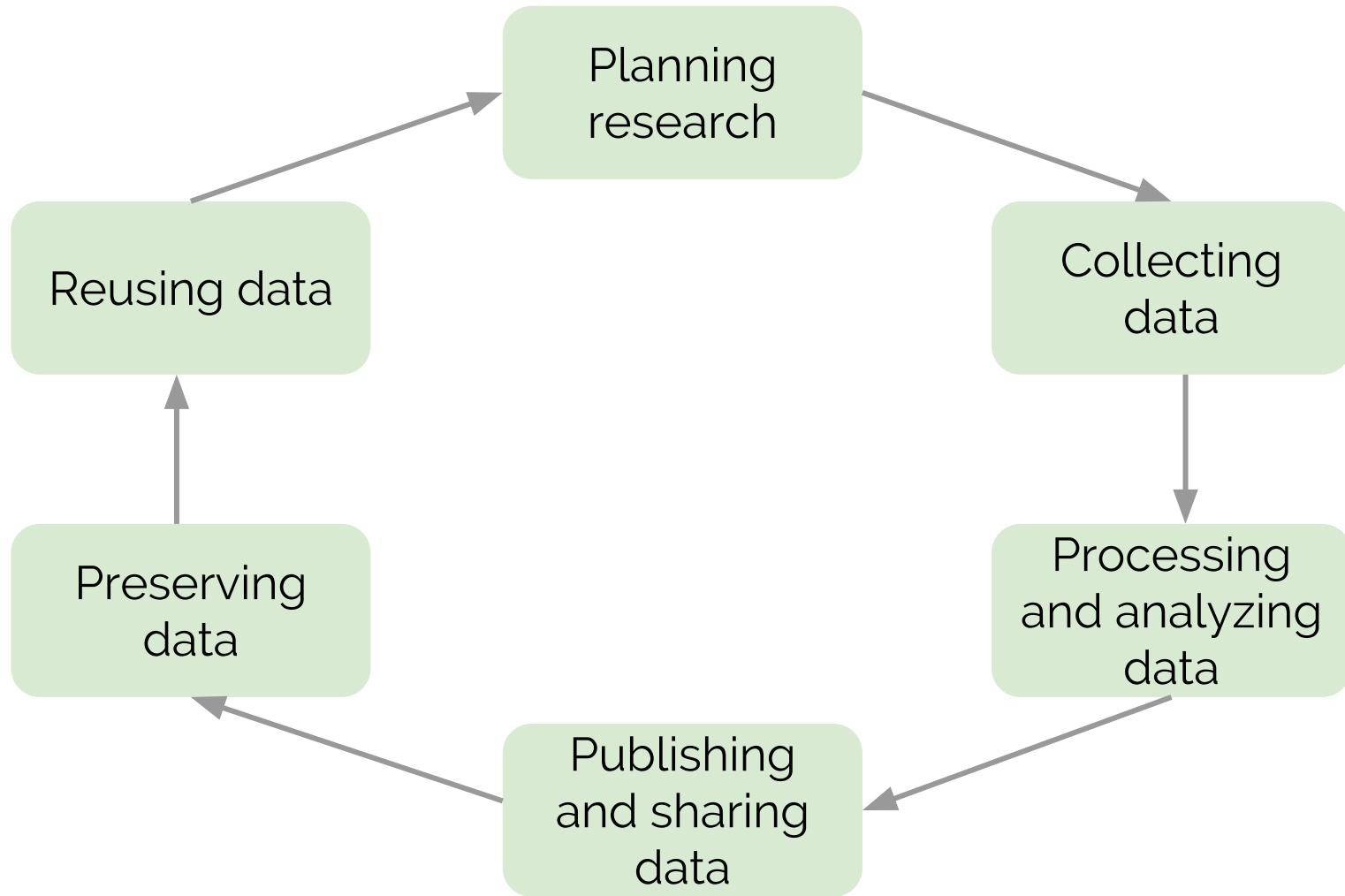
## For others:

- Easier for collaborators to understand and use
- Other scientists can easily find and use

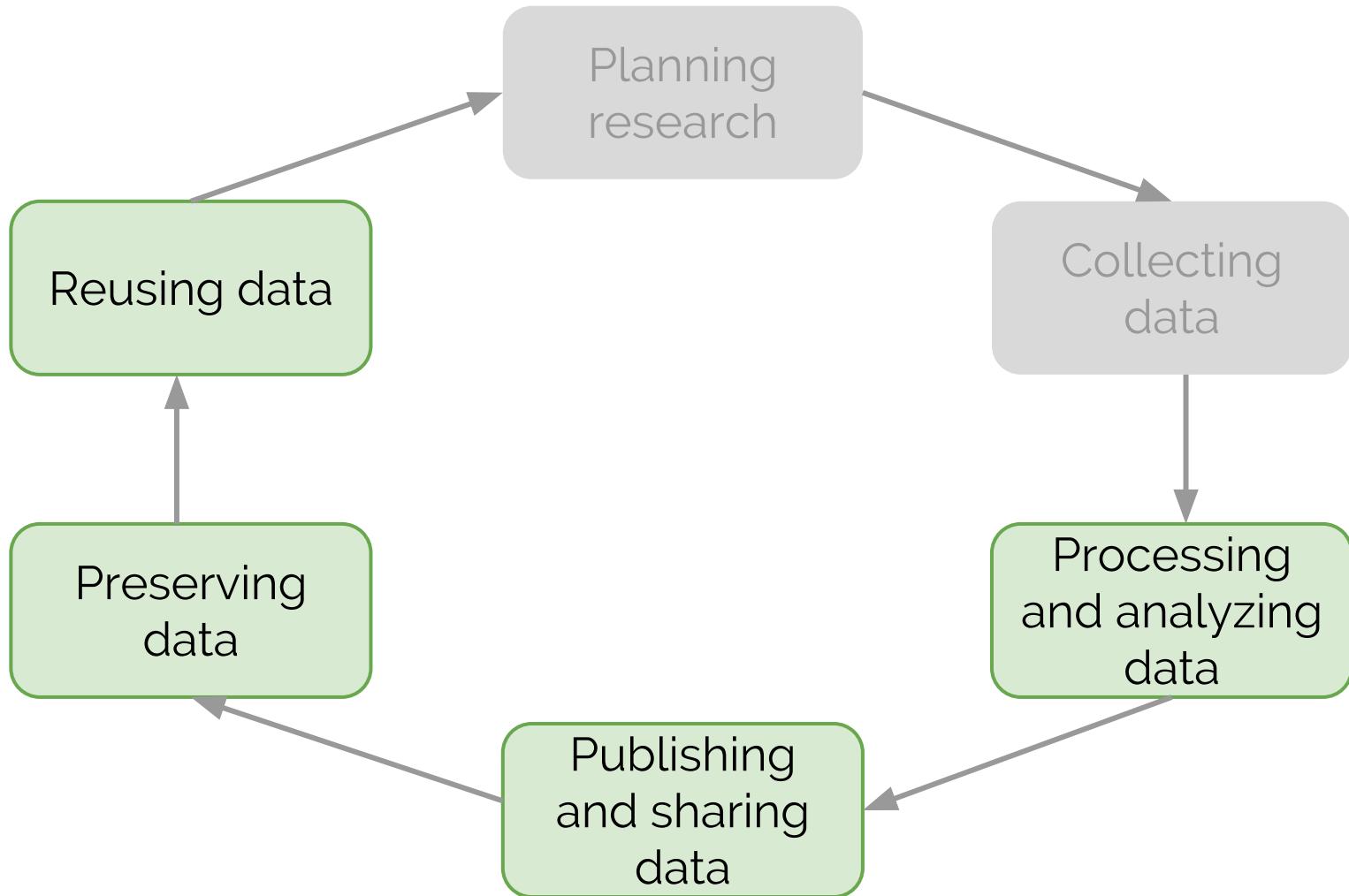
## For society:

- Loss of data deprives society of benefits of research
- Journals and funders increasingly require data management and sharing

# Research data lifecycle



# Research data lifecycle



# Not covered

- Confidentiality
- Ethics

Issues with “anonymous” data:

<https://datascience.berkeley.edu/anonymous-data>



# Data sharing

Making all underlying data available along with publication

In supplemental files, a data repository, or a field-specific database

# Why share data?

For you:

- Get credit (citations) for your data
- Feel good about contributing to science

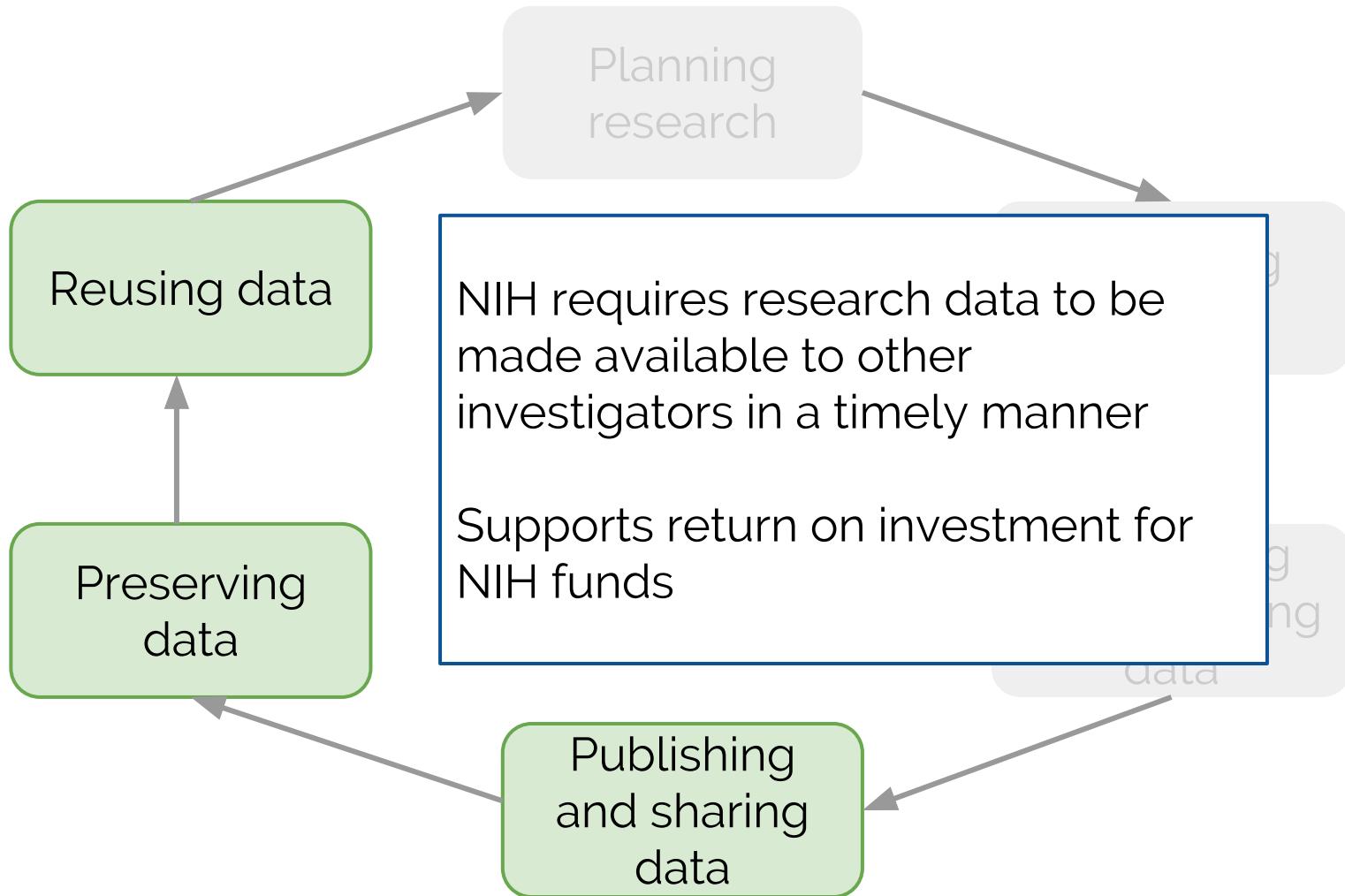
For others:

- Other scientists can easily find and use

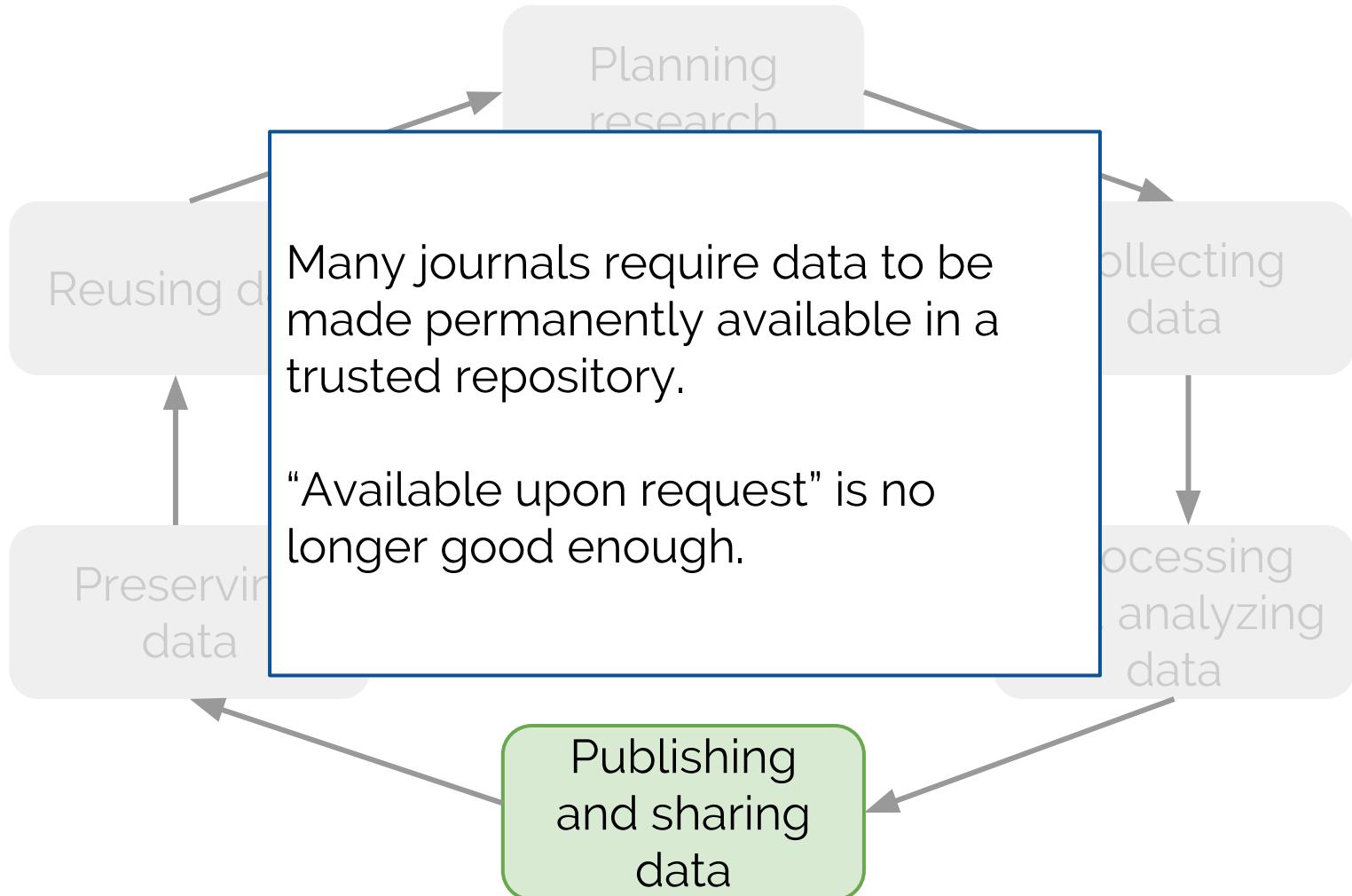
For society:

- Increased value from additional use of data
- Journal and funder requirements

# Funding agency requirements



# Journal requirements



# Data repositories

“A trusted digital repository is one whose mission is to provide reliable, long-term access to managed digital resources to its designated community, now and in the future.”

- *Trusted Digital Repositories*, RLG-OCLC report

Types of repositories:

- General purpose
- Disciplinary
- Institutional

# Data citation



“Citing research data in a manner similar to traditional scholarly works can help ensure proper attribution, improve reproducibility, improve discoverability, and help provide credit for data as a scholarly output.”

- Princeton Research Data Management

# Data citation

Repository provides persistent digital object identifier (DOI)

## PANGAEA

Willmes, S et al. (2009): Onset dates of annual snowmelt on Antarctic sea ice in 2007/2008. doi:10.1594/PANGAEA.701380 

## Dryad

Kingsolver JG, Hoekstra HE, Hoekstra JM, Berrigan D, Vignieri SN, Hill CE, Hoang A, Gibert P, Beerli P (2001) Data from: The strength of phenotypic selection in natural populations. Dryad Digital Repository.  
doi:10.5061/dryad.166 

## Dataverse

Frederico Girosi; Gary King, 2006, ‘Cause of Death Data’,  
<http://hdl.handle.net/1902.1/UOVMCPSWOL>  
UNF:3:9JU+SmVyHgwRhAKclQ85Cg== IQSS Dataverse Network  
[Distributor] V3 [Version].

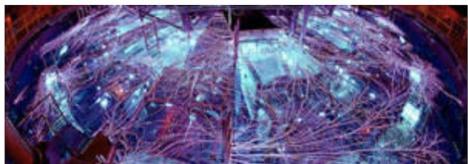
**Figure 2:** Data citation formats suggested by repositories

## Example: PLoS journals

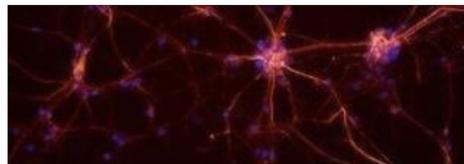
"PLOS is refining the definition, evaluation and recognition of scientific work, expanding toward a more Open Science ethos."



# Example: PLOS journals



PLOS encourages posting [preprints](#) to accelerate the dissemination of important results. Our partnership with bioRxiv makes posting life sciences preprints easy and convenient.



With their own experts at the helm, [PLOS Channels](#) provide scientific communities a single destination featuring curated research articles, commentary, blogs, news and discussions.



[Ambra](#)™, our journal and collections publishing platform and [Aperta](#)™, a submission platform developed here at PLOS, are available open source to the community under an MIT License.



Our partnership with [Protocols.io](#) enables authors to more effectively share methodological details about their research, supporting reproducibility.



A forward-thinking [Data Policy](#) fosters reproducibility, transparency and accelerates discovery.



Operating in tandem, [ORCID](#) and [CRediT](#) ensure that authorship is properly attributed, and each individual authors' contributions are fully described.

# PLOS data policy

“PLOS journals require authors to make all data underlying the findings described in their manuscript fully available without restriction, with rare exception.”

# Acceptable data sharing methods

| Method                                       | Notes                                          |
|----------------------------------------------|------------------------------------------------|
| Data deposition                              | Best method, strongly recommended              |
| Supporting information files                 | For smaller datasets                           |
| Available upon request or from a third party | Only if above methods are unethical or illegal |

# Recommended data repositories

## Cross-disciplinary repositories

- [Dryad Digital Repository](#)
- [figshare](#)
- [Harvard Dataverse Network](#)
- [Open Science Framework](#)
- [Zenodo](#)

## Repositories by type

|                                     |                                   |                                                   |
|-------------------------------------|-----------------------------------|---------------------------------------------------|
| <a href="#">Biochemistry</a>        | <a href="#">Neuroscience</a>      | <a href="#">Social Sciences</a>                   |
| <a href="#">Biomedical Sciences</a> | <a href="#">Omics</a>             | <a href="#">Structural Databases</a>              |
| <a href="#">Marine Sciences</a>     | <a href="#">Physical Sciences</a> | <a href="#">Taxonomic &amp; Species Diversity</a> |
| <a href="#">Model Organisms</a>     | <a href="#">Sequencing</a>        | <a href="#">Unstructured and/or Large Data</a>    |



Diverse perspectives on science and medicine

Search PLOS Blogs



STAFF BLOGS ▾

BLOGS BY TOPIC ▾

ABOUT PLOS BLOGS

CONTACT



# The Official PLOS Blog

About This Blog

Contact

Search This Blog



Boosting Open Science Hardware in an academic context: opportunities and challenges

Posted January 17, 2019 by PLOS in Guest blog, Open Science, Publishing



PLOS Board Appointments

Posted January 9, 2019 by Alison Mudditt in In the News, Open Access, Publishing

Sign up for PLOS Updates

Email Address (required)

I have read and agree to the terms of the [PLOS Privacy Policy](#) and hereby consent to send my personal information to PLOS.

Sign Up



peer review science publishing  
discovery collaboration PLOS journals  
impact diversity Open Access PLOS  
featured transparency Open  
Science advocacy researchers early  
career researcher preprints  
recognition bioRxiv science

<https://blogs.plos.org/plos/>

# PLOScasts



Episode List

About This Blog

Search This Blog



## Welcome to PLOScast

Here is the full list of PLOScasts, podcasts about scholarly publishing, the future of academia, and the changing experiences of scientists.

Explore PLOScast episodes:

- Episode 27: Ivan Oransky on the Value of Tracking Retractions
- Episode 26: How ECRs like Jessica Polka are reinventing science publishing
- Episode 25: Exploring the World of Medical Research: An interview featuring Lauren Maggio
- Episode 24: Check Your Stats: An interview featuring Michèle Nuijten
- Episode 23: Big Data in the Social Sciences: An interview featuring Ian Mulvany
- Episode 22: Building Taxonomies: An interview featuring Bob Kasenbach
- Episode 21: Why Open Research: An interview featuring Erin McKiernan
- Episode 20: Science Communication and Critique: An interview featuring Hilda Bastian
- Episode 19: Accessing Academic Sources on Wikipedia: An interview featuring Jake Orlowitz
- Episode 18: Science of Peer Review: An interview with Eamon Duede Part 2
- Episode 17: The Science of Science Part 1: An interview with Eamon Duede
- Episode 16: Open Source Science Communities: An interview with Abigail Cabunoc Mayes
- Special crossover episode with Open Science Radio
- Episode 15: The Power of Preprints: An interview with James Fraser

Sign up for PLOS Updates

Email Address (required)

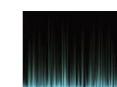
I have read and agree to the terms of the [PLOS Privacy Policy](#), and hereby consent to send my personal information to PLOS.

Sign Up

## Popular Posts



PLOScast Highlights from 2017  
posted on December 29, 2017



The History of Scientific...

# File naming and organization

“File names should provide context and be easily understood by humans and computers, now and in the future.”

- Meghan Frazer



# MIT Libraries recommendations

- Spend time planning folder hierarchy and file naming conventions in the beginning of a project
- Establish appropriate folder hierarchy for the project
- Develop file naming scheme that includes important metadata e.g. [Date]\_[Run]\_[Sample]
- Consider sorting. Dates should be YYYYMMDD
- Provide a README (more on this next)

# Metadata

Supporting documentation of the data

Without documentation, data can become unusable

Simplest way: keep in a README.txt file in the folder with the data

# **MIT Libraries recommendations: metadata attributes**

- Name of project or dataset
- Who created it
- Key dates
- Subject or content
- Funders
- Any intellectual property rights
- Methodology (how data was generated)
- Explanation of directory and file naming structure

# File formats

Try not to use proprietary formats (e.g. Microsoft Office)

If you must, then also create a plain text version or other appropriate open format

# GitHub and data

Never add protected data to GitHub, even  
a private repository

You might make the repo public and  
forget

Git makes it hard to permanently remove  
something from the history of a repo

# Use `.gitignore` to keep files out of version control

#discuss  20

## Show me your `.gitignore`



Maximilian Koch • Mar 3

DEV

# **Private alternatives to GitHub**

GitLab is like GitHub but private and internal to organizations

# Backups

Best way: institution's drive with regular backups

Don't store master copies on your laptop

External drives not recommended for long term (loss, bit rot)