






Sharing
Tools and
Artefacts for
Reproducible
Simulations
in healthcare

Protocol for assessing the computational reproducibility of simulation models on STARS

 Amy Heather¹,  Thomas Monks¹,  Alison Harper²,
 Navonil Mustafee², and  Andy Mayne³

¹University of Exeter Medical School, Exeter, UK

²University of Exeter Business School, Exeter, UK

³Taunton and Somerset NHS Foundation Trust, UK

Publication date: Add prior to submission

This protocol outlines how we plan to reuse available artefacts to reproduce results from published simulation studies. This forms part of the project STARS: "Sharing Tools and Artefacts for Reproducible Simulations in healthcare". It will be utilised to assess the computational reproducibility of published discrete-event (DES) simulation models in Python and R.

This protocol is archived as a pre-registration. Haroz 2022 identifies the Open Science Framework (OSF, <https://osf.io/>) and Zenodo (<https://zenodo.org/>) as suitable platforms for pre-registration.¹ In this case, Zenodo will be used as this is where other materials already exist for the STARS project, and so it can be stored alongside them in a Zenodo "community".

This work is supported by the Medical Research Council [grant number MR/Z503915/1].

Contents

| | | |
|----------|---|-----------|
| 1 | Summary diagram | 3 |
| 2 | Introduction | 4 |
| 3 | Logbook | 5 |
| 4 | Timing | 6 |
| 5 | Assessment of computational reproducibility | 7 |
| 5.1 | Set-up | 7 |
| 5.2 | Scope of reproduction | 9 |
| 5.3 | Familiarise with artefacts | 10 |
| 5.4 | Set up environment | 10 |
| 5.5 | Attempt to reproduce items in scope | 11 |
| 5.6 | Finishing up | 13 |
| 6 | Evaluation against guidelines | 14 |
| 6.1 | Best practice for sharing of research artefacts | 14 |
| 6.2 | Badges | 15 |
| 6.3 | Reporting guidelines | 15 |
| 7 | Summary report and research compendium | 16 |
| 7.1 | Computational reproducibility report | 16 |
| 7.2 | Research compendium | 16 |
| 7.3 | Archive on Zenodo | 18 |
| 7.4 | Inform the authors | 18 |
| A | Appendix: Badges | 19 |
| A.1 | "Open objects" badges | 20 |
| A.2 | "Object review" badges | 21 |
| A.3 | "Reproduced" badges | 22 |
| | References | 23 |

1 Summary diagram

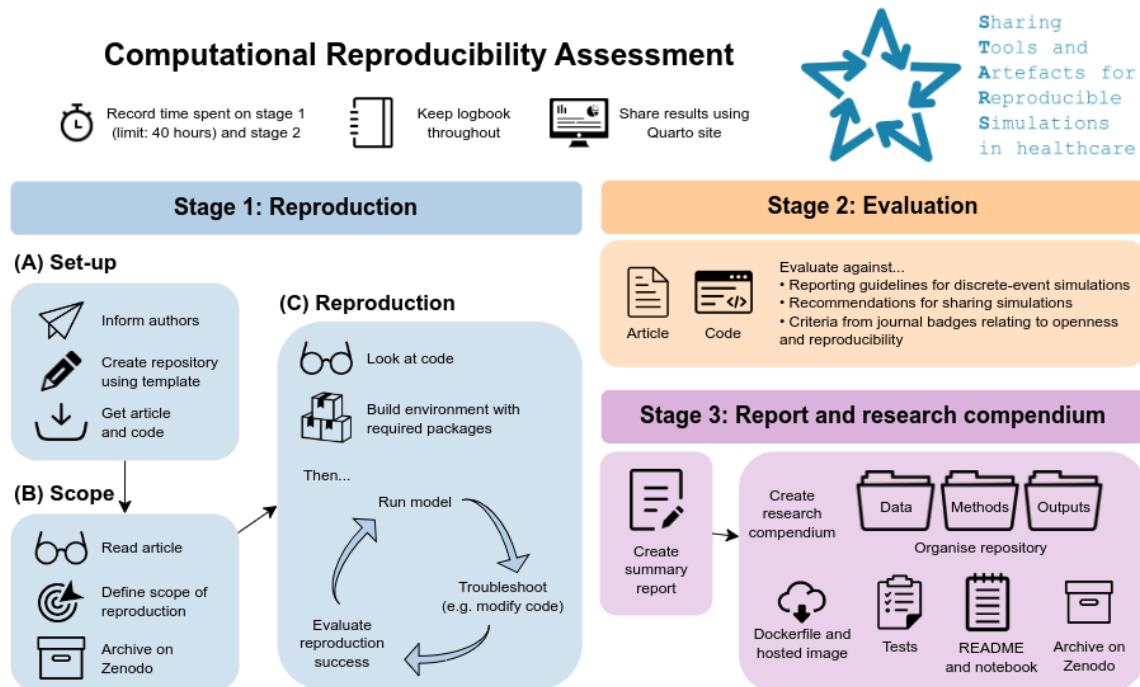


Figure 1: Workflow for assessing the computational reproducibility of discrete-event simulation models on STARS.

2 Introduction

In this protocol, we are focused on the **computational reproducibility** of simulation models. This is defined as the ability to get consistent results with a prior study when using the same data and methods as that study. We are focusing on models developed using **Python and R**, as these are popular free and open-source software (FOSS) for the development of models like discrete-event simulation (DES).²

This protocol will be used to assess the computational reproducibility of **six DES models**. These are selected from the models identified by **Monks and Harper 2023**.² The selection criteria are that: (i) the model code is publicly available, (ii) the model is created using Python or R, and (iii) the code has an open license (either already published, or added upon request from the STARS team).

Throughout the study, results will be openly available and shared via a **Quarto website**. This will compile information on the reproduction of the article. This includes the notebooks (.ipynb or .Rmd) producing the items in the scope, as well as a chronological log of work using Quarto blog posts, and then later, the reproduction report and detailed study results.

The protocol will often refer to our **template repository** which can be viewed here -https://github.com/pythonhealthdatascience/stars_reproduction_template. As you will see, one of the first steps for researchers will be to use this template to set-up their repository for the reproducibility assessment. **Change to link to Zenodo before publishing**

There are **three key stages** to this protocol which you should **work through in the following order**:

1. Assessment of computational reproducibility
2. Evaluation against guidelines
3. Summary report and research compendium

However, we first introduce two important processes that will need to be **completed during some or all of the stages**:

- Keeping a detailed record of work using a **logbook**
- **Timing** how long tasks take to complete

3 Logbook

Throughout **all of the stages** below, you should keep a **logbook**.

This is a detailed record of work recorded using **Quarto blog posts**, using the template provided. As suggested by Ayllón et al. 2021³ in their guidelines for keeping modelling notebooks, these posts will be **daily**, dated, chronological entries. **Tags** will be used to help indicate the activity on each day, and enable posts to be filtered by activity. Keeping a detailed log will support later understanding of what was done, and support preparing of final documents like the summary report.

Each entry in the logbook should contain the:

- **Researcher name** and **date**
- **Tags** (e.g. setup, scope, read, reproduce, guidelines, compendium, report)
- **Time** spent on tasks (if applicable)
- **Comprehensive** record of work. This should include record of working through each stage in the protocol, detailing **successes** and any **issues** faced, any **solutions** found to problems, and any **changes made to the model code** (noting where and how the code was changed). It may be relevant to include links to particular versions of a file or repository such as via the Git commit history.
- Clear statement if and when each item in the scope is considered to have been **successfully reproduced**.

4 Timing

During the first and second stages of your study (**assessment of computational reproducibility**, and **evaluation against guidelines**) you should **time how long each task takes**. Whilst timing, it is important you timestamp:

- When you finish reproducing each item within the scope
- When you have finished working on evaluation of artefacts (badges and recommendations on sharing), and when you have finished working on evaluation of the article (STRESS-DES and ISPOR-SDM)

These times should be recorded within the logbook alongside each activity (e.g. 12:10 to 12:45). The times should be monitored with a **maximum of forty hours** allowed for the first stage (attempting to reproduce the study), as in Krafczyk et al. 2021.⁴ This cut-off is implemented as we anticipate there would be little more to learn from spending longer than that time on reproducing a single study.

The only exceptions to this timing are:

- **Computation time.** This is at the researcher's discretion. For example, you should include short run times where you are still continuously working on the study. You should exclude longer run times where you are no longer working towards the assessment of computational reproducibility - for example, if you set the simulation to run for five minutes whilst going to make a cup of tea.
- **Time spent by other researchers.** You should record the time spent by the primary researcher on completing tasks, and time spent in discussion with other researchers. However, if other researchers spend time preparing for those discussions (e.g. reading the article, looking over the work), this does not need to be recorded.

5 Assessment of computational reproducibility

Remember! Record progress in your **logbook** and the **time** spent on each task, as described above.

5.1 Set-up

5.1.1 Inform the authors

You may have already completed this step (as we will have informed the authors about the study if we had to email them to request the addition of an open license to their repository). If not, you should:

1. **Email the corresponding author** about the study, including a **copy of the study protocol**.
2. If email rebounds, search online for a **recent email address** for any of the study authors.

If you do not hear back from the study authors, you do not need to follow-up on this email asking for a response, as the email simply serves to notify the study authors and does not require a reply.

A template for this email is provided:

Subject: *Reproduction of [Title of paper]*

Dear [Title (e.g. "Dr")] [Full name]

I am contacting you regarding your paper "[Title]" (doi: [DOI]).

*I am a [Postdoctoral research associate/research fellow/other position] based at [Institution], and I would like to **finish draft, or Tom to share***

Best wishes,

[Full Researcher Name including organisation]

5.1.2 Create repository using template

To set up the repository:

1. Go to https://github.com/pythonhealthdatascience/stars_reproduction_template and select the "**Use this template**" button.
2. Set-up the repository in 'pythonhealthdatascience' with the name 'stars-reproduce-surname-year' and description 'Assessing the computational reproducibility of surname et al. year as part of STARS.' Here, 'surname' and 'year' refer to the article being reproduced.

Update the template so that it is specific to the study being evaluated (e.g. titles, links, CITATION.cff). You should set up the provided **conda environment** for creation of the Quarto site.

5.1.3 Upload code to the repository

We require that the model code has been shared under an open license, and so, we should be free to upload this to our repository. You should upload any **code** and related artefacts (i.e. make a full copy of their code repository).

5.1.4 Upload journal article to the repository

Our next step will be to upload the article to the repository, but we first need to check which **license the article has been distributed under**. This is likely to be in a "copyright" section of the article. If it is distributed under a permissive license (e.g. CC-BY), then you should be free to upload the article and artefacts to the repository. However, if the article does not have a license enabling reuse in this manner (which is likely if it is restricted access/behind a paywall), then you will likely be unable to. Although you can request permission, this will often incur costs.

In these cases, you should search for whether there is a **green open access version** of the article. For example, copies of the article available on a pre-print site or on an institutional. In which case, these can be uploaded to the repository. If using a pre-print, you should check

If you cannot find a version of the article that can be shared, then you should not upload the article or images, but instead just provide a link to the article. However, otherwise, you should **upload all available materials** from the article to the 'original_study/' folder. This includes:

- The **journal article** and any **supplementary materials**
- Each **table and figure** from the main journal article as individual digital objects (e.g. .jpeg, .png) (not including those from the appendices/supplementary materials)

You should amend the provided **template page** ('quarto_site/study_publication.qmd') to display the PDFs or links to the journal article, and provide a link to the code, as well as providing links for where these uploaded materials were sourced from.

5.1.5 Update license for repository if required

Check the type of license used by the study authors. By default, the template includes an MIT license, but **you may need to change this if the authors used a different license**, so it is compatible.

Following the guidance of the Turing Way⁵ and R packages book,⁶ the license does not need to be modified according to the packages used, unless code from that package is embedded within the work (for example, copy+and+paste a function) or if it is distributed as a binary with the work (i.e. bundled and stored with the work, rather than setting it to be exported from somewhere like PyPI or CRAN). If the code simply imports and uses functions from packages, this is not assumed to be derivative work, and hence a permissive license can be chosen. Likewise, the license can be permissive when using Docker (as set up during the research compendium stage), as we do not distribute the Docker image itself, but instead distribute the Dockerfile or refer to someone else to distribute (such as the GitHub container registry).⁷

5.2 Scope of reproduction

5.2.1 Read the journal article

Read through the journal article (but not yet looking into the code or data). If you want to take notes, include these within the logbook. These can be within a collapsible callout to aid readability of the log.

5.2.2 Define scope of reproduction

The next step is to define the scope of the reproducibility study - in other words, what parts of the paper you intend to reproduce. This should be focused on the **results of the simulation** (rather than other results like description of the sample). To identify the scope you should:

1. Look through each of the **tables and figures** in the article (excluding the supplementary material) and identify whether they are within scope or not (i.e. do they present results of the simulation).
2. Identify any **'key results' in the text** of the article that are within scope. These are results that are highlighted within the **abstract or results** section. This may include items from the supplementary materials if referred to in the text.
3. Evaluate whether the identified 'key results' are **already covered by the tables and figures** from the article. If not, they should be included in the scope.
4. Make **consensus decision on scope with at least one other team members**.

You can include notes from thoughts and discussions around the potential scope within the **logbook**.

The criteria for what is considered part of the scope is adapted from Wood et al. 2018.^{8,9}

5.2.3 Compile items in scope

Once the scope has been decided, upload each item to the 'original_study/' folder.

- For **tables**, download a **CSV** version of each if available. Otherwise, convert the tables into CSV format.
- For **figures**, download the **highest-quality** version of each figure that is available. You should have already done this in an earlier step, but may still need to do this step if any figures are chosen from the supplementary materials as being in scope.
- For results described in the **text** (but not captured in a table or figure), record in a format appropriate to then later compare against (for example, within a **CSV**).

You should amend the provided **template page** ('evaluation/scope.qmd') to display each of the items.

If the article does not have permissions to enable this, then you can upload these materials to a separate **private repository** for your own reference, but cannot display these images or tables within the public repository/quarto site.

5.2.4 Archive scope on Zenodo

With the organisation linked to Zenodo, **toggle Zenodo to preserve** that repository, and then create a **release on GitHub**, which Zenodo will then automatically download and register with a DOI. The release should be recorded within the **changelog**.

This release should serve as a public registration of the intended scope of the reproducibility study, and archiving the repository at this point (prior to having started using or really looking at the the code). As stated in the "Guide for Accelerating Computational Reproducibility (ACRe) in the Social Sciences", it is important that the scope is defined at the start of the study, and publicly archived so as not to be amended during the course of the study.[10]

5.3 Familiarise with artefacts

5.3.1 Look over the code/data

Browse through any code and data that you uploaded to the repository. The aim of this step is to familiarise with the materials before setting up the environment or running the code. There are no requirements for how you should do this, but some suggestions include:

- In logbook, recording a one-sentence description of each file and tree of the uploaded files (as suggested by Ayllón et al. 2021³)
- Looking for sections of code that produce items within the scope and recording this within the logbook (as in Krafczyk et al. 2021[4]).

5.4 Set up environment

Identify the **software packages and their versions**, as used by the original study authors. This may be provided in an **environment file or within the article**. If not, the researcher should:

- Identify packages from those named to **import** within each of the scripts
- Select versions by looking at the **version history** for that package on a package repository (e.g. PyPI, CRAN), and identifying a version whose release date is **closest to but still prior to** the date of the code archive or paper publication (whichever is earliest).

Use or set up an environment file with the identified dependencies, and **create the environment**. Researchers should use simple methods for environment management such as Conda or VirtualEnv in Python, and renv in R. For simplicity, we are not requiring that you match the operating system used.

Important! This should be set up within the '**reproduction/**' folder, whilst the '**original_study/**' folder should remain untouched. This means, if an environment file is provided, you should copy it into the '**reproduction/**' folder.

5.5 Attempt to reproduce items in scope

This stage is an **iterative** process of running the code and attempting to reproduce the items in the scope. For this stage, the researcher should:

- Leave the 'original_study/' untouched - simply **copy over any relevant files** into the 'reproduction/' folder before running or modifying them.
- Use a **notebook** (.ipynb or .Rmd) when running the code, as this enables you to easily share the code and produced outputs from the scope, hence following a literate programming approach. This notebook can be made available to view within the Quarto site by setting it as part of the toctree in '_quarto.yml', as in the template.
- Continue attempting to reproduce each item until you feel it is **successfully reproduced** (as detailed below) - or you run out of time (from the maximum 40 hours allowed).
- **Troubleshoot issues**, contacting the study authors if necessary (as detailed below)

Remember! Continue taking detailed notes and timings in your logbook, including each success and issue, and making note of any changes made to the model code. Whilst keeping notes in the logbook, it is recommended that you **copy over in-progress outputs** to the blog post folder (e.g. .png file for a figure), so that you can easily and visually share progress in reproduction within the log.

5.5.1 Successful reproduction

For **each item in the scope**, the researcher should decide whether it has been **successfully reproduced**. A binary decision should be made for each item (and none should be labelled as 'partial success').

This is a **subjective** decision. A successful reproduction does **not require that exactly the same results** are found. An item can be considered successfully reproduced if **minimal variation** is observed from the original results.

As an example, if it is possible to produce a table with some numbers being a match or very similar, but some numbers being substantially different, then this would be classed as having **not** been successfully reproduced. If however all aspects of the item were reproduced with reasonable similarity, this can be classed as **successful** reproduction.

Further recommendations:

- **Figures** - these are compared **by eye**. Researchers should be unconcerned by **minor differences in presentation**, with regards to evaluating reproduction success.
- **Numbers** - researchers should calculate and report the **percentage difference** in results between the manuscript and the researcher. As reported by Wood et al. 2018,^{9,8} a meaningful difference in a value will vary between studies, and so it is difficult to set a single rule on what is or is not a minor difference. As such, researchers should follow a similar approach to Schwander et al. 2021,¹¹ considering whether the figure is reproducible at varying levels of percentage error (5%, 10% and 20%). However, they **should then use their judgement to decide whether the**

item has been reproduced. This is similar to one of the definitions proposed by McManus et al. 2019¹² - "*Results... vary only by XX% compared to the original, AND are consistent with the original conclusions*" - incorporating both numerical comparison and allowance for variability in whether this constitutes a meaningful difference from the original results.

Before concluding the reproduction, this decision should be run by at least one other researcher on the project, to ensure there is a **consensus decision** on whether the items were successfully reproduced. This conversation and the decision made should be including in the timing and recorded in the logbook.

In assessing reproduction success, it is important to note (as in Laurinavichyute et al. 2022¹³ and Wood et al. 2018⁹) that the focus is **not** on the quality or robustness of the original results, or whether the main claims of the study are consistent. Instead, the focus is on whether it was possible to reproduce the article's results **within a reasonable margin of error** (given that we do expect a little variation, since discrete-event simulations are stochastic models, and may not have been fully controlled using random seeds or with any environment differences).

It is important to **clearly timestamp** once the decision of "successful reproduction" has been made for each of the items in the scope.

5.5.2 Troubleshooting

Researchers should **troubleshoot** any issues encountered (including **making changes to the provided code**). In allowing modification and writing of code, our intention is that researchers try **as much as possible** to attempt to reproduce from the scope. The allowance of writing new code is similar to the approach of Krafczyk et al. 2021⁴ and the ACRE project.¹⁰ Examples of changes you may need to make include:

- Correcting paths to files
- Correcting the versions of software, or adding missing packages or libraries
- Fixing errors in the code
- Adding code to produce an item in the scope, if not otherwise provided
- Adding a method for **controlling randomness** in the simulation (if not otherwise set up), so you can get consistent results with yourself between re-runs of your notebook
- Adding a warm-up period (if suspected but not included in the code)

Troubleshooting can include asking **advice from other members of the STARS team**. In these cases, the researcher should ensure that they include a record the time spent, everything discussed, and any recommendations made.

5.5.3 Contacting the authors

Despite troubleshooting, you may remain unable to run the code, or have large discrepancies with the original paper. In this case, once troubleshooting is exhausted, you should **contact the original**

author. This email should:

- **Recap** the project (since our last email, when we informed them about the study).
- Link to the **Quarto site** with the documented reproduction attempt and list of issues that require resolution. Make sure the description of problem is specific (e.g. identifying line in paper and place in code where think something is missing, or where an issue is occurring).
- **Ask for suggestions** on an alternative course of action for issues, or for the complete code/data if missing.

If there is no response in two weeks, the researcher should contact them again. If there is still no response two weeks later, this can be marked as non-response. When emailing authors, it is suggested to follow the guidance on language and adapt from the **email templates** provided by ACR in the chapter "Guidance for Constructive Communication Between Reproducers and Original Authors" from their guide.¹⁰ The allowance of contacting authors is similar to the approaches of several studies,^{4,9,10,14,15} with a maximum of four weeks for responses as in Konkol et al. 2018.¹⁵ This approach does however differ from Laurinavichyute et al. 2022¹³ who did not contact authors, since they considered reproducibility to be only about the available data and procedures and not anything shared privately.¹³

5.5.4 Running out of time

If forty hours has passed, you should stop working on the reproduction. However, there are two exceptions:

1. If the model was not set up - and you had not yet implemented - control for randomness (enabling you to reproduce your own results), then you should choose an alternative method for getting stable simulation results. This can be by doing a very large number of replications, or by assessing the required number of replications for a stable simulation.
2. Decide on the reproduction success of each item.

5.6 Finishing up

5.6.1 Tidy up notebook and create reproduction success page

Tidy the reproduction notebook, so it simply produces each of the items in the scope, and clearly state how each section relates to the original article (e.g. captioning 'Reproduction attempt for Figure 2').

Using the **template page** ('evaluation/reproduction_success.qmd'), show each item from the scope (as in the original article) alongside our best reproduction attempts. Include the decision on the reproduction success for each item (along with any justification for this decision).

6 Evaluation against guidelines

This section is completed **after** the attempted reproduction (so as to not interfere with timings).

Remember! Record progress in your **logbook** and **time spent** on each task.

Important: This evaluation is based on the **original** journal article or repository from the author (as in original_study/), and not on the repository you made whilst reproducing this study (reproduction/). If the original study had multiple repositories to choose from (e.g. development and archived code, both prior to publication date), remember to **refer to both of them** if there are any differences between them.

Getting a second opinion: If you are uncertain on any criteria, you should note these in the **logbook**. Any criteria that were **unmet or uncertain** should then be discussed with at least one other researcher on the project to get a second opinion. Record the discussion (and its timing) in the logbook, and explain and justify the choices for uncertain items.

6.1 Best practice for sharing of research artefacts

The artefacts (repository) associated with the original study will be evaluated against two sets of criteria/ recommendations on the sharing of research artefacts for simulation models:

- Criteria from the **best practice audit** conducted by Monks and Harper 2023² as part of their review. The audit is described in detail in the repository associated with their review.¹⁶ The items used in this audit were based on guidance from the Turing Way,⁵ Taylor et al. 2017,¹⁷ and the Open Modelling Foundation (OMF) minimal and ideal reusability standards,¹⁸ focussing on items that were relevant to the modelling and simulation community.
- The **STARS framework** (developed following the above review) recommends essential and optional components when sharing healthcare simulation studies.

You should use the **provided template** ('evaluation/artefacts.qmd') to assess whether the artefacts from the original study meet the criteria/recommendations from each of these sources. Each criteria are evaluated as being "fully", "partially" or "not met".

For the best practice audit, all studies will have previously been evaluated by Monks and Harper 2023². After you have done your assessment, you should add the results from their review of that study to the **logbook** and **compare** against your assessment, as this can provide a simple sense-check/second opinion.

6.2 Badges

Several organisations and journals have developed badges which can be displayed alongside a research article to indicate how open and potentially reproducible it is, as detailed in Appendix A. We will evaluate the original study artefacts (repository) against **badges that relate to code** (and not those specific to data), due to the nature of DES models (where "data" is often just parameters as part of the model script, with perhaps a few additional parameters in a separate data file within the repository). These badges are:

- "Open objects" badges: NISO "Open Research Objects (ORO)", ACM "Artifacts Available", COS "Open Materials" and "Open Code", and IEEE "Datasets Available"
- "Object review" badges: ACM "Artifacts Evaluated" (rated as "Functional or "Reusable") and IEEE "Code Reviewed"
- "Reproduced" badges: NISO "Results Reproduced (ROR-R)", ACM "Results Reproduced", IEEE "Code Reproducible" and Psychological Science "Computational Reproducibility"

You should use the **provided template** (evaluation/badges.qmd) to assess whether the artefacts from the original study meet the criteria for each of these badges. A **binary** decision is made for each criteria (as being either met or not met).

6.3 Reporting guidelines

The **journal article** will be evaluated against two reporting guidelines for discrete-event simulation studies:

- STRESS-DES: Strengthening The Reporting of Empirical Simulation Studies (Discrete-Event Simulation)¹⁹
- The generic reporting checklist for healthcare-related discrete event simulation studies derived from the the International Society for Pharmacoeconomics and Outcomes Research Society for Medical Decision Making (ISPOR-SDM) Modeling Good Research Practices Task Force reports.²⁰

You should use the **provided template** ('evaluation/reporting.qmd') to assess whether the criteria from these guidelines are met by the journal article (including the supplementary material, although not including the code unless the article specifically refers to it for providing particular information). Each criteria are evaluated as being "fully", "partially" or "not met", with detailed evidence provided to support these claims (such as quotations from the article). If a criteria is not met by the original study, you are welcome to make a suggestion in the evidence column of what you think the likely answer for that criteria might be.

7 Summary report and research compendium

7.1 Computational reproducibility report

Use the **provided template** ('evaluation/reproduction_report.qmd') to produce a simple summary report for the reproducibility assessment. This template should guide you to include:

- Short study description
- Citation to original study
- Number and percentage of items reproduced from scope and time elapsed
- Required troubleshooting steps
- Display of reproduced items alongside original study items
- Percent stacked bar chart displaying the proportion of criteria met for each of the evaluations against guidelines

7.2 Research compendium

Once the computational reproduction has been completed, the repository will be restructured into a "**research compendium**". This a term first introduced by Gentleman and Lang 2007²¹ which they define as "*both a container for the different elements that make up the document and its computations (i.e. text, code, data, . . .), and as a means for distributing, managing and updating the collection.*"²¹ Marwick et al. 2018 define a research compendium as having three key components:

1. Files organised according to convention
2. Seperate data, methods and outputs
3. Specifying the environment used for the analysis.²²

A research compendium might also be referred to as a "**reproducibility file bundle**",²³ or as a "**reproduction package**".⁴ Although not required to be structured as a package, this can be helpful in providing a structure for dependency management and file organisation, and for continuous integration of automated code testing.²² The repository will be modified as per the proposed structure below. Some of these modifications may already be in place from having gone through the reproduction steps above.

Our primary motivation in doing this step is to make it **easy and clear for someone to re-run** our reproduction, whilst making relatively **minimal changes** to the code itself. Hence, we are not necessarily amending the repository to meet all of the recommendations for best practice of sharing DES models.

7.2.1 Modify repository

Make the following changes to the **reproduction/** folder, if not already implemented:

- **Have separate folders for data, methods and outputs** - as recommended by Marwick et al. 2018.²² The exception for this change is parameters coded into the scripts (since these would require a large amount of work and restructuring to separate from the scripts, contrary to our motivation in this stage).
- Create **tests** which check whether a user is able to get the same reproduce results as we obtained during the reproduction, based on comparison of csv files.
- Create a **Dockerfile** and double-check it works (build image and run model notebook/s).
- Enable the GitHub action to publish the Docker image on the **GitHub container registry**.
- Make sure that **model notebook/s** contain:
 - Run time recorded within the model notebook.
 - Clearly states which parts of the notebook produce each item from the reproduction scope.
- Ensure that the **README** contains:
 - Original study citation.
 - Simple summary of the model (potentially incorporating any diagrams of the model that were provided).
 - Scope of the reproduction (including images of the figures/tables from the original study).
 - Overview of the repository.
 - Instructions for setting up the environment.
 - Instructions for running the model (and reproducing items from the scope).
 - Instructions for running the test/s.
 - Hardware and software specs for the computer used for reproduction (machine, RAM, operating system and version).
 - Run time for the model.
 - Instructions for citation.
 - Short description of license.
- Ensure **Quarto site** displays the reproduction README and notebook/s.

7.2.2 Test-run with second team member

Once the research compendium is complete, a **second researcher** on the team should attempt to use it and confirm if they were able **reproduce** the results of the first researcher, and to check the compendium for **clarity**. This is similar to the approach of Krafczyk et al. 2021.⁴ It should be recorded within the **logbook**.

7.3 Archive on Zenodo

Once modification of the repository is completed, a new GitHub release should be created to archive the repository on Zenodo (with record of this in the changelog).

7.4 Inform the authors

Email the authors again to:

- Let them know we have finished the assessment
- Include link to GitHub and Zenodo.
- Let them know how we are going to use the results from this work (i.e. lessons from reproduction, and guide framework design, and that this was not about validity of results).

A Appendix: Badges

A badge is a label or image that is displayed alongside a published research article. There are currently several different badges from various organisations that can be used to indicate how open and potentially reproducible an article and its artefacts are. We have identified journal badges provided by the following organisations/journals:

- National Information Standards Organisation (NISO)
- Association for Computing Machinery (ACM)
- Institute of Electrical and Electronics Engineers (IEEE)
- Center for Open Science (COS)
- Springer Nature
- Psychological Science

Badges related to reproducibility are typically across three categories - "open objects", "object review" and "reproduced". However, there are other badges available that are related to reproducibility. These include badges for pre-registration, like the COS "Preregistered" badge.²⁴ There are also badges for replication, which is when an independent study on the same question find consistent results (potentially with new artefacts and methods). Examples of this are the NISO "Results Replicated (RER)" badge²⁵ and IEEE "Code Replicated" and "Dataset Replicated" badges.²⁶

In some journals, these criteria are set as requirements for publication with the journal, rather than as badges. An example of this is the Psychological Science journal, which recently transitioned from awarding COS badges to making them requirements.²⁷

We are focused on badges awarded by journals, but there are examples of badges that can be added to a repository if authors have followed a particular framework, either by self-allocating the badge or going through a review process. An example of this is Van Lissa et al. 2021²⁸ who provide a package to facilitate use of a reproducible workflow in R, and suggest that a badge can be added to the README.md file of that repository if the package is used.²⁸

A.1 "Open objects" badges

"Open objects" badges relate to permanently archiving digital objects in a public repository with a persistent identifier and open license.²⁵ Examples include:

- NISO "Open Research Objects (ORO)" and "ORO-A" (if all relevant objects available)²⁵
- ACM "Artifacts Available"²⁹
- COS "Open Data", "Open Materials" and "Open Code"²⁴
- IEEE "Code Available" and "Datasets Available"²⁶
- Springer Nature "Badge for Open Data"³⁰

The table below compares the criteria for each of these badges (based on the sources cited above). As our focus is on DES models, we have excluded the badges that just relate to data, since simulation studies typically do not have "data", but instead have parameters (which often form part of the code itself). We have likewise excluded the COS "materials" badge since this is described as sharing the "components of the research methodology" - which in our case, would typically just be the article and the code.

The remaining badges either specifically relate to code, or are more broad: the NISO badge is relevant to "author-created digital objects used in the research (including data and code)",²⁵ and the ACM badge is relevant to "artifacts associated with the research".²⁹

Table 1: "Open objects" badge criteria

| Criteria | NISO | ACM | COS (code) | IEEE (code) |
|--|------|-----|------------|-------------|
| Stored in permanent archive that is publicly and openly accessible | ✓ | ✓ | ✓ | X |
| Has a persistent identifier | ✓ | ✓ | ✓ | X |
| Includes an open license | ✓ | X | ✓ | X |
| Complete set of materials are shared (as would be needed to fully reproduce article) | ✓* | X | ✓ | ✓ |
| Artefacts are sufficiently documented, for researcher to understand how the code is used and relates to the reported methodology (e.g. package versions) | X | X | ✓ | X |

* The standard "ORO" badge does not require this - but if all relevant research objects are available, the badge is modified to "ORO-A".

Abbreviations: ACM, Association for Computing Machinery; COS, Center for Open Science; IEEE, Institute of Electrical and Electronics Engineers; NISO, National Information Standards Organisation; ORO, Open Research Objects.

A.2 "Object review" badges

"Object review" badges relate to the digital objects (i.e. data, code) being reviewed according to the criteria of the badge issuer.²⁵ Examples include:

- NISO "Research Objects Reviewed (ROR)"²⁵
- ACM "Artifacts Evaluated - Functional" and "Artifacts Evaluated - Reusable"²⁹
- IEEE "Code Reviewed" and "Datasets Reviewed"²⁶

Their criteria are summarised and compared below. The NISO badge is not included as it does not have criteria, but just states that badges in this category would evaluate against a specified set of criteria. The IEEE datasets badge is also excluded (as above for the "open objects" badges).

Table 2: "Open review" badge criteria

| Criteria | ACM (functional) | ACM (reusable) | IEEE (code) |
|---|---------------------|-------------------|-------------|
| Artefacts are sufficiently documented, to enable them to be run | ✓ | ✓ | X |
| Artefacts are very carefully documented, to the extent that reuse is facilitated | X | ✓ | X |
| Artefacts are relevant to and contribute to the article's results | ✓ | ✓ | X |
| Complete set of materials are shared (as would be needed to fully reproduce article) | ✓ | ✓ | ✓ |
| Scripts can be successfully executed | ✓ | ✓ | ✓ |
| Artefacts are well structured, to the extent that reuse is facilitated, adhering to norms and standards of research community | X | ✓ | X |

Abbreviations: ACM, Association for Computing Machinery; IEEE, Institute of Electrical and Electronics Engineers;

A.3 "Reproduced" badges

"Reproduced" badges are awarded when an independent party regenerates the article results using author objects.²⁵ Examples include:

- NISO "Results Reproduced (ROR-R)"²⁵
- ACM "Results Reproduced"²⁹
- IEEE "Code Reproducible" and "Dataset Reproducible"²⁶
- Psychological Science "Computational Reproducibility"^{27,31}

The criteria are summarised in the table below. It should be noted that ACM specify that results are reproduced "in part" using artefacts from the author, and that exact reproduction is not required but that results should be within an acceptable range for experiments of that type.²⁹ Whether deviation in results or any modification of the author code (such as minor troubleshooting) is permissible is not detailed within the viewed criteria for the other badges.

Table 3: "Reproduced" badge criteria

| Criteria | NISO | ACM | IEEE (code) | Psychological Science |
|--|------|-----|-------------|-----------------------|
| Independent party regenerated results using the authors' research artefacts | ✓ | ✓ | ✓ | ✓ |
| Reproduced within approximately one hour (excluding compute time) | X | X | X | ✓ |
| Artefacts are well-organised | X | X | X | ✓ |
| Artefacts are clearly documented and accompanied by a README file with step-by-step instructions on how to reproduce results in the manuscript | X | X | X | ✓ |

Abbreviations: ACM, Association for Computing Machinery; IEEE, Institute of Electrical and Electronics Engineers; NISO, National Information Standards Organisation.

References

- [1] Steve Haroz. *Comparison of Preregistration Platforms*. Feb. 2022. DOI: <https://doi.org/10.31222/osf.io/zry2u>. URL: <https://osf.io/preprints/metaarxiv/zry2u> (visited on 05/10/2024).
- [2] Thomas Monks and Alison Harper. "Computer model and code sharing practices in healthcare discrete-event simulation: a systematic scoping review". In: *Journal of Simulation* 0.0 (2023). Publisher: Taylor & Francis .eprint: <https://doi.org/10.1080/17477778.2023.2260772>, pp. 1–16. ISSN: 1747-7778. DOI: 10.1080/17477778.2023.2260772. URL: <https://doi.org/10.1080/17477778.2023.2260772> (visited on 05/10/2024).
- [3] Daniel Ayllón et al. "Keeping modelling notebooks with TRACE: Good for you and good for environmental research and management support". In: *Environmental Modelling & Software* 136 (Feb. 2021), p. 104932. ISSN: 1364-8152. DOI: 10.1016/j.envsoft.2020.104932. URL: <https://www.sciencedirect.com/science/article/pii/S1364815220309890> (visited on 05/13/2024).
- [4] M. S. Krafczyk et al. "Learning from reproducing computational results: introducing three principles and the Reproduction Package". In: *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 379.2197 (Mar. 2021). Publisher: Royal Society, p. 20200069. DOI: 10.1098/rsta.2020.0069. URL: <https://royalsocietypublishing.org/doi/10.1098/rsta.2020.0069> (visited on 05/10/2024).
- [5] The Turing Way Community. *The Turing Way: A handbook for reproducible, ethical and collaborative research (1.0.2)*. 2022. URL: <https://doi.org/10.5281/zenodo.7625728> (visited on 05/15/2024).
- [6] Hadley Wickham and Jennifer Bryan. *12 Licensing*. Apr. 2023. URL: <https://r-pkgs.org/license.html> (visited on 05/21/2024).
- [7] The Linux Foundation. *Docker containers: What are the open source licensing considerations?* URL: <https://www.linuxfoundation.org/resources/publications/docker-containers-what-are-the-open-source-licensing-considerations> (visited on 06/06/2024).
- [8] Benjamin Wood et al. "Replication Protocol for Push Button Replication (PBR)". en-us. In: *OSF* (Jan. 2018). Publisher: Open Science Framework. DOI: <https://doi.org/10.17605/OSF.IO/YFBR8>. URL: <https://osf.io/yfbr8/> (visited on 05/10/2024).
- [9] Benjamin D. K. Wood, Rui Müller, and Annette N. Brown. "Push button replication: Is impact evaluation evidence for international development verifiable?" en. In: *PLOS ONE* 13.12 (Dec. 2018). Publisher: Public Library of Science, e0209416. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0209416. URL: <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0209416> (visited on 05/10/2024).
- [10] Berkeley Initiative for Transparency in the Social Sciences. *Guide for Advancing Computational Reproducibility in the Social Sciences*. Sept. 2022. URL: <https://bitss.github.io/ACRE/> (visited on 05/15/2024).
- [11] Björn Schwander et al. "Replication of Published Health Economic Obesity Models: Assessment of Facilitators, Hurdles and Reproduction Success". In: *Pharmacoeconomics* 39.4 (2021), pp. 433–446. ISSN: 1170-7690. DOI: 10.1007/s40273-021-01008-7. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8009773/> (visited on 05/15/2024).

- [12] Emma McManus, David Turner, and Tracey Sach. "Can You Repeat That? Exploring the Definition of a Successful Model Replication in Health Economics". en. In: *PharmacoEconomics* 37.11 (Nov. 2019), pp. 1371–1381. ISSN: 1179-2027. DOI: 10.1007/s40273-019-00836-y. URL: <https://doi.org/10.1007/s40273-019-00836-y> (visited on 05/15/2024).
- [13] Anna Laurinavichyute, Himanshu Yadav, and Shravan Vasishth. "Share the code, not just the data: A case study of the reproducibility of articles published in the Journal of Memory and Language under the open data policy". In: *Journal of Memory and Language* 125 (Aug. 2022), p. 104332. ISSN: 0749-596X. DOI: 10.1016/j.jml.2022.104332. URL: <https://www.sciencedirect.com/science/article/pii/S0749596X22000195> (visited on 05/15/2024).
- [14] Tom E. Hardwicke et al. "Analytic reproducibility in articles receiving open data badges at the journal Psychological Science: an observational study". In: *Royal Society Open Science* 8.1 (Jan. 2021). Publisher: Royal Society, p. 201494. DOI: 10.1098/rsos.201494. URL: <https://royalsocietypublishing.org/doi/10.1098/rsos.201494> (visited on 05/15/2024).
- [15] Markus Konkol, Christian Kray, and Max Pfeiffer. "Computational reproducibility in geoscientific papers: Insights from a series of studies with geoscientists and a reproduction study". In: *International Journal of Geographical Information Science* 33.2 (Feb. 2019). Publisher: Taylor & Francis .eprint: <https://doi.org/10.1080/13658816.2018.1508687>, pp. 408–429. ISSN: 1365-8816. DOI: 10.1080/13658816.2018.1508687. URL: <https://doi.org/10.1080/13658816.2018.1508687> (visited on 05/15/2024).
- [16] Thomas Monks and Alison Harper. "Supplementary Materials: Computer model and code sharing practices in healthcare discrete-event simulation: a systematic scoping review. v1.2.0." In: *Zenodo* (June 2024). DOI: 10.5281/zenodo.11490636. URL: <https://zenodo.org/records/11490636> (visited on 06/06/2024).
- [17] Simon J. E. Taylor et al. "Open science: Approaches and benefits for modeling & simulation". In: *2017 Winter Simulation Conference (WSC)*. ISSN: 1558-4305. Dec. 2017, pp. 535–549. DOI: 10.1109/WSC.2017.8247813. URL: <https://ieeexplore.ieee.org/document/8247813> (visited on 06/04/2024).
- [18] The Open Modeling Foundation (OMF). *Reusability Standards*. OMF. May 2024. URL: <https://www.openmodelingfoundation.org/standards/reusability/> (visited on 06/04/2024).
- [19] Thomas Monks et al. "Strengthening the reporting of empirical simulation studies: Introducing the STRESS guidelines". In: *Journal of Simulation* 13.1 (Jan. 2019). Publisher: Taylor & Francis .eprint: <https://doi.org/10.1080/17477778.2018.1442155>, pp. 55–67. ISSN: 1747-7778. DOI: 10.1080/17477778.2018.1442155. URL: <https://doi.org/10.1080/17477778.2018.1442155> (visited on 05/21/2024).
- [20] Xiang Zhang, Stefan K. Lhachimi, and Wolf H. Rogowski. "Reporting Quality of Discrete Event Simulations in Healthcare—Results From a Generic Reporting Checklist". In: *Value in Health* 23.4 (Apr. 2020), pp. 506–514. ISSN: 1098-3015. DOI: 10.1016/j.jval.2020.01.005. URL: <https://www.sciencedirect.com/science/article/pii/S1098301520300401> (visited on 05/21/2024).
- [21] Robert Gentleman and Duncan Temple Lang. "Statistical Analyses and Reproducible Research". In: *Journal of Computational and Graphical Statistics* 16.1 (Mar. 2007). Publisher: Taylor & Francis .eprint: <https://doi.org/10.1198/106186007X178663>, pp. 1–23. ISSN: 1061-8600. DOI: 10.1198/106186007X178663. URL: <https://doi.org/10.1198/106186007X178663> (visited on 05/17/2024).

- [22] Ben Marwick, Carl Boettiger, and Lincoln Mullen. "Packaging Data Analytical Work Reproducibly Using R (and Friends)". In: *The American Statistician* 72.1 (Jan. 2018). Publisher: Taylor & Francis. eprint: <https://doi.org/10.1080/00031305.2017.1375986>, pp. 80–88. ISSN: 0003-1305. DOI: 10.1080/00031305.2017.1375986. URL: <https://doi.org/10.1080/00031305.2017.1375986> (visited on 05/20/2024).
- [23] Florio Arguillas et al. "10 Things for Curating Reproducible and FAIR Research". eng. In: (June 2022). Publisher: Zenodo. DOI: <https://doi.org/10.15497/RDA00074>. URL: <https://zenodo.org/records/6797657> (visited on 05/20/2024).
- [24] Ben B. Blohowiak et al. *Badges to Acknowledge Open Practices*. en. Publisher: OSF. Sept. 2023. URL: <https://osf.io/tvyxz/> (visited on 05/20/2024).
- [25] NISO Reproducibility Badging and Definitions Working Group. *Reproducibility Badging and Definitions*. en. Jan. 2021. DOI: 10.3789/niso-rp-31-2021. URL: <https://www.niso.org/publications/rp-31-2021-badging>.
- [26] Institute of Electrical and Electronics Engineers (IEEE). *About Content in IEEE Xplore*. URL: <https://ieeexplore.ieee.org/Xplorehelp/overview-of-ieee-xplore/about-content> (visited on 05/20/2024).
- [27] Tom E. Hardwicke and Simine Vazire. "Transparency Is Now the Default at Psychological Science". In: *Psychological Science* 0.0 (2023). DOI: <https://doi.org/10.1177/09567976231221573>.
- [28] Caspar J. Van Lissa et al. "WORCS: A workflow for open reproducible code in science". en. In: *Data Science* 4.1 (Jan. 2021). Publisher: IOS Press, pp. 29–49. ISSN: 2451-8484. DOI: 10.3233/DS-210031. URL: <https://content.iospress.com/articles/data-science/ds210031> (visited on 05/20/2024).
- [29] Association for Computing Machinery (ACM). *Artifact Review and Badging Version 1.1*. en. Aug. 2020. URL: <https://www.acm.org/publications/policies/artifact-review-and-badging-current> (visited on 05/20/2024).
- [30] Springer Nature. *Springer Nature Open data badge*. July 2018. URL: <https://badgr.com/public/badges/xhW4FLHBRe6Tzz2Cj4Q1tA> (visited on 05/20/2024).
- [31] Association for Psychological Science (APS). *Psychological Science Submission Guidelines*. Dec. 2023. URL: https://www.psychologicalscience.org/publications/psychological_science/ps-submissions (visited on 05/20/2024).