# Sharing in the Open

NHSR conference 2022

Jonny Pearson,

Lead Data Scientist,

Digital Analytics and Research Team (DART)

NHS England

jonathanpearson@nhs.net

16th November 2022

# Current Materials

## The benefits of coding in the open

I was often asked what the value of coding in the open is to the teams themselves, those who are opening the code. There is a lot of value. I've shared that in various formats:

- Blog post on GDS blog: The benefits of coding in the open
- 30 min conference talk: Coding in the open in government
- 2 min video: Why we code in the open

## How to make your code open

- Official guidance on the Government Service Design Manual: Making source code open and re
- Blog post with examples: How to open up closed code
- GOV.UK guidelines for licensing GDS code
- GDS open source guidelines for open code that GDS explicitly intends to support

## GDS open code

- GDS's GitHub organisation: Alphagov
- GOV.UK have documented all of the GOV.UK code: Developer docs
- GOV.UK Verify code on GitHub
- Curated list of GOV.UK Frontend code and ecosystem
- All posts about open code on the GDS technology blog

## Coding in the open across government

- List of UK central government code on GitHub
- Blog post by Ministry of Justice: Why we code in the open
- Blog post by DWP: Doing the hard work to make things open
- GCHQ's most popular open source project: CyberChef

## Security when coding in the open

- Blog post on GDS technology blog: Don't be afraid to code in the open: here's how to do it securely
- Guidance: When code should be open or closed
- Guidance: Security considerations when coding in the open

From the security perspective, it's also worth knowing that while GCHQ don't code in the open, they have released quite a bit of open source code.

**If you take one thing from this talk then take this**
In 2018, Anna Shipman as Open Source Lead at GDS published these relevant materials:
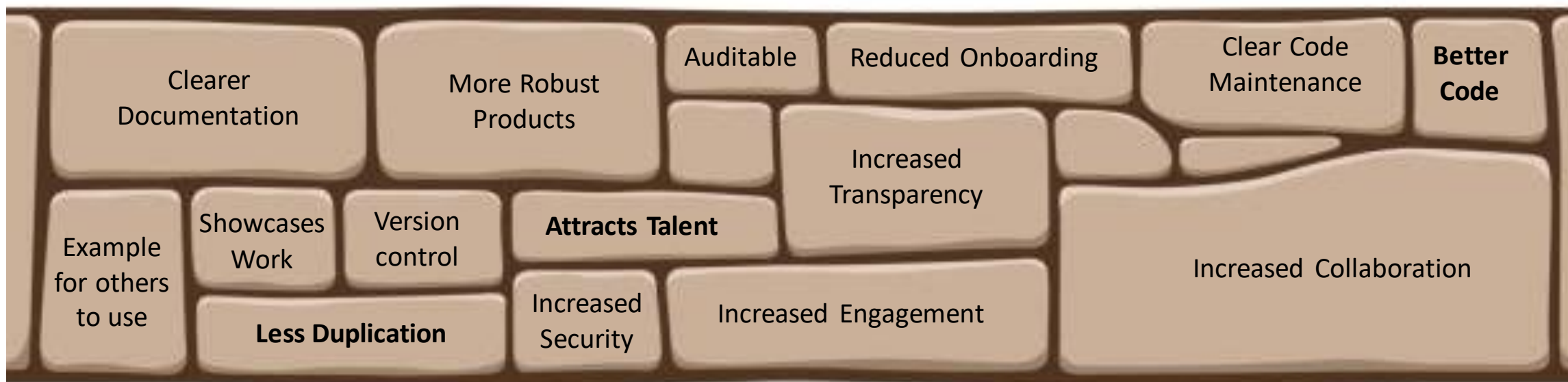https://www.annashipman.co.uk/jfdi/open-code-resources.html

# Why Should I (Mandate & Benefits)?

For a formal mandate then see the **Goldacre Review** & **12. Make new source code open - Service Manual - GOV.UK** copied into **NHS service standard - 12. Make new source code open**.
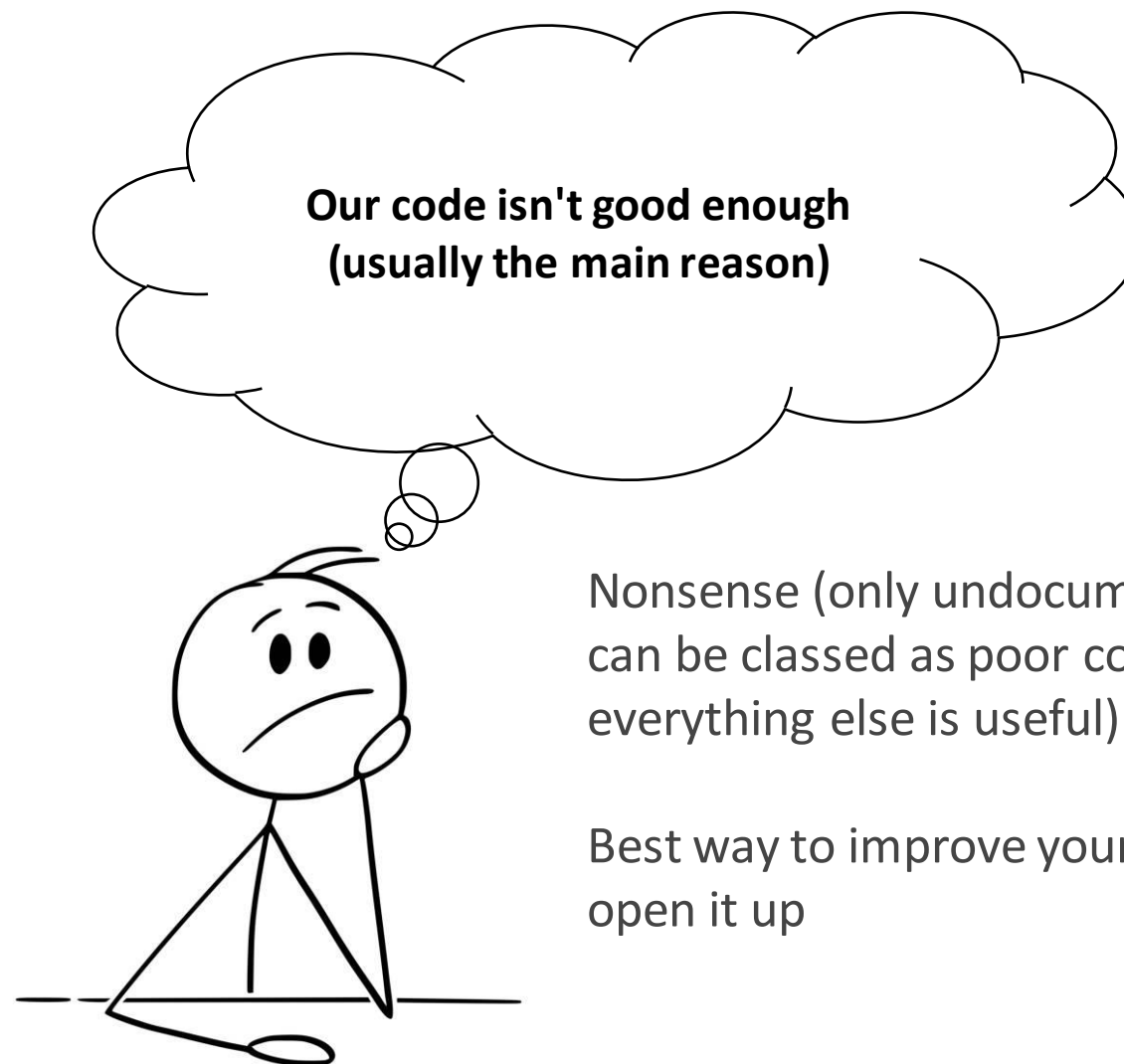Also see **Be open and use open source - GOV.UK** & **Open Data Charter - GOV.UK**

*"Public services are built with public money. So unless there's a good reason not to, the code they're based should be made available for other people to reuse and build on.*

*Open source code can be reused by developers working in government, avoiding duplication of work and reducing costs for government as a whole. And publishing source code under an open license means that you're less likely to get locked in to working with a single supplier."*

# What stops us?

*Based off Terence Eden talk on overcoming barriers to open code*

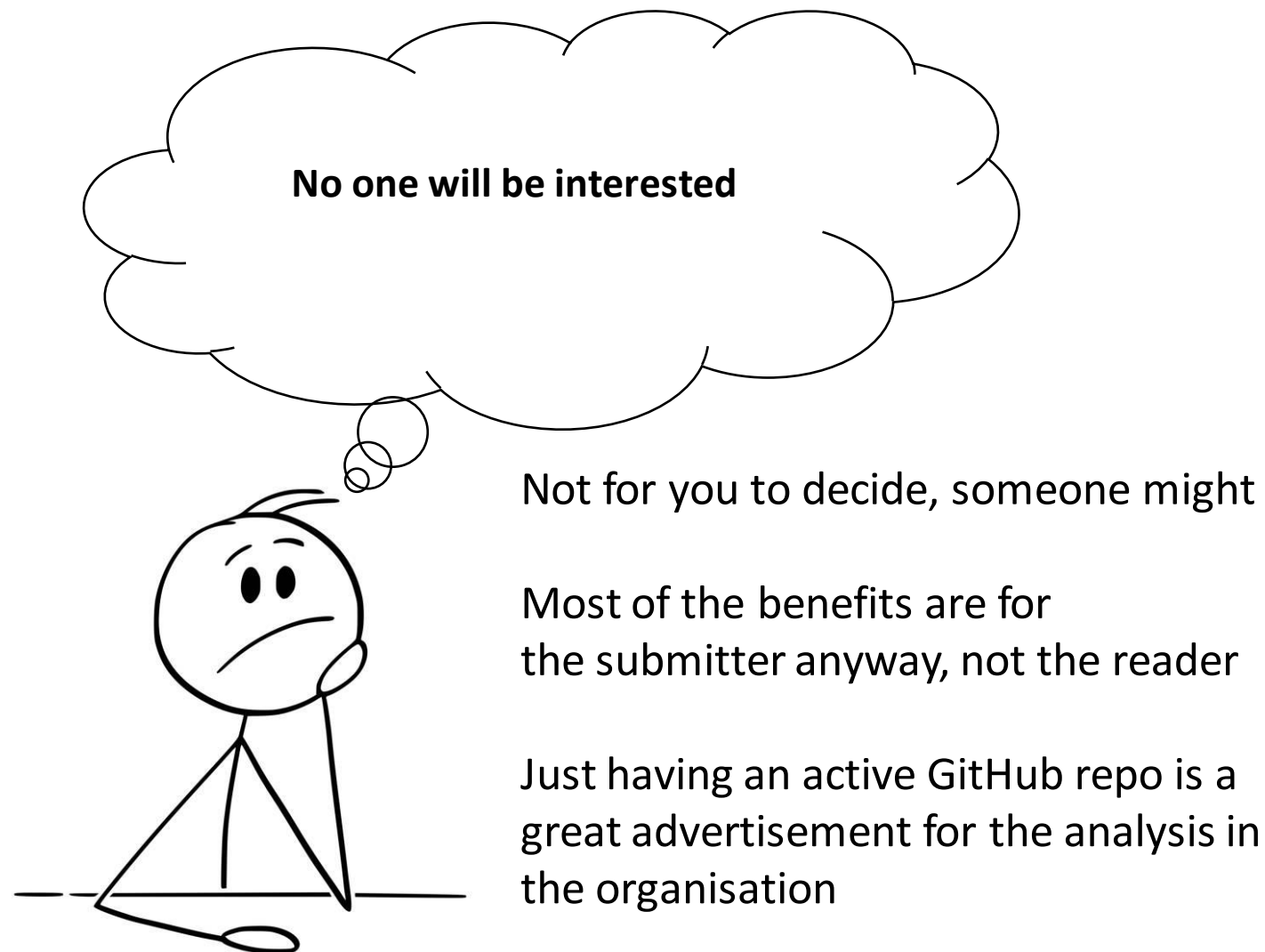**Our code isn't good enough
(usually the main reason)**

Nonsense (only undocumented chaos can be classed as poor code, everything else is useful)

Best way to improve your coding is to open it up

# What stops us?

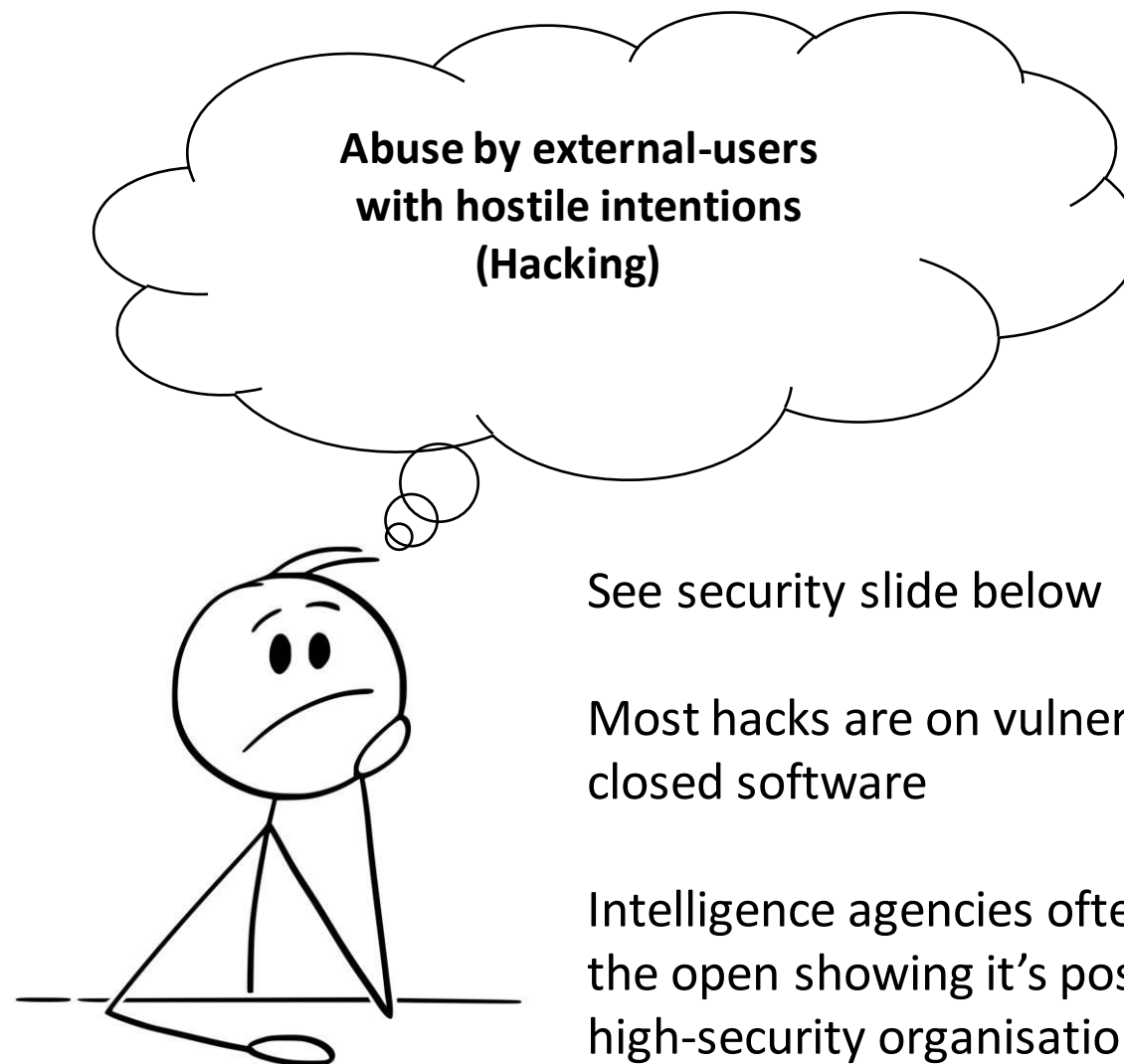*Based off Terence Eden talk on overcoming barriers to open code*

- ~~Our code isn't good enough (usually the main reason)~~

**No one will be interested**

Not for you to decide, someone might

Most of the benefits are for the submitter anyway, not the reader

Just having an active GitHub repo is a great advertisement for the analysis in the organisation

# What stops us?

*Based off Terence Eden talk on overcoming barriers to open code*

- ~~Our code isn't good enough (usually the main reason)~~
- ~~No one will be interested~~

**Abuse by external-users with hostile intentions (Hacking)**

See security slide below

Most hacks are on vulnerabilities of closed software

Intelligence agencies often work in the open showing it's possible for high-security organisations

# What stops us?

*Based off Terence Eden talk on overcoming barriers to open code*

- Our code isn't good enough (usually the main reason)
- No one will be interested

- Abuse by external-users with hostile intentions (Hacking)

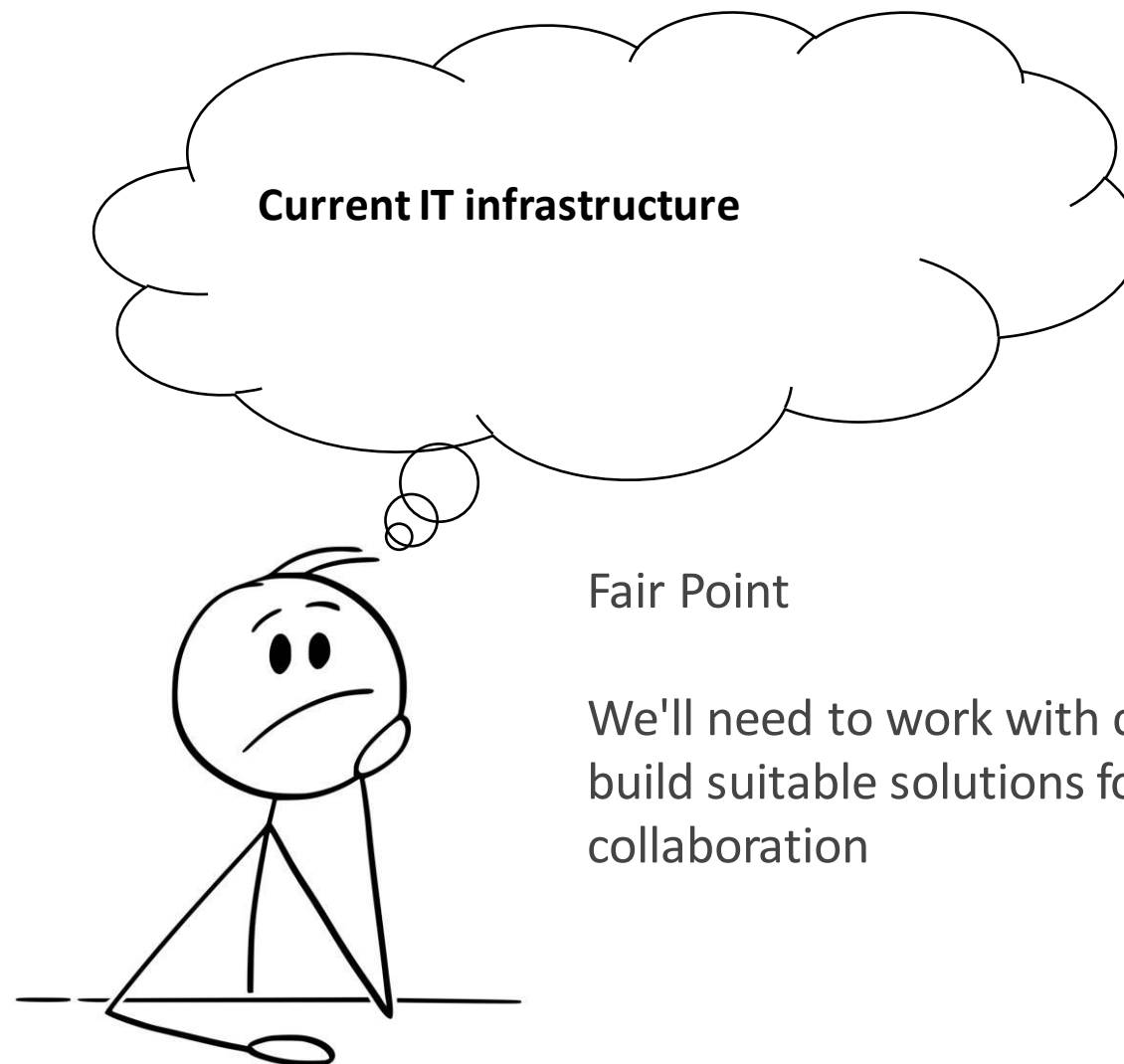**Abuse by external-users with good intentions**

Use of git branching strategy to isolate new features allows for appropriate review before merging.

Can control collaborators and clearly see who has done what and revert to earlier versions if required.

Clear licenses also needed for clarity about use and responsibility

# What stops us?

*Based off Terence Eden talk on overcoming barriers to open code*

- ~~Our code isn't good enough (usually the main reason)~~
- ~~No one will be interested~~

- ~~Abuse by external-users with hostile intentions (Hacking)~~
- ~~Abuse by external-users with good intentions~~

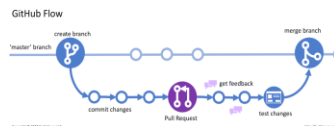**Accidental publication of sensitive information**

Need to assume a sensitive leak will occur at some point if coding in the open <u>or not</u> and so be ready to act!

Level of risk should be relative to likelihood and impact

# What stops us?

*Based off Terence Eden talk on overcoming barriers to open code*

- ~~Our code isn't good enough (usually the main reason)~~
- ~~No one will be interested~~

- ~~Abuse by external-users with hostile intentions (Hacking)~~
- ~~Abuse by external-users with good intentions~~

- ~~Accidental publication of sensitive information~~
- ~~~~

**It will take too long**

Benefits higher than cost

Gets faster the more you do it

# What stops us?

*Based off Terence Eden talk on overcoming barriers to open code*

- Our code isn't good enough (usually the main reason)
- No one will be interested

- Abuse by external users with hostile intentions (Hacking)
- Abuse by external users with good intentions

- Accidental publication of sensitive information
- It will take too long

**Current IT infrastructure**

Fair Point

We'll need to work with colleagues to build suitable solutions for open collaboration
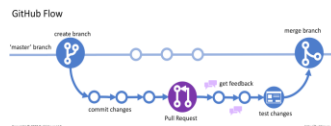
# What do I need to know / have?
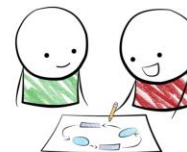


**Required**

| Understanding of Sensitive Data | Code with Version Control | Peer Review | Understanding of Licenses |

| Understanding of Security Considerations | Knowledge of ig escalation routes | Open Source Software and push/pull permissions |

**Useful**

| Awareness of What Others are Doing | Non-sensitive useful starter project | Open data for testing and demo | Community of Support |

# What do I need to know / have?

**NHS England**

**Required**

**Useful**

### Understanding of Sensitive Data

### Code with Version Control
GitHub Flow

### Peer Review

### Understanding of Licenses

### Understanding of Security Considerations

### Awareness of What Others are Doing

### Non-sensitive starter

More than just sensitive, secret, top secret information

Ideally no data stored alongside the code but if so, then written permission needs to be obtained by the data owner.

Consider checking the data content, code content, notebook outputs, commit messages and git history for:
- Credentials,
- Connection strings,
- SQL server addresses,
- Secret keys
- Unreleased Policy
- Business sensitive algorithms

Best to write config code separately.

If unsure about the git history then play it safe and create blank repo and copy across before making public

# What do I need to know / have?

**Required**

**Understanding of Sensitive Data**

**Code with Version Control**

GitHub Flow

**Peer Review**

**Understanding of Licenses**

**Understanding of Se...** **Considerations**

Even if you think you'll be the only person to ever use the code or that it will remain static, still use a versioning system

Establishing a branching strategy - recommend github flow

Consider templates (e.g. government cookie cutter)

Consider Semantic Versioning

Good commit notes - See Maintaining version control in coding - Service Manual

**Useful**

**Awareness of What Others are Doing**

...ity of Support

# What do I need to know / have?

**Understanding of Sensitive Data**

**Code with Version Control**

GitHub Flow

**Peer Review**

**Understanding of Licenses**

**Require**

Beyond quality assurance and accountability peer reviews increase knowledge sharing and produce better code.

Colleagues need to be ready to do informal and formal reviews of the code and approach used.

Recommend using the quality assurance of code for analysis and research checklist

A minimum level of testing includes "can a colleague run the code?". If the data is sensitive then simple fake data can be used to allow a smooth initial run

...edge of ig escalation routes

**Open Source Software and push/pull permissions**

Pull
Push
add commit
Local repo
Remote repo

**Useful**

**Awareness of What Others are Doing**

**Non-sensitive useful starter project**

START

**Open data for testing and demo**

**Community of Support**

# What do I need to know / have?

**Required**

**Useful**

U... S...

Review

**Understanding of Licenses**

**Open Source Software and push/pull permissions**

Pull

Push

add commit

Local repo

Remote repo

Awaren...

r testing and mo

**Community of Support**

---

MIT Licence

Allow unrestricted use but protects developer from liability and acknowledges contributions - Default licence for all new code

APLv2

If your code has regulatory requirements then an accompanying legal notice can be used here

GPLv3

If you want to prevent proprietary or closed re-use of code

Open Government 3.0 Licence

Recommended as default licence for all documentation

You should also include a copyright notice

# What do I need to know / have?

**Required**

**Understanding of Sensitive Data**

**Code with Version Control**

GitHub Flow

**Peer Review**

**Understanding of Licenses**

**Understanding of Security Considerations**

Ensure good development practice:

- Open the code early

- a set way of managing changes - GDS pull request guidance

Ensure you have checked any third party tools used for data transfers:

- Assess against the NCSC's cloud security principles

Ensure the libraries used are reputable:

- Do the developers behind the library have a track record?

- Is the development team adequately supported or just an interested individual?

- Is the library kept up-to-date or static?

**Useful**

**Awareness of What Others are Doing**

**Non-sensitive use starter project**

START

# What do I need to know / have?

**Required**

**Understanding of Sensitive Data**

Assume that any accidental leakage of data or sensitive content is a breach.

Thus, a route for reporting and dealing with this possibility needs to be clear from the offset. This should already be in place for any organisation dealing with sensitive data.

**Understanding of Licenses**

**Understanding of Security Considerations**

**Knowledge of ig escalation routes**

**Open Source Software and push/pull permissions**

Pull
Push
add commit
Local repo
Remote repo

**Useful**

**Awareness of What Others are Doing**

**Non-sensitive useful starter project**

START

**Open data for testing and demo**

**Community of Support**

# What do I need to know / have?

**Required**

| Understanding of Sensitive Data | Code with Version Control | Peer Review | Understanding of Licenses |
|---|---|---|---|

GitHub Flow

This gap in knowledge, combined with a culture of risk-aversion greatly hampers the process of adoption of new software in the NHS, leading many developers to look for workarounds.

Used a stepwise approach to engage properly.

Share both success and failures

**Open Source Software and push/pull permissions**

Pull

Push

add commit

Local repo          Remote repo

---

**Useful**

| Awareness of What Others are Doing | Non-sensitive useful starter project | Open data for testing and demo | Community of Support |
|---|---|---|---|

START

# What do I need to know / have?

**Required**

**Understanding of Sensitive Data**

Code

GitHub Flow

'master' branch

Copyright © 2018 Buildchent LLC

**Understanding of Secu... Considerations**

Balance between:

- Useful for the business

- Achievable

- Interesting

Suggestions:

1. Comparison of multiple standard models on open prescribing data

2. Novel visualization of GBD data

3. Geospatial mapping to highlight inequalities across multiple hierarchical boundaries

4. Sentiment analysis of survey data

**Useful**

**Awareness of What Others are Doing**

**Non-sensitive useful starter project**

START

**Open data for testing and demo**

**Community of Support**

# What do I need to know / have?

**Structured Activity and HER records**

- Opendata.nhsbsa
- PHE Fingertips
- NHS England Data Catalogue
- OpenPrescribing
- **Global Burden of Disease (GBD)**

Text

- **MIMIC III**
- **n2c2 PII tasks**
- **Diameter Health - GPT-2 generated notes**
- **Nottinghamshire Healthcare Foundation Trust (NHFT) - Friends & Family Test (FFT) Feedback Dataset**

Images

- **MIMIC CXR**
- **OASIS**
- **The Cancer Genome Atlas Lung Adenocarcinoma data**
- **NIHCC - DeepLesion**
- **CDAS**

**Requir**

**Usefu**

| ...ion Control | Peer Review | Understanding of Licenses |
|---|---|---|

| Knowledge of ig escalation routes | Open Source Software and push/pull permissions |
|---|---|

| ...e useful ...oject | Open data for testing and demo | Community of Support |
|---|---|---|

# Open Code Checklist

We use https://github.com/nhsengland/analyticsunit-template/blob/main/OPEN_CODE_CHECKLIST.md

# Using Git and Github with RStudio

Great materials on connecting git and github to RStudio can be found in happy git with r created by Jenny bryan (also see her other work on use this and tidyverse addins)

Helen Richardson (NHSD) is running a workshop on an Introduction to git and github on the 23rd November.

On NHS-R github there are git training materials also created by Helen.

# Example - https://github.com/nhsx/stm-survey-text

**NHS England**

---

≔ README.md ✎

## Structural Topic Modelling for NHS survey data

### NHSX Analytics Unit - PhD Data Science Internship Project

#### About the Project

`status experimental`

An exploration of methods and R libraries that can support information extraction from survey and free text responses.

---

## Roadmap

See the open issues for a list of proposed features (and known issues).

## Contributing

Contributions are what make the open source community such an amazing place to learn, inspire, and create. Any contributions you make are **greatly appreciated**.

1. Fork the Project
2. Create your Feature Branch ( `git checkout -b feature/AmazingFeature` )
3. Commit your Changes ( `git commit -m 'Add some AmazingFeature'` )
4. Push to the Branch ( `git push origin feature/AmazingFeature` )
5. Open a Pull Request

See *CONTRIBUTING.md* for detailed guidance.

## License

Distributed under the MIT License. See *LICENSE* for more information.

## Contact

To find out more about the Analytics Unit visit our project website or get in touch at analytics-unit@nhsx.nhs.uk.

---

### Project Structure

- The project code is found in the `R` folder of the repository (see Usage below for more information).
- The data used for in this analysis is found in the `data` folder of the repository.
- Exemplar outputs of this analysis is found in the `outputs` folder of the repository.
- The accompanying report is also available in the `reports` folder.

### Built With

`r v3.6.1`

- quanteda v3.0.0
- vader v0.2.1
- stm v1.3.6

## Installation

### `cran` installation

Launch the `stmnhsx.Rproj` file in a suitable IDE (e.g. RStudio).

The required packages are stored in `libraries.R`. Currently R files in R/main/ source this module and will

---

## Running the code

The folder `R/main/` contains the core code for the stm analysis and visualisation. The folder `R/experiments/` contains exploratory code used in additonal experiments in this project. To run the main code:

Update `main.R` User Inputs for the specific task and then run. Suggest running single sections at a time. The code starts off by loading the data before text-preprocessing of the text (removing punctuation and digits, stemming, tokenisation etc.), sentiment analysis and converting it to an stm data format. The best STM models are then determined by running a search over the number of topics. The outputs are then visualised in static and interactive ways. Finally, the last section of code allows an interactive term serach capability.

# Questions?

# Contact

See our open work [here](here)
Contact the team using [analytics-unit@nhsx.nhs.uk](analytics-unit@nhsx.nhs.uk)