# Code Review in Academia Workshop

Thursday
**13th March 2025**

**Room G.03
Bayes Centre**

The University of
Edinburgh,
EH8 9BT

```
∨ 11 ▪▪▪▪▪ README.md ⧉                               ⟨⟩  📄  ···
```

```
···   ···     @@ −1,2 +1,4 @@
```

**10am - 11am**

**+ Introductory presentation**

(attendee laptops not required)

**11.00 - 11:30**

**+ Tea/coffee break**

(please do not bring food/drink into this room G.03)

**11:30 - 12:30**

**+ Practice reviewing R or Python in small groups/pairs**

(attendee laptops required)

# Who we are



Amelia Edmondson-Stait
(post-doc)

Emily Ball
(post-doc)

Hannah Casey
(PhD student)

Ella Davyson
(PhD student)

Poppy Grimes
(PhD student)

- Psychiatry Department, in Prof Andrew McIntosh & Prof Heather Whalley's groups

- Apply epidemiological and statistical methods to large population cohorts and electronic health records.

- Advocators of good coding practices in academia.

- This workshop aims to initiate discussion and teach some concepts of code review.

- Funding for this workshop comes from the **Improving Research Community Builder Award**

# Workshop Expectations
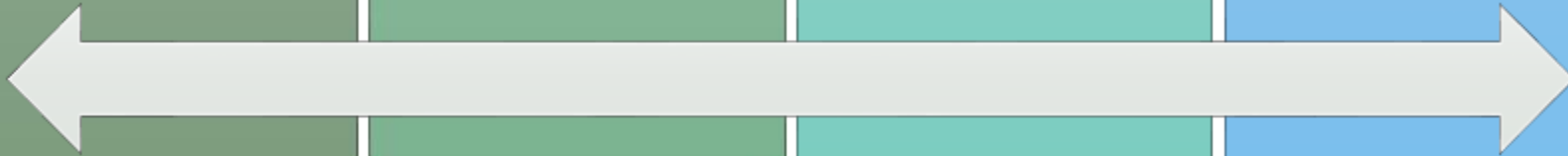
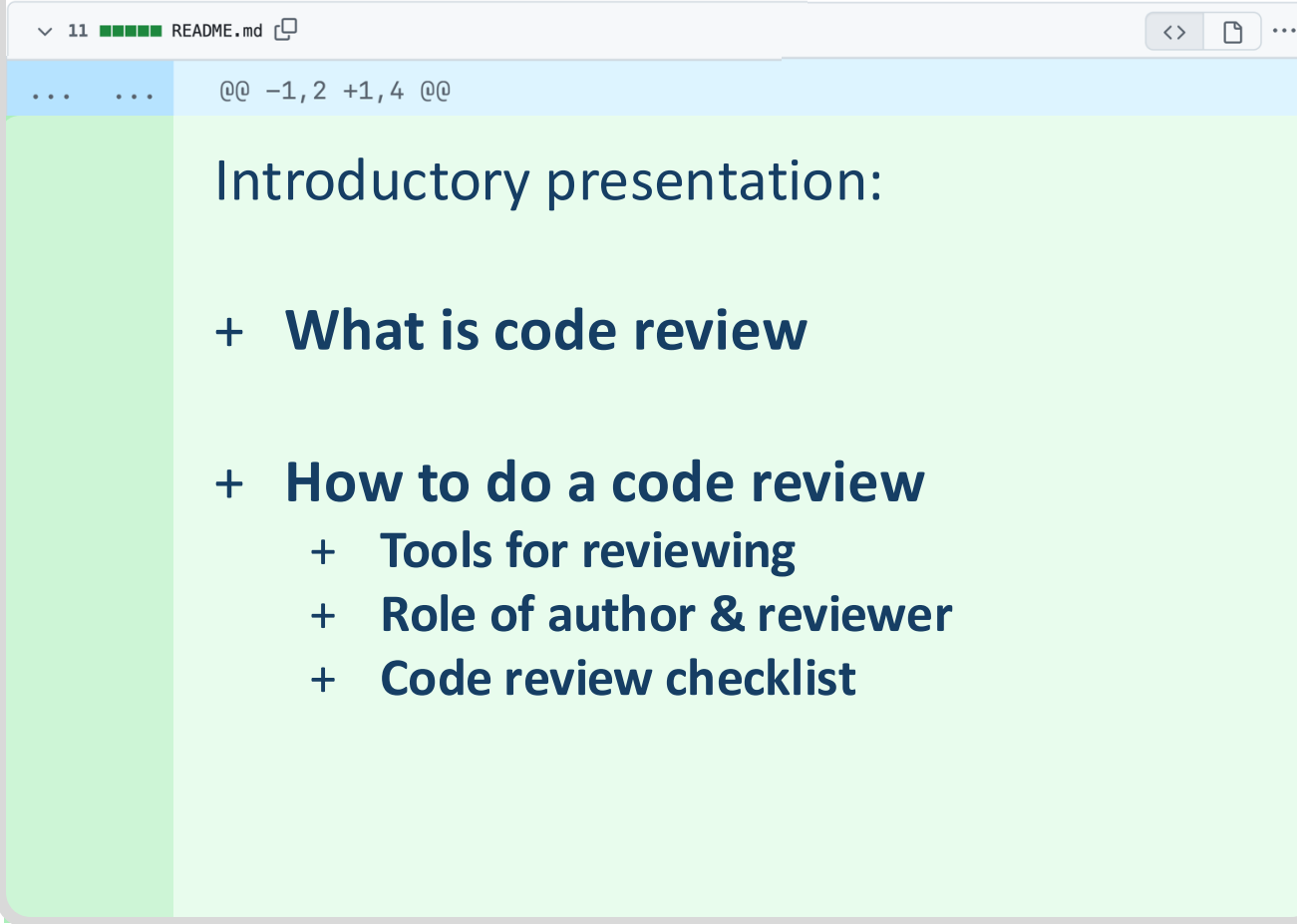This is a safe, supportive, respectful space

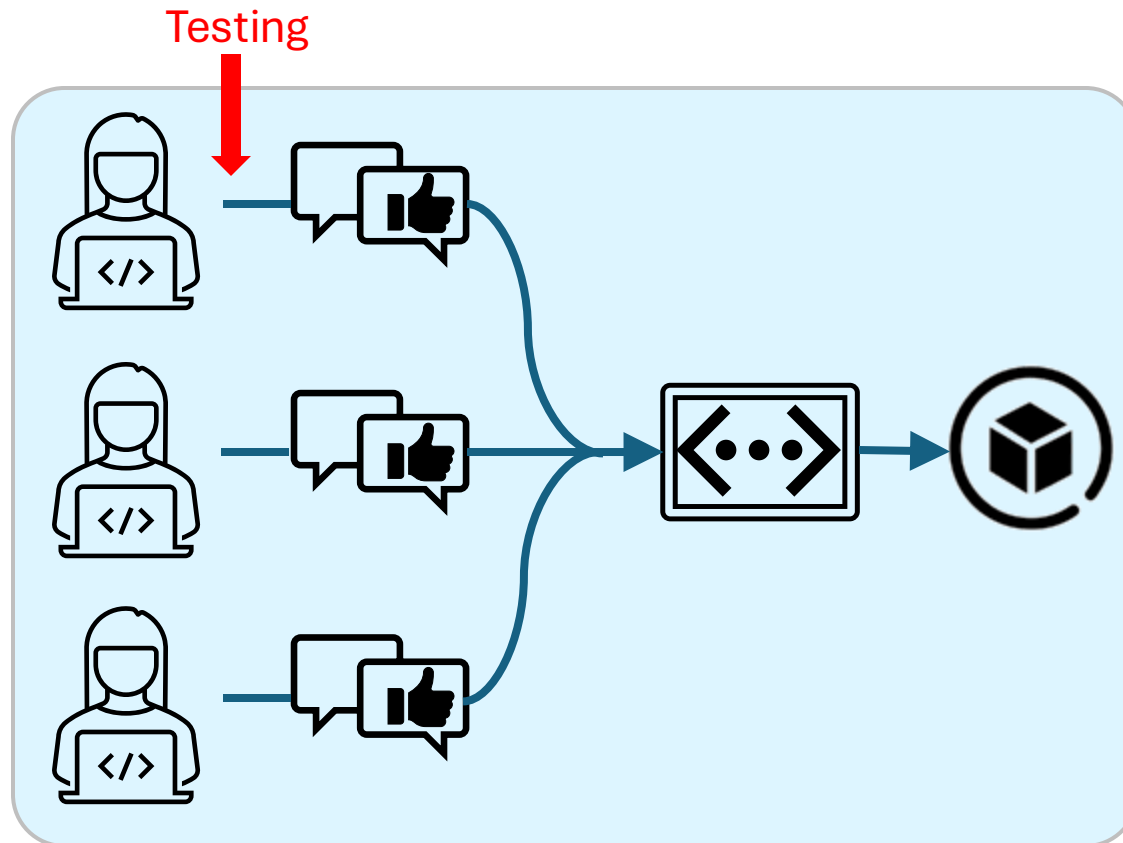Participation is encouraged

It's ok to make mistakes

We'll use a "Parking Sheet" for items that are off topic or taking up too much time
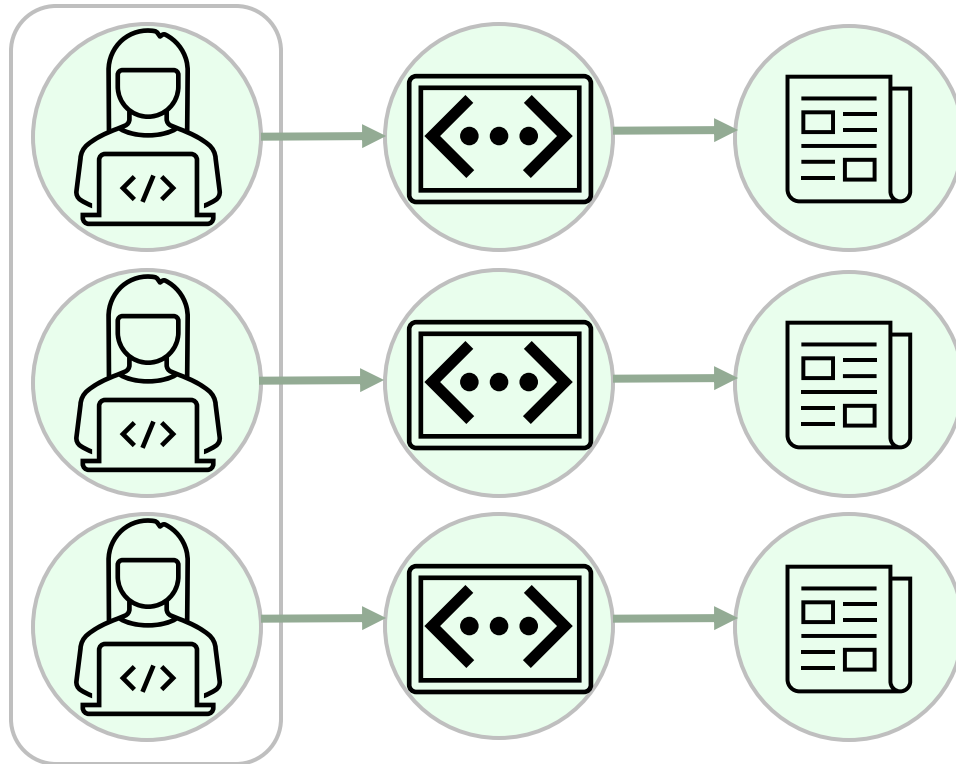
# Code Review in Academia Workshop

∨ 11 ▪▪▪▪▪ README.md  ⧉                                        ‹› ▯ ···

··· ···   @@ -1,2 +1,4 @@

Introductory presentation:

+ **What is code review**

+ **How to do a code review**
    + **Tools for reviewing**
    + **Role of author & reviewer**
    + **Code review checklist**
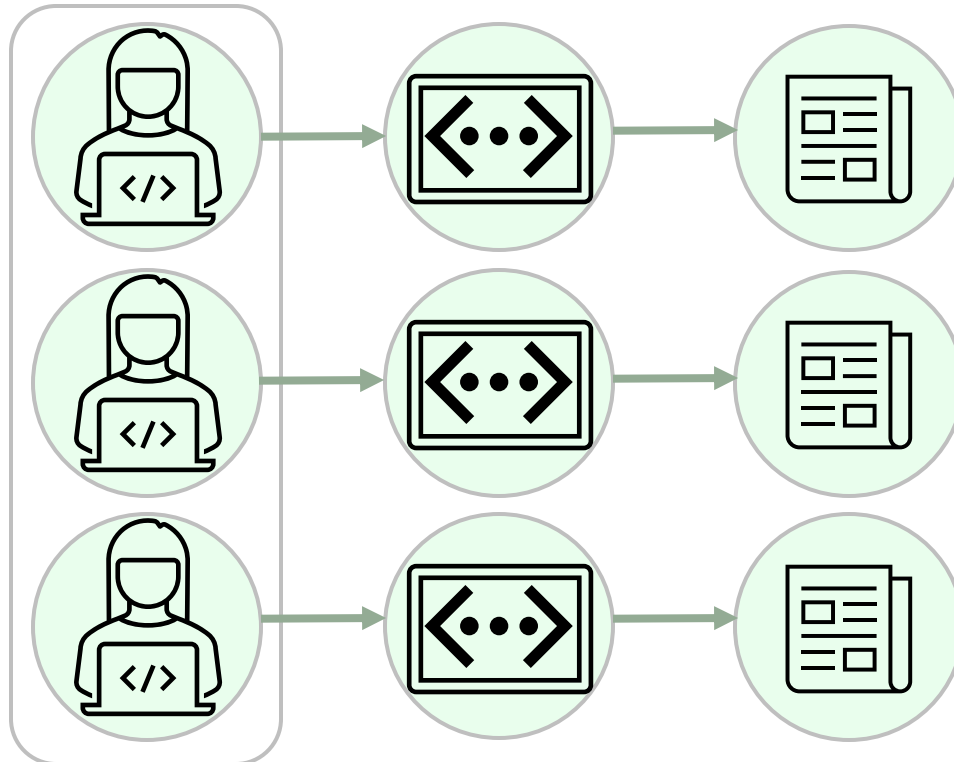
# Coding practices in industry



Testing

**Why code review is important in industry**

- Mistakes happen

- ↑ code quality

- Opportunity to learn

- Discover bugs earlier

- ↑ maintainability of code
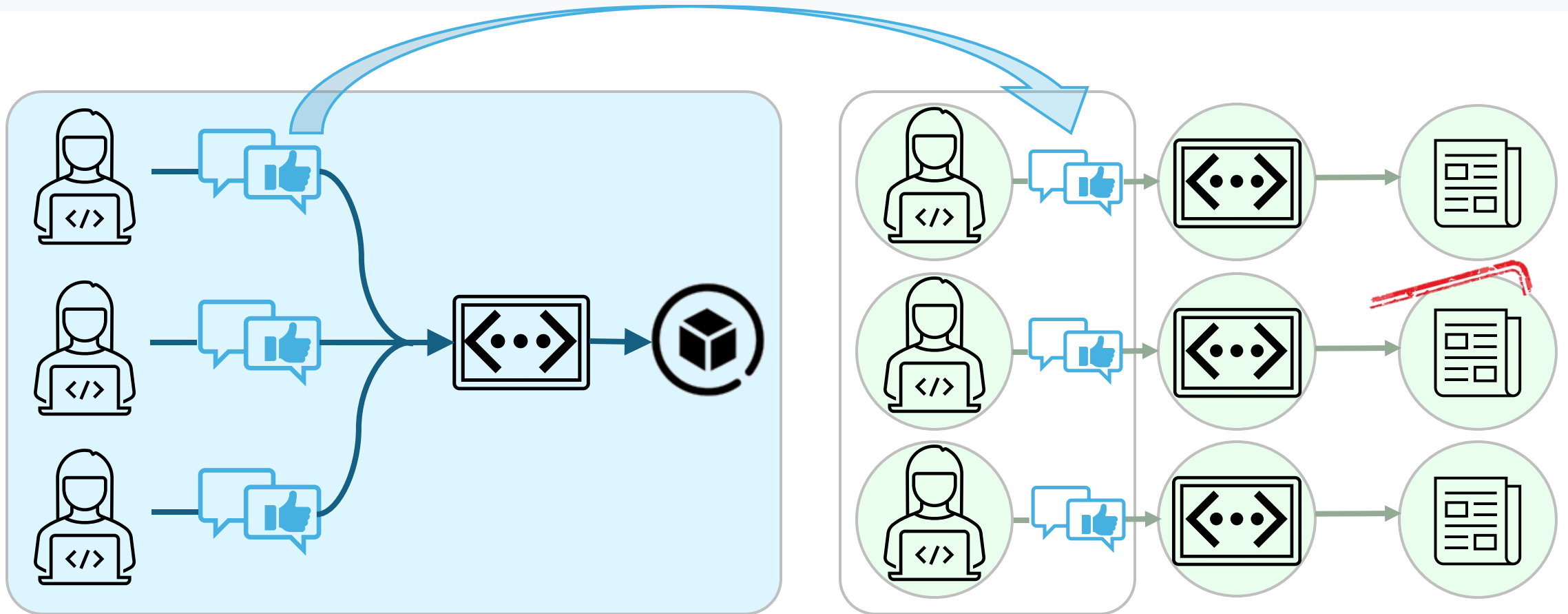
Coding practices in academia

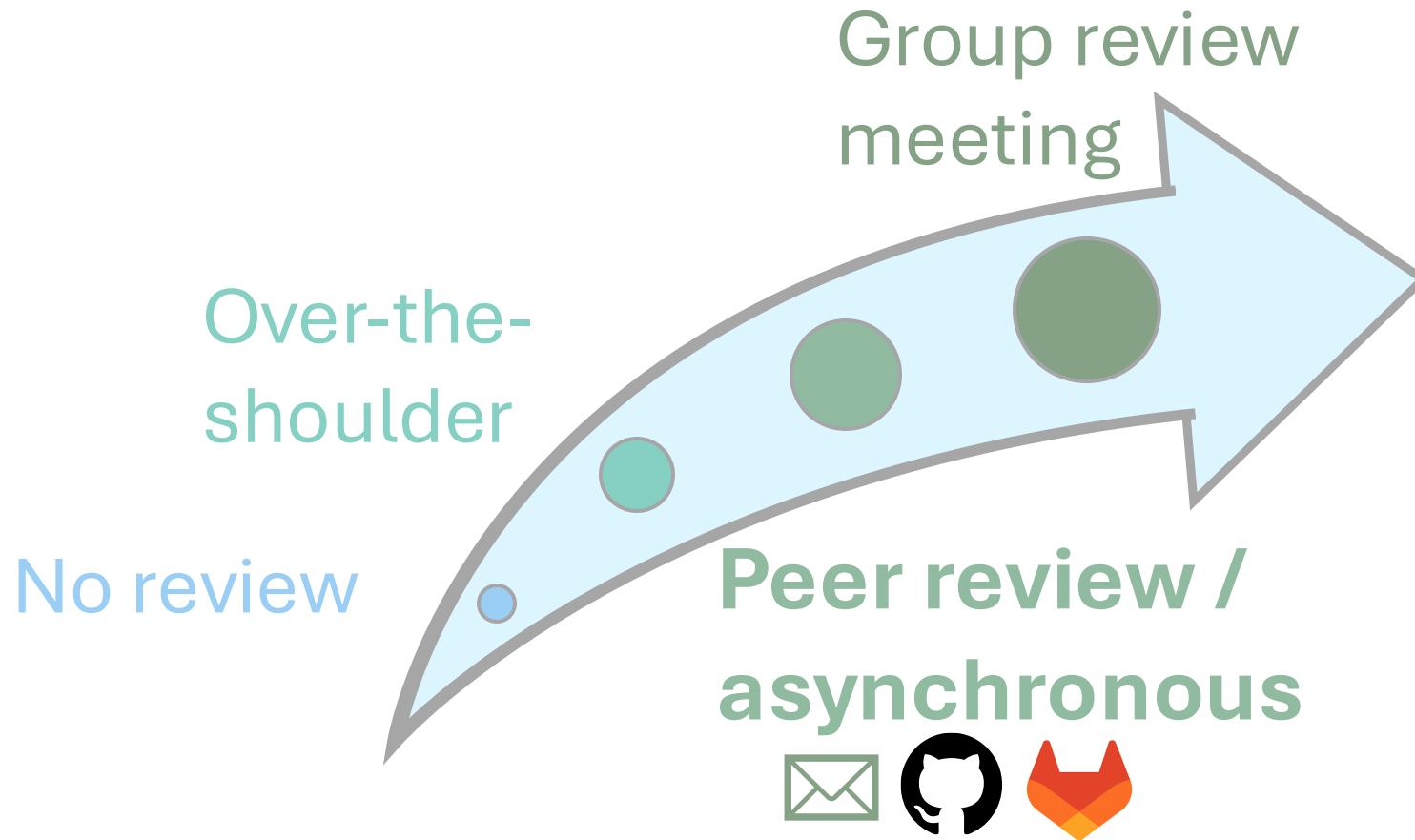# Coding practices in academia



**Why code review is important in academia**

- Mistakes happen

- ↑ code quality → reproducibilty

- ↑ research quality

- Opportunity to learn

- Discover bugs earlier

- ↑ maintainability → reproducibilty

# Apply industry methods to academia

# Types of code review



Group review meeting

Over-the-shoulder

No review

**Peer review / asynchronous**

# Summary so far...

- Coding in industry vs academia is different

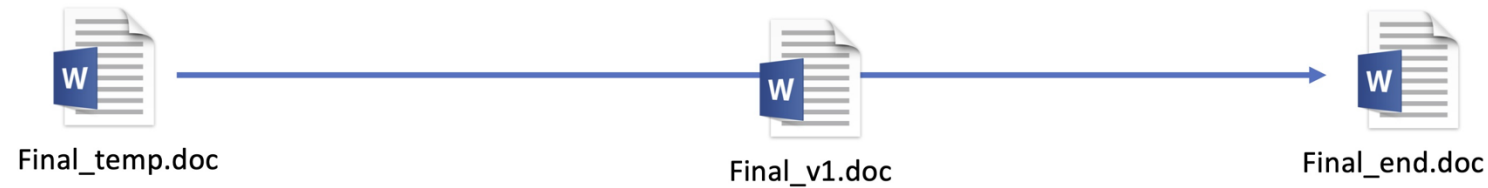- Code review is probably also going to be different

Next...
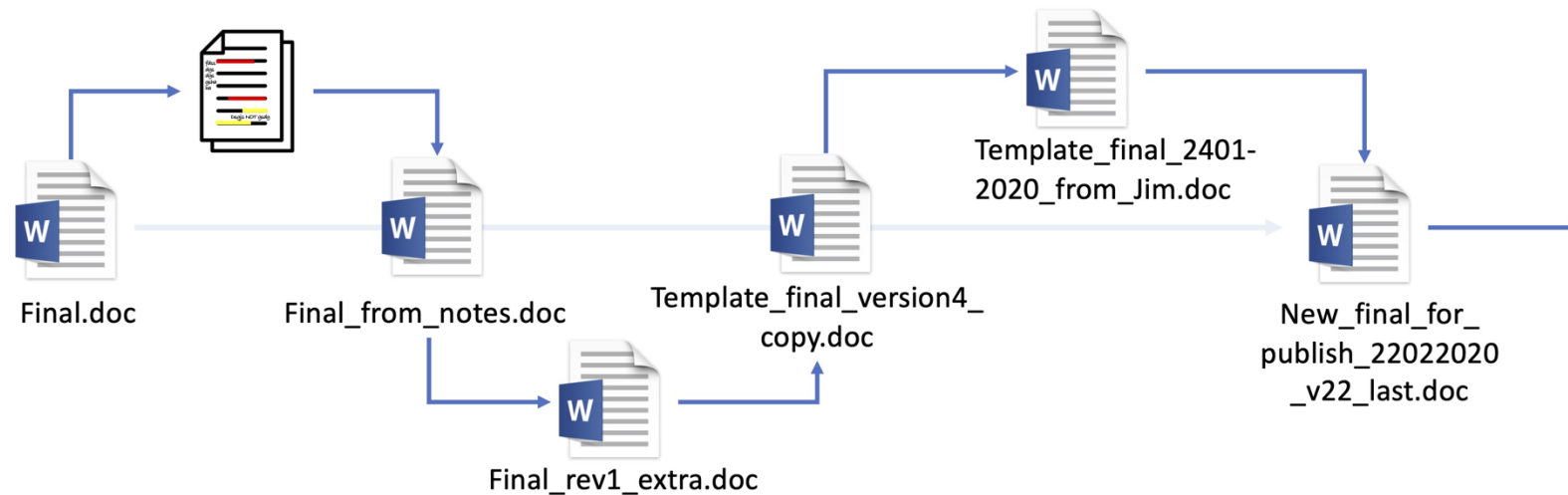
-    Methods used for asynchronous code review

# Industry uses "pull requests" for peer review

# What is git... git branches... pull requests?

**How I think simple version control would be**



Final_temp.doc      Final_v1.doc      Final_end.doc

**How version control actually is**



Final.doc    Final_from_notes.doc    Template_final_version4_copy.doc    Template_final_2401-2020_from_Jim.doc    New_final_for_publish_22022020_v22_last.doc

Final_rev1_extra.doc
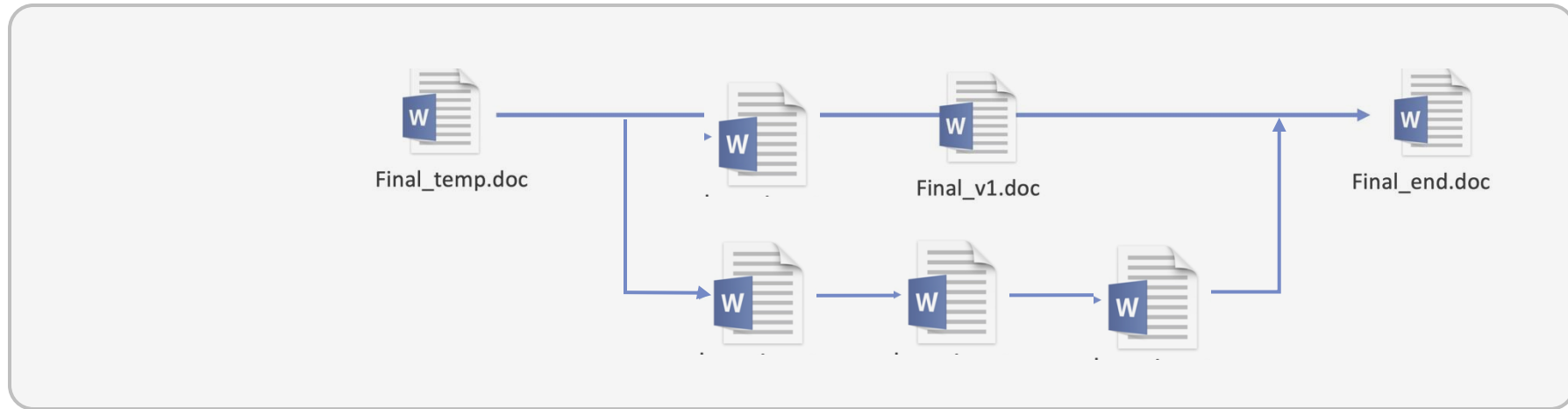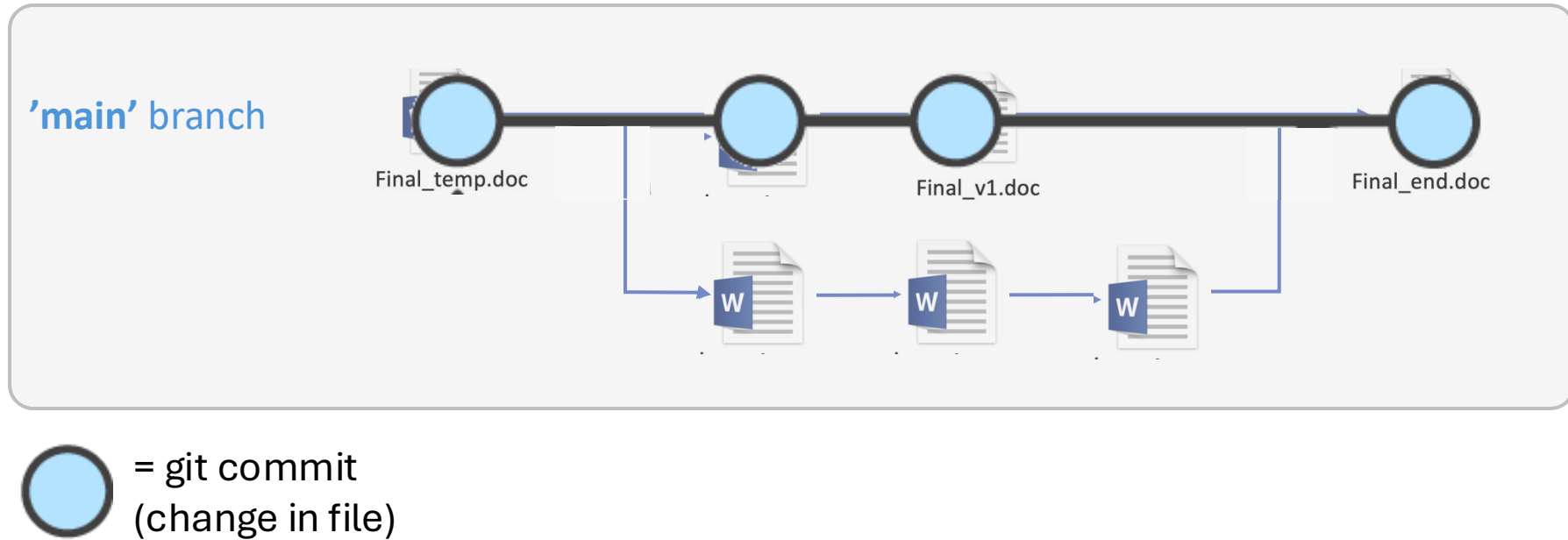
# What is git... git branches... pull requests?

**git = version control system**

# What is git... git branches... pull requests?

**git = version control system**

**'main'** branch

Final_temp.doc  Final_v1.doc  Final_end.doc



= git commit
(change in file)

# What is git... **git branches...** pull requests?

**git = version control system**



**'main'** branch

Final_temp.doc   Final_v1.doc   Final_end.doc

**'feature'** branch

◯ = git commit
(change in file)

# What is git… git branches… **pull requests?**

**'main'** branch

Create 'feature' branch from 'main'

Commit changes

**Submit Pull Request**

**Merge 'feature' branch into 'main'**

**Code reviewed/ commit changes**

**Pull requests are useful for reviewing changes made to an existing code base**

# Industry uses "pull requests" for peer review

**'main'** branch

Create 'feature' branch from 'main'

**Merge 'feature' branch into 'main'**

Commit changes

**Submit Pull Request**

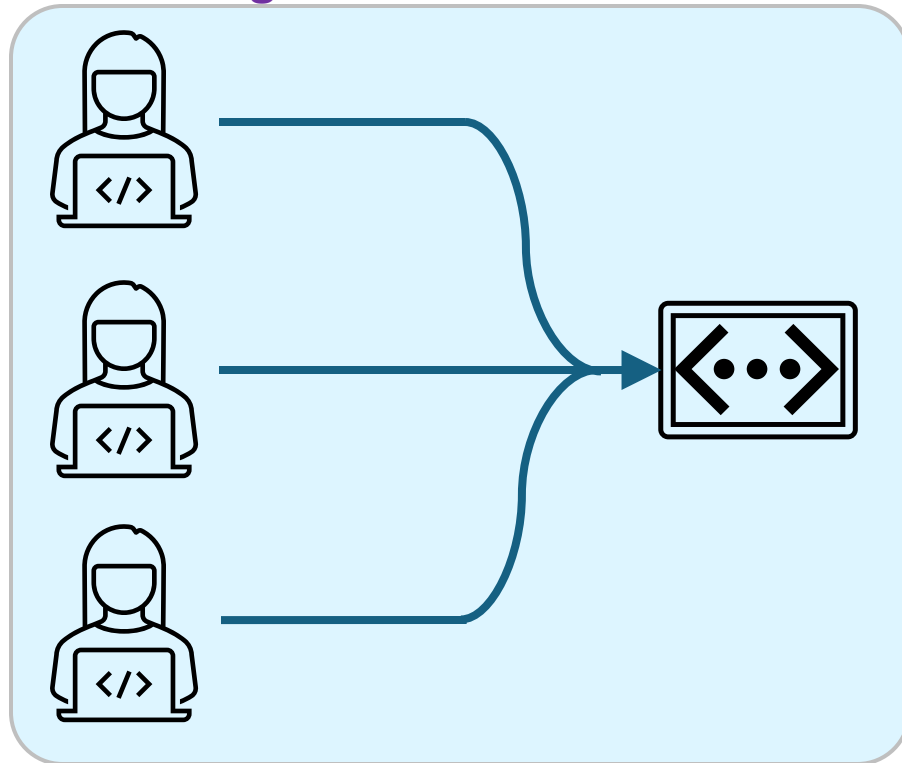**Code reviewed/ commit changes**

**Pull requests are useful for reviewing changes made to an existing code base**

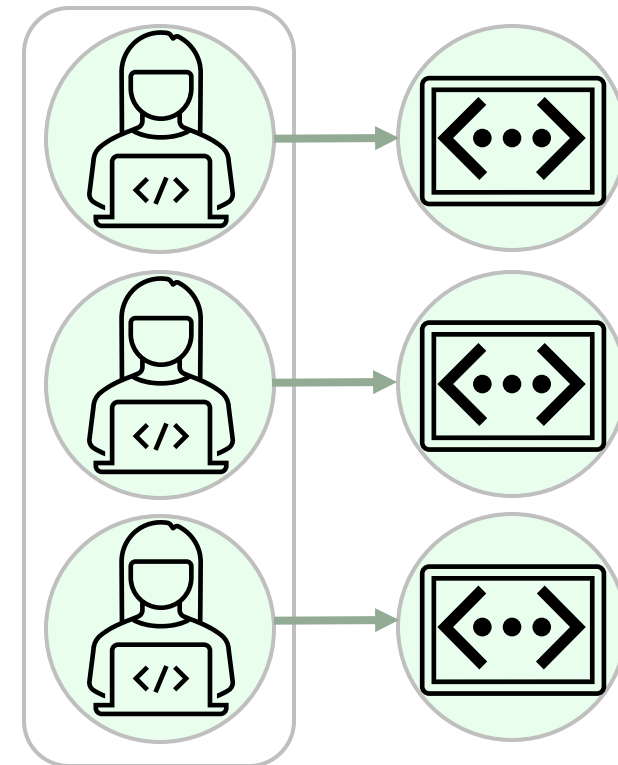# "Pull requests" for peer review in academia?

**Industry:**
Changes made to an existing code base
**Changes to code are reviewed**

**Academia:**
Each researcher/project = new code base
**Not just changes that need reviewing**

# Do we need something less advanced and more suitable for academia?

# Oxford Code Review Network

https://github.com/OxfordCodeReviewNet/forum

# Barriers for code review

# Barriers for code review



**Reviewing the code is more important**

**than the tools used to review**

# Personal experience of code review

- Didn't use pull requests or GitHub issues...

- My code was on a shared GitHub repository

- New coding language to me - nextflow

- My colleague looked at the script I was working on

- We had regular meetings where I would walkthrough/explain the code

- Restructured logic of code execution, learnt new functions, etc.

# Summary so far...

- Tools for code review balanced with reducing barriers

Next...

- Explain how a code review is done (role of author & reviewer)

- Things to look for in code when reviewing (checklist)

# Overview - How to do a code review

Role of code author

Role of code reviewer

```
Prepare code ──┐
               ├──→  Send code to reviewer
Find reviewer ─┘         ✉  GitHub
```

→ **Initial meeting?**

- Walk-through code
- Agree on expectations & objectives

→ Review code → Provide feedback → Improvements & follow up

# Preparing code for review

- Identify **50 to 200 lines** of code you want reviewed

- Make it easy to review: include **comments and documentation**

- Decide the **focus of the review** in advance

- Ask for **specific feedback** if needed

- Be prepared to **walk through your code** in the logical order of execution

</>  </>  </>

# Finding a reviewer

- Any researcher that codes is a good candidate

- Any coding experience is useful

- Colleagues in your research group,
      collaborators,
          mentors,
              students...

# How to do a code review

Role of code reviewer

# How to do a code review

- **Simply ask questions**

- Get the code author to explain their code in detail..
  *"Can you help me understand what this line of code does?"*

- Leads to them finding bugs & areas of improvements themselves

*Slide adapted from materials: https://github.com/OxfordCodeReviewNet/forum/blob/master/guidelines_for_reviewers.md*

# How to do a code review

- **Check code logic**

- Does the code do what is expected?

# How to do a code review

Be kind, give **personal opinions**
rather than imperative statements

*"I think this function's name could be improved"*
*NOT*
*"You should rename this function"*

# Variable names

Use descriptive names that convey intent



```
1  def calc(a, b):
2      return a * b
```



```
1  def calculate_area(width, height):
2      return width * height
```

*Slide adapted from materials: https://github.com/OxfordCodeReviewNet/forum/blob/master/guidelines_for_reviewers.md*

# Hard coded values

Avoid "magic numbers"

👎
```
1   for (int i = 0; i < 26; i++) {
```

👍
```
1   alphabetLength = 26;
2   for (int i = 0; i < alphabetLength; i++) {
```

👍👍
```
1   alphabetLength = alphabetData.size();
2   for (int i = 0; i < alphabetLength; i++) {
```

# **Duplicated code**

Don't repeat yourself (DRY)

```
1  dfA <- filter(df, group == "A")
2  analysisA <- t.test(dv~condition, data = dfA)
3
4  dfB <- filter(df, group == "B")
5  analysisB <- t.test(dv~condition, data = dfB)
```

```
1  subtest <- function(data, level) {
2    sub_df <- filter(data, group == level)
3    t.test(dv~condition, data = sub_df)
4  }
5
6  analysisA <- subtest(df, "A")
7  analysisB <- subtest(df, "B")
```

*Slide adapted from materials: CC-BY: DeBruine, Lisa (2022) Intro to Code Review. https://debruine.github.io/code-review/*

# Complex if else statements

Flatten nested conditional statements with guard clauses

```
 1  calculate_value <- function(x) {
 2    if (x > 0) {
 3      if (x < 10) {
 4        return(x * 2)
 5      } else {
 6        if (x < 20) {
 7          return(x * 3)
 8        } else {
 9          return(x * 4)
10        }
11      }
12    } else {
13      return(0)
14    }
15  }
```

```
 1  calculate_value <- function(x) {
 2    if (x <= 0) return(0)
 3    if (x < 10) return(x * 2)
 4    if (x < 20) return(x * 3)
 5    return(x * 4)
 6  }
```

# Long functions

Functions should be short and do one thing

```
1  long_function <- function() {
2    # Step 1: Do a lot of things
3    # Step 2: More things
4    # Step 3: Even more things
5  }
```

```
1  step_one <- function() { ... }
2  step_two <- function() { ... }
3  step_three <- function() { ... }
4
5  long_function <- function() {
6    step_one()
7    step_two()
8    step_three()
9  }
```

*Slide adapted from materials: https://github.com/OxfordCodeReviewNet/forum/blob/master/guidelines_for_reviewers.md*

# Obscure lines

Resist clever one-liners



```
1  result = reduce(sum, map(square, filter(positive, map(double, data))))
```



```
1  filtered_data = filter(positive, data)
2  doubled_data = map(double, filtered_data)
3  squared_data = map(square, doubled_data)
4  result = reduce(sum, squared_data)
```

# File paths

All file references should use relative paths, not absolute paths.

👎
```
1  patientRecords <- read.csv("C:/Users/Username/project/data/patientRecords.csv")
```

👍
```
1  patientRecords <- read.csv("data/patientRecords.csv")
```

"here" R package: https://here.r-lib.org/

# File names

Name files so both people and computers can easily find things

👎

myabstract.docx
Joe's Filenames Use Spaces and Punctuation.xlsx
figure 1.png
fig 2.png
JW7d^(2sl@deletethisandyourcareerisoverWx2*.txt

three principles for (file) names

machine readable

human readable

👍

2014-06-08_abstract-for-sla.docx
joes-filenames-are-getting-better.xlsx
fig01_scatterplot-talk-length-vs-interest.png
fig02_histogram-talk-attendance.png
1986-01-28_raw-data-from-challenger-o-rings.txt

plays well with default ordering

Use YYYY-MM-DD format for dates

# Unintended behaviour

Validate inputs to prevent unintended behaviour or errors

👎
```
1  calculate_scaled_log <- function(value) {
2    log_value <- log(value)
3    scaled_value <- log_value * 10
4    return(scaled_value)
5  }
```

👍
```
1  calculate_scaled_log_good <- function(value) {
2    if (value <= 0) {
3      stop("Input must be a positive number for log()")
4    }
5    log_value <- log(value)
6    scaled_value <- log_value * 10
7    return(scaled_value)
8  }
```

*Slide adapted from materials: https://github.com/OxfordCodeReviewNet/forum/blob/master/guidelines_for_reviewers.md*

# Comments

Explain the "why" not the "what"

- Redundant comments
- Complicated comments

```
1   # Subtract the mean age from age
2   centeredAge <- data$age - mean(data$age)
```

- Warnings of consequences
- Assumptions made

```
1   # Mean-center age to improve interpretation,
2   # reduce multicollinearity, and better model
3   # individual age-related changes over time
4   # in longitudinal trajectories.
```

*Further info on writing good comments here:* https://stackoverflow.blog/2021/12/23/best-practices-for-writing-code-comments/

# Documented code

Functions and classes should contain docstrings

```
 1  #' Descending order
 2  #'
 3  #' Transform a vector into a format that will be sorted in descending order.
 4  #' This is useful within [arrange()].
 5  #'
 6  #' @param x vector to transform
 7  #' @export
 8  #' @examples
 9  #' desc(1:10)
10  #' desc(factor(letters))
11  #'
12  #' first_day <- seq(as.Date("1910/1/1"), as.Date("1920/1/1"), "years")
13  #' desc(first_day)
14  #'
15  #' starwars %>% arrange(desc(mass))
16  desc <- function(x) {
17    obj_check_vector(x)
18    -xtfrm(x)
19  }
```

*Slide adapted from materials: https://github.com/OxfordCodeReviewNet/forum/blob/master/guidelines_for_reviewers.md*

# Documented code

Functions and classes should contain docstrings

```
1  def desc(x):
2      """
3      Descending order
4      Transform a vector into a format that will be
       sorted in descending order.
5      This is useful within [arrange()].
6
7      Parameters:
8      x (array-like): The vector to transform.
9
10     Returns:
11     array-like: A transformed version of `x` for
       descending sorting.
12
13     Example:
14     >>> desc([1, 2, 3])
15     [-1, -2, -3]
16     """
17     # function code
18     return x
```

*Slide adapted from materials: https://github.com/OxfordCodeReviewNet/forum/blob/master/guidelines_for_reviewers.md*

# Documented data

Be clear, consistent, and provide context

- Names (i.e., the column names)
- Labels/description
- Codings (e.g., 1 = always, 5 = never)
- Data type (e.g., binary, continuous)
- Descriptives (e.g., min, max)
- Data units (e.g., mg/L, months)
- Missing values (e.g., NA, 999)

# Coding Style

"Good coding style is like correct punctuation: you can manage without it, butitsuremakesthingseasiertoread."

- Notation and naming



snake_case

kebab-case

- Syntax (spacing, indentations, line length)
- Commenting guidelines
- And more...
  Google style guide on many languages: *https://google.github.io/styleguide/*
  *R style guide: https://style.tidyverse.org/*

*Slide adapted from materials: https://github.com/OxfordCodeReviewNet/forum/blob/master/guidelines_for_reviewers.md*

# Automate what can be automated

Code author can use a linter to automate coding style checks before review



- **R –** lintr R package

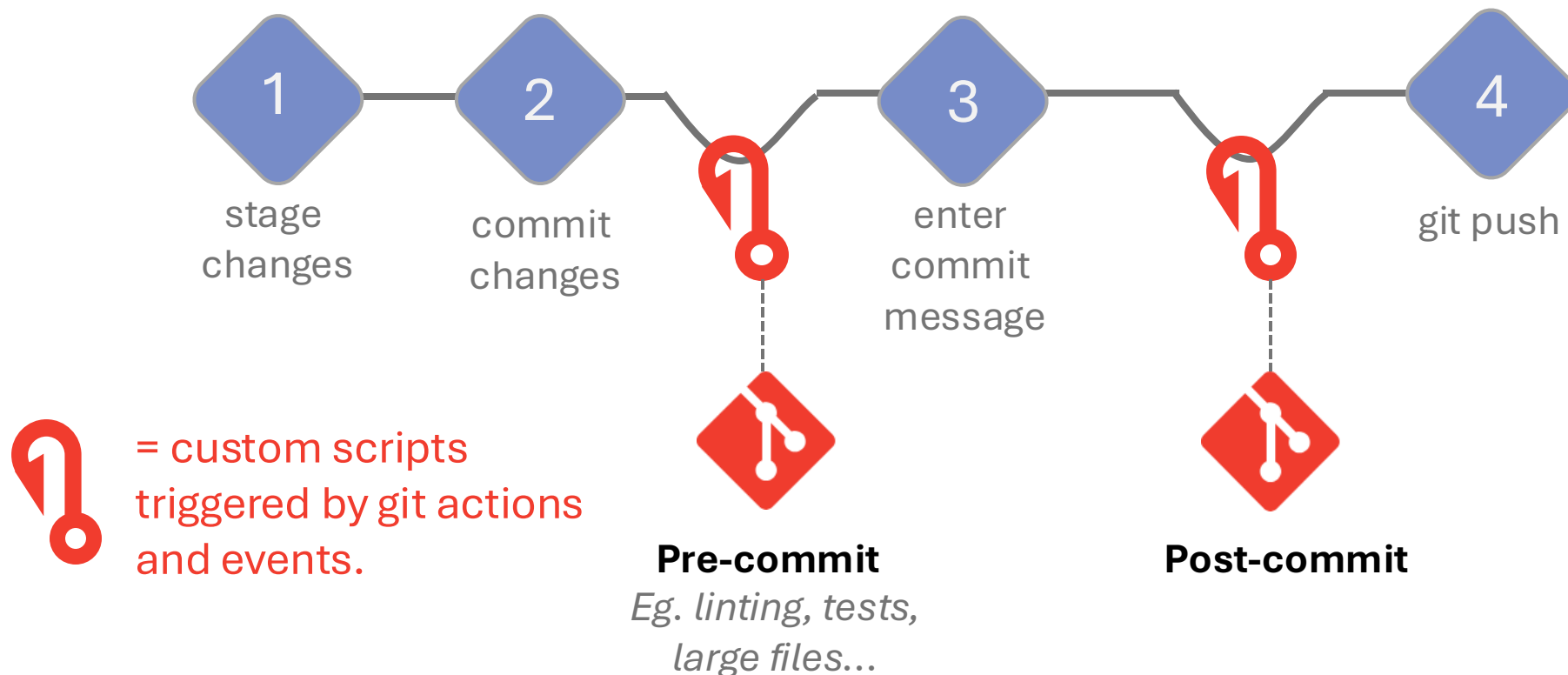 *"lintr provides static code analysis for R. It checks for adherence to a given style, identifying syntax errors and possible semantic issues, then reports them to you so you can take action."*

- **Python** - list of linters:

 *https://github.com/vintasoftware/python-linters-and-code-analysis*

# Using git hooks and pre-commits

Code author can use automated pre-commit checks before review



1 — stage changes

2 — commit changes

3 — enter commit message

4 — git push

**Pre-commit**
*Eg. linting, tests, large files...*

**Post-commit**

= custom scripts triggered by git actions and events.

https://git-scm.com/book/ms/v2/Customizing-Git-Git-Hooks
https://pre-commit.com/

# Summary of checklist

- Variable names
- Hard coded values/magic numbers
- Duplicated code
- Complex if else statements
- Long functions
- Obscure lines
- File paths
- File names

- Unintended behavior
- Comments
- Documented code
- Documented data
- Coding style (with automated checks using linter/git hooks)

# Other things for a reviewer to check?