

Ten Simple Rules for Writing Dockerfiles for Reproducible Research

Daniel Nüst^{*}, Vanessa Sochat[†], Ben Marwick[‡], Stephen J. Eglen[§], Tim Head^{||}, Tony Hirst

^{*} Corresponding author: daniel.nuest@uni-muenster.de

Abstract

Containers are greatly improving computational science by packaging software and data dependencies. In a scholarly context, transparency and support of reproducibility are the largest drivers for using these containers. It follows that choices that are made with respect to building containers can make or break a workflow's reproducibility. The build for the container image is often created based on the instructions in a plain-text file. For example, one such container technology, Docker, provides instructions using a **Dockerfile**. By following the rules in this article researchers writing a **Dockerfile** can effectively create containers suitable for sharing with fellow scientists, for including in scholarly communication such as education or scientific papers, and for an effective and sustainable personal workflow.

Introduction

[SJE: this needs to be snappier as an introduction. Something like... Computing infrastructure has advanced to the point where not only can we share data underlying research articles, we can share the code that processes these data and even entire workflows. This improves the reproducibility of our research. A research compendium ...] With access to version control systems (VCS, [1]) and collaboration platforms based on these, such as GitHub or GitLab, it has become increasingly easy to not only share abstract algorithms and study designs, but also instructions for executing whole research workflows, including building and testing of the software used. The publication of these instructions and related configuration files are a response to the increasing complexity of computer-based research, which is too complicated for “papers” based on the traditional journal article format to fully communicate the details of research [2], but where the actual contribution to knowledge includes the full software environment that produced a result [3]. A *research compendium* is a compilation of data, code, and documentation that accompanies, or is itself, a scholarly publication for increased transparency and reproducibility (cf. <https://research-compendium.science>). Within such a compendium, it is desirable to include instructions, i.e. a human- and machine-readable recipe, for building containers that capture the computing environment, i.e., all software and data dependencies. By providing this recipe, authors of scientific articles greatly improve their work's level of documentation, transparency, and reusability. Such practice is one important part of common practices for scientific computing [4,5], with the result that it is much more likely both the author and others are able to reproduce and extend an analysis workflow. The containers built from these recipes are portable encapsulated snapshots of a specific computing environment. Such

containers have been demonstrated for capturing scientific notebooks [6] and reproducible workflows [7].

While there are several tutorials for using containers for reproducible research [8–12], there is no detailed manual for how to write the actual instructions to create the containers beyond generic best practices outside of the scientific domain [13]. Several platforms for facilitating reproducible research are built on top of containers [14–18], but they hide most of the complexity from the researcher. Because “*the number of unique research environments approximates the number of researchers*” [18], sticking to conventions helps every researcher to understand, modify, and eventually write container recipes suitable for their needs. Even if they are not sure how the technology behind them actually works, researchers can leverage containerisation, but should craft their own definition of computing environments following good practices. The practices outlined here are strongly related to software engineering in general and research software engineering (RSEng) in particular, which is concerned with quality, training, and recognition of software in science [19]. Research Software Engineers (RSEs) are not the target audience for this work, but we want to encourage you to reach out to your local or national RSE community if your needs go beyond the contents detailed here.

While there are many different container technologies, this article focuses on Docker [20]. Docker is a highly suitable tool for reproducible research (e.g., [21]) and our observations indicate it is the most widely used container technology in academic data science. The goal of this article is to guide you to write a **Dockerfile** so that it best facilitates interactive development and computer-based research, as well as the higher goals of reproducibility and preservation of knowledge. Such practices are generally not part of generic containerization tutorials and are rarely found in published **Dockerfiles**, which are often used as templates by novices. The differences between a helpful, stable **Dockerfile** and one that is misleading, prone to failure, and full of potential obstacles, are not obvious, especially for researchers who do not have extensive software development experience or formal training. A commitment to this article’s rules can ensure that workflows are reproducible and reusable, that computing environments are understandable by others, and researchers can collaborate effectively. Their application should not be triggered by the publication of a finished project, but be weaved into day-to-day habits (cf. thoughts on openness as an afterthought by [22] and on computational reproducibility by [3]).

To start with, we assume you have a scripted scientific workflow, i.e. you can, at least at a certain point in time, execute the full process with a single command, for example `make my_workflow`, `Rscript analysis.R`, or `python3 paper-plots.py`. (A key constraint is Docker containers can only run open source software, so environments like Mathematica and Matlab cannot currently be easily used within containers.) This execution should be able to be triggered by command-line instruction, which is possible for all generic programming languages and widely used workflow tools, but may also be opening and starting a process with a graphical user interface. A workflow that does not support scripted execution is out of scope for reproducible research, and does not fit well with containerization.

Docker

Docker is a container technology that is widely adopted and supported on many platforms, and has become highly useful in science. They are distinct from virtual machines (VM) or hypervisors as they do not emulate hardware or operating system kernels, and thus do not require the same system resources. To create Docker containers, we write text files that follow a particular format called **Dockerfile** [23]. **Dockerfiles** are machine- and human-readable recipes for building images. Images are

inert, immutable, read-only files that include the application, e.g., the programming language interpreter needed to run a workflow, and the environment required by the application to run. A **Dockerfile** consists of a sequence of instructions to copy files and install software, which add layers to the image. [SJE: “See figure 1 for an example of a complete Dockerfile that we will refer to throughout this article”. Remember: this might be the first paper that someone reads about Docker; if it is critical that the reader has already seen Dockerfiles, we should point them somewhere suitable??.] These layers can be used for caching across image builds, which is important for minimizing build and download times. The images have a main executable that is started when they are run as stateful containers, which are the running instances of Docker images. Containers can be modified, stopped, restarted and purged. Dockerfiles, like Makefiles, are meant to be useful for both humans to read and for computers to execute.

While Docker was the original technology to support this format, other container technologies support the format including podman/buildah supported by RedHat, kaniko, img, and buildkit. The Singularity container software [24] is optimized for high performance computing and although it uses its own format, the *Singularity recipe*, it can import Docker images directly from a Docker registry. Although the Singularity recipe format is different, the rules here are transferable to some extent. While some may argue for reasons to not publish reproducibly, e.g., lack of time and incentives, reluctance to share (cf. [25]) and there are substantial technical challenges to maintain software and documentation, providing a **Dockerfile**, a pre-built Docker image, or other type of container should become an increasingly easier task for the average researcher. If a researcher is able to find and create containers or write a **Dockerfile** to address their most common use cases, then arguably it will not be extra work after this initial set up (cf. README of [26]). In fact, the **Dockerfile** itself can be a powerful documentation to show from where data and code was derived, i.e. downloaded or installed, and consequently where a third party might obtain them again.

1. Consider tools to assist with Dockerfile generation

Rule 1 could informally be described as “Don’t bother to write a Dockerfile!” Writing a **Dockerfile** from scratch can be difficult and even experts sometimes take shortcuts. Thus, it is a good strategy to first look to tools that can help to generate a **Dockerfile** for you. Such tools have likely thought about and implemented good practices, and they may have added newer practices when reapplied at a later point in time. Therefore the most important rule is to apply a multi-step process for your specific use case. First, you want to determine if there is an already existing container that you can use, and in this case, use it and add to your workflow documentation instructions for doing so. As an example, you might be doing some kind of interactive development. For interactive development environments such as notebooks and development servers or databases, you can readily find containers [SJE: where do you look to find containers?] that come installed with all the software that you need. Second, in the case that there isn’t an existing container for your needs, you might next look to well-maintained tools to help with **Dockerfile** generation. These tools can add required software packages to an existing image without you having to manually write a **Dockerfile** at all. Well-maintained includes the tool’s own stability and usability, but also that suitable base images are used, likely from the official Docker library [27], to ensure that the container has the most recent security fixes for the operating system in question. Third, you may want to write your own **Dockerfile** from scratch if these tools do not suffice for your needs. *In this case, follow the remaining rules.*

Tools for container generation

Repo2docker [17] is a tool that is maintained by Project Jupyter that can help to transform a [SJE github/gitlab? any other repos?] repository with common simple configuration files for defining software dependencies and version into a container. As an example, we might install `jupyter-repo2docker` and then run it against a repository with a `requirements.txt` file, an indication of being a Python workflow with dependencies on the Python Package Index (PyPI) with the following command.

```
jupyter-repo2docker https://github.com/norvig/pytudes
```

The resulting container image installs the dependencies listed in the requirements file, along with providing an entrypoint to run a notebook server to interact with any existing workflows in the repository. A precaution that needs to be taken is that the default command above will create a home for the current user, meaning that the container itself would not be ideal to share, but rather any researchers interested in interaction with the code inside should build their own container. For this reason, it is good practice to look at any help command provided by the tool and check for configuration options for user names, user ids, and similar. It's also recommended to add custom labels to your container build to define metadata for your analyses even if you are using tools (see Rule 6).

Additional tools to assist with writing Dockerfiles include `containerit` [28] and `dockta` [29]. `containerit` automates the generation of standalone Dockerfiles for workflows in R. It can provide a starting point for users unfamiliar with writing Dockerfiles, or together with other R packages provide a full image creation and execution process without having to leave an R session. `dockta` supports multiple programming languages and configurations files, just as `repo2docker`, but attempts to create readable Dockerfiles compatible with plain Docker and to improve user experience by cleverly adjusting instructions to reduce build time.

2. Use versioned and automatically built base images

[SJE: this feels WAY too long. can't we say it much shorter, i.e. it is good to pin version numbers if you want a stable archive of reproducible code, or use "latest" to check that your code continues to work as other software environments evolve. I think having a table here with some recommendations for e.g. python, R and julia, base images would go a long way to practically serve about 80% + of data science users? The bullet pointed list at the end of the section is probably most useful.]

It's good practice to use base images that are maintained by the Docker library, so called "*official images*" [30], which benefit from a review for best practices and vulnerability scanning [13]. You can identify these images by the missing the user portion of the image name, which comes before the /, e.g., `r-base` or `python`. There are also repositories that are maintained by organisations or fellow researchers. While some organizations can be trusted to update containers with security fixes (e.g., the R Project or Jupyter), for most individual accounts that provide ready to use images, it is likely that these will not be updated regularly. It is even possible that images or Dockerfiles may disappear, or that images are published with malicious intent. Images furthermore have *tags* associated with them with specific meanings, usually as a version number but sometimes also variants of images. Tags are defined at image built time and appear in image name after the : when you use an image, e.g., `python:3.6`. By *convention* a missing tag is assumed to be the word `latest`. When you re-build you images locally using the same tag, the tag is always associated with the newest built.

But the image also has an automatically generated image name using a hash, which will *not* be updated. So if you tag an image, use the tag when you start a container.

Because of these options, a good understanding of how base images and image tags work is crucial, as the image and tag that you choose has important implications for your own image. You should know how to use **FROM** instructions to trace all the way up to the original container base, an abstract base called **scratch**. Doing this kind of trace is essential because it makes clear all of the steps that were taken to generate your container and who took them. If you want to build on a third party image, carefully consider the author/maintainer and *never* use an image without a published **Dockerfile**. If you want to use an unofficial image, do save a copy of the **Dockerfiles** created by others. Alternatively, copy relevant instructions into your own **Dockerfile**, acknowledge the source in a comment, and configure an automated build for your own image. The official images and images by established organisations should be your preferred source of snippets if you craft your own **Dockerfile**. Automated builds can be complex to set up, and details are out of scope of this article. The documentation for, and toolsets provided by container registries such as Docker Hub or GitLab as well as CI platforms such as GitHub actions or CircleCI can help you get started; tools may also be available to automate builds using intermediate tools, such as the **repo2docker** Github action [31,32]. You should avoid publishing an image in a public registry with **docker push**, because it is opaquely breaking the linkage between **Dockerfile** and image.

A tag like **latest** is good in that security fixes and updates are likely present, however it is bad in that it is a moving target so that it is more likely that an updated base could break your workflow. Other tags that should be avoided are **dev**, **devel**, or **nightly** that may provide possibly unstable versions of the software. Official builds are often tagged with corresponding release version numbers, the newest of which also has the **latest** tag. For example, if you build a container with **python:latest** that only supports Python version 3, when the base is updated to Python 4 it's likely that your software won't work as expected. Therefore a tag like **python:3.5** is a preferable choice, where the lack of the third component of the version, i.e., **3.5.x**, ensures that security patches and bugfixes *will* be included. These *bugfix* version changes should not break your code, i.e., are backwards compatible (cf. semantic versioning, [33]), and if they break your code you might rely on a bug in the used software. When you choose a base image, try to choose one with a Linux distribution that supports the software stack you are using, and also take into account the bases that are widely used by your community. As an example, Ubuntu is heavily used for geospatial research, and so the **rocker/geospatial** image would be a good choice for spatial data science with R, or **jupyter/tensorflow-notebook** could be a good choice for machine learning with Python.

Container tags may also provide an indication as the size of a container. Smaller images are often indicated by having **slim** or **minimal** as part of the tag. If you need to build smaller containers, consider a **busybox** base or a *multi-stage build* [34], which allows you to selectively copy files from one build step to another. Also take into account the software that you actually need, but mostly for clarity. In scientific contexts, the image size will be often the least of your concerns, because data size is much larger or image portability is not time critical.

Base images that have complex software installed (e.g. machine learning tool stacks, specific BLAS library) are helpful and fine to use. If you do not want to rely on a third party or other individual to maintain the recipe and container, you should copy the **Dockerfile** to your own repository, and also provide a deployment for it via your own automated build. Trusting that the third party **Dockerfile** will persist for your analyses may be risky because that **Dockerfile** could disappear without warning. Here is a selection of communities that produce widely used regular builds and updates:

- Rocker for R [21] 221
- Docker containers for Bioconductor for bioinformatics 222
- NeuroDebian images for neuroscience [35] 223
- [Jupyter Docker 224
Stacks](https://jupyter-docker-stacks.readthedocs.io/en/latest/index.html for 225
Notebook-based computing 226
- Taverna Server for running Taverna workflows 227

For example, here is how we would use a base image `r-ver` with tag 3.5.2 from the `rocker` organization on Docker Hub (`docker.io`). 228
229

FROM rocker/r-ver:3.5.2 230

3. Use formatting, document within, and favor clarity 231 232

It is good practice to think of the `Dockerfile` as a human *and* machine readable file. 233
This means that you should use indentation, new lines, and comments to make your 234
`Dockerfiles` well documented and readable. Specifically, carefully indent commands 235
and their arguments to make clear what belongs together, especially when connecting 236
multiple commands in a `RUN` instruction with `&&`. Use `\` at the end of a line to break a 237
single command into multiple lines. This will ensure that no single line gets too long to 238
comfortably read. Use long versions of parameters for readability (e.g., `--input` instead 239
of `-i`). When you need to change a directory, use `WORKDIR`, because it not only creates 240
the directory if it doesn't exist but also persist the change across multiple `RUN` 241
instructions. 242

You can use a linter [36] to avoid small mistakes and follow good practices from 243
software development communities. The consistency added by linting also helps keeping 244
your edits to a `Dockerfile` in a VCS meaningful (see Rule 9). Note however that a 245
linter's rules may not primarily serve the intention of reproducible scientific workflows. 246

As you are writing the `Dockerfile`, be mindful of how other people will read it. Are 247
your choices and commands being executed clear, or is further comment warranted? To 248
assist others in making sense of your `Dockerfile`, you can add comments that include 249
links to online forums, code repository issues, or VCS commit messages to give context 250
for your specific decisions, for example this `Dockerfile` by Kaggle does a good job at 251
explaining the reasoning and steps of the contained instructions. Comments should 252
include helpful usage guides and links for readers inspecting the `Dockerfile`, including 253
future you. Dependencies can be grouped in this fashion, which also makes it easier to 254
spot changes if your `Dockerfile` is managed in a VCS (see Rule 9). It can even be 255
helpful to include comments about commands that did not work so you do not repeat 256
past mistakes. If you find that you need to remember an undocumented step, that is an 257
indication that it should be documented in the `Dockerfile`. 258

Labels are useful to provide a more structured form of documentation about 259
software, especially for users of supporting infrastructure that might not expose the 260
actual `Dockerfile` (see Rule 6). It is often helpful to provide commented lines with 261
`docker build` and `docker run` within the `Dockerfile` to show how to build and run 262
the image, even though these lines are not necessary for a valid `Dockerfile`, this is a 263
convenient place to record them. These comments can be especially relevant if volume 264
mounts or ports are important for using the container, and by putting them at the end 265
of the file, they are more likely to be seen, easy to copy-paste after a container build, 266
and to be in consistent state compared to documentation in another file. Here is an 267
example of a commented section to show build and usage. 268

```

# Build the images with
##> docker build --tag great_workflow .
# Run the image:
##> docker run --it --port 80:80 --volume ./input:/input --name gwf great_workf
# Extract the data:
##> docker cp gwf:/output/ ./output

```

[SJE: the run command above spills over the right hand side of the printed page so you can't see the 'workflow' word at the end of the command. Reformat?]

If you were to discover a previously written **Dockerfile** and not remember the container identifier you used, it would be represented in the **Dockerfile** in the **--name** parameter, preserved in a comment. Following a common coding aphorism, we might say *"A Dockerfile written three months ago may just as well have been written by someone else"*. Here is a selection of typical kinds of comments that are useful to include in a **Dockerfile**:

```

# apt-get install specific version, use 'apt-cache madison <pkg>'
# to see available versions
RUN apt-get install python3-pandas=0.23.3+dfsg-4ubuntu1

# RUN command spreading several lines
RUN R -e 'getOption("repos")' && \
    install2.r \
    fortunes \
    here

# this library must be installed from source to get version newer
# than in sources
RUN git clone http://url.of/repo && \
    cd repo && \
    make build && \
    make install

# following commands from instructions at <LINK HERE>

```

Clarity is nearly always more important than brevity. For example, if your container uses a script to run a complex install routine, instead of removing it from the container upon completion, which is commonly seen in production **Dockerfiles** aiming at small image size, you should keep the script in the container for a future user to inspect. Depending on the programming language used, your project may already contain files to manage dependencies and you may use a package manager to control this aspect of the computing environment. This is a very good practice and helpful, though you should consider the externalization of content to outside of the **Dockerfile** (see Rule 5). A single long **Dockerfile** with sections and helpful comments can be complete and thus more understandable than a collection of separate files.

Generally, aim to design the **RUN** instructions so that each performs one scoped action, e.g., download, compile, and install one tool. Each instruction will result in a new layer, and reasonably grouped changes increase readability of the **Dockerfile** and facilitate inspection of the image, e.g., with tools like **dive** [37]. Convolved **RUN** instructions can be acceptable to reduce the number of layers, but careful layout and consistent formatting should be applied. A **RUN** instruction longer than a page requires scrolling, diminishing readability. This may be challenging for the next reader to digest and you should consider splitting it up, being aware of the extra layers added to the image.

[SJE: I thought the following point about caching to be significant enough to be its own point. Ordering the dockerfile to minimise rebuild times etc during development] When you install several system libraries, it is good practice to add comments about why the dependencies are needed. This way, if a piece of software is removed from the container, it will be easier to remove the system dependencies that are no longer needed. If you intend to build the image more than once (perhaps during development) and you can take advantage of build caching to avoid execution of time-consuming instructions, e.g., install from a remote resource or a file that gets cached. You should list commands *in order* of least likely to change to most likely to change and use the `--no-cache` flag to force a re-build of all layers. Docker will execute the instructions in the order as they appear in the **Dockerfile**. When one instruction is completed, the result is cached, and the build moves to the next one. If you change something in the Dockerfile, and rebuild the container, each instruction is inspected in turn. If it hasn't changed, the cached layer is used and the build progresses. If the line has changed, that build step is executed afresh and then every following instruction will have to be rebuilt in case the line that changed influences something for a following instruction. A recommended ordering based on these considerations is as follows.

1. system libraries
2. language-specific libraries or modules
3. from repositories (binaries)
4. from source
5. own software/scripts (if not mounted)
6. labels
7. RUN/ENTRYPOINT

Finally, as a supplement to content inside the **Dockerfile**, it is good practice to also write a section in a **README** alongside the **Dockerfile** for exactly how to build, run, and otherwise interact with the container. [SJE: really? Doesn't that mean you end up writing the instructions twice – once in the Dockerfile as a comment, and once in the README? Surely the README could just say: “Please read the end of the Dockerfile for instructions on how to run the docker image.”] If a pre-built image is provided on Docker Hub, you should direct the user to it in your **README**.

4. Define version numbers for reproducible builds

The reproducibility of your **Dockerfile** heavily depends on how well you define the versions of software to be installed in the image. The more specific, the better, because using the desired version leads to reproducible builds. The practice of specifying versions of software is called *version pinning* (e.g., on **apt**: <https://blog.backslasher.net/my-pinning-guidelines.html>). For stable workflows in a scientific context, it is generally advised to freeze the computing environment explicitly and not rely on the “current” or “latest” software, which is a moving target.

System libraries

System library versions can largely come from the base image tag that you choose to use, e.g., **ubuntu:18.04**, because the operating system's software repositories are very unlikely to introduce breaking changes, but predominantly fix errors with newer versions. However, you can also install specific versions of system packages with the respective package manager. For example, you might want to demonstrate a bug, prevent a bug in an updated version, or pin a working version if you suspect an update could lead to a problem. Generally, system libraries are more stable than software modules supporting

analysis scripts, but in some cases they can be highly relevant to your workflow. *Installing from source* is a useful way to install very specific versions, however it comes at the cost of needing to build libraries. Here are some examples of terminal commands that will list the currently installed versions of software on your system:

- Debian/Ubuntu: `dpkg --get-architecture`
- Alpine: `apk -vv info|sort`
- CentOS: `yum list installed` or `rpm -qa`

Extension packages and programming language modules

In the case of needing to install packages or dependencies for a specific language, package managers are a good option. Package managers generally provide reliable mirrors or endpoints to download software, many packages are tested before release, and most importantly they provide access to specific versions. Most package managers have a command line interface that can be used from `RUN` commands in your `Dockerfile`, along with various flavors of “freeze” commands that can output a text file listing all software packages and versions (cf. <https://markwoodbridge.com/2017/03/05/jupyter-reproducible-science.html> cited by [6]) The biggest risk with using package managers with respect to `Dockerfiles` is outsourcing configuration. As an example, here are configuration files supported by commonly used languages in scientific programming:

- Python: `requirements.txt` (pip tool, [38]), `environment.yml` (Conda, [39])
- R: `DESCRIPTION` file format [40] and `r` (“little R”, [41])
- JavaScript: `package.json` of `npm` [42]
- Julia: `Project.toml` and `Manifest.toml` [43]

In some cases (e.g., Conda) the package manager is also able to make decisions about what versions to install, which is likely to lead to a non-reproducible build. In all of the above, the user is required to inspect the file or the build to see what is installed. For this reason, in the case of having few packages, it is suggested to write the install steps and versions directly into the `Dockerfile` (also for clarity, see Rule 3). For example, the `RUN` instruction here:

```
RUN pip install geopy==1.20.0 && \
    pip install uszipcode==0.2.2
```

serves as more clear documentation in a `Dockerfile` than a `requirements.txt` file that lists the same:

```
RUN pip install -r requirements.txt
```

This modularisation is a potential risk for understandability and consistency, which can be mitigated by carefully organizing all these files in the same version-controlled project. You can also use package manager to install software from source code `COPY`d into the image (see Rule 5). And finally, you can use many package managers to install software from source obtained from code management repositories, e.g., installing a specific tool identified by a GitHub version tag or commit hash. Be aware of the risk of depending on such repositories, especially if they are out of your control, although the installation command with a full URL should allow readers of your `Dockerfile` to dig deeper if problems arise and such an installation from source gives you a lot of freedom (including the one to shoot yourself in the foot). The version pinning capabilities of these file formats and package managers are described in their respective documentation.

As a final note on software installation, you should be aware of the **USER** instruction in a **Dockerfile**. It is common to install all software as the default user **root**. In fact it is *required* for system dependencies. You can switch the user running subsequent **RUN**, **CMD** and **ENTRYPOINT** instructions with the **USER** instruction, but you have to make sure yourself that the user exists within the image. This is useful if a software aborts installations as the **root** user, or to avoid permission problems when using the container (see Rule 5). Be sure to document a non-default user in detail and to examine the used users in base images if you encounter errors on user permissions. [SJE: okay, this para confused me – just tell me what to do! should I not bother creating a user, and run as root, or should I create a user? rstudio complicates matters here I think, don't they create a “rstudio” username?]

5. Add user scripts and data into containers by mounting directories

The role of containers is to provide the methods, not to encapsulate datasets. It is better to insert data files and scripts files from the local machine into the container at run time, and using the container image primarily for the software and dependencies. This insertion is achieved by using *bind mounts*. Mounting these files is preferable to using the **ADD/COPY** instructions in the **Dockerfile**, because files persist when the container instance or image is removed from your system, and the files are more accessible to other users if published in an online research compendium outside of the image. You should use **COPY** if you actually copy files because it is explicit. **ADD** supports additional features and you should only use it if you leverage them: (i) a URL as a source, for which you may use it instead of a **RUN** instruction to download files, and (ii) a local tar file, which you should avoid because such an archive file will be harder for future readers of your **Dockerfile** to make sense of. Unless specific features are needed, bind mounts are also preferable to storage volumes.

Standalone *script files* are distinct from other software as they are not packaged and versioned, and are treated as files, not as managed software. If you developed software for a specific analysis in form of a software package, you should publish in a public source code or software repository and follow Rule 4 for installing it. You should avoid installing software packages from source after **COPY**ing the code into the image, because the connection between the file outside of the image and the one copied in is easily lost (cf. Rule 5). Consider using the suitable package system with pinned versions even for small scripts if the functions are reusable across datasets or workflows by you or others. If you cannot publish a software project albeit it being in a language's package structure, you should make sure the files are always published together with the **Dockerfile** (see Rule 9).

Storing *data files* outside of the container further allows handling of very large datasets and datasets with data worthy of protection, e.g., proprietary data or private information. Do not include such data in an image. To avoid publishing sensitive data by accident, you can add the data directory to the **.dockerignore** file, which excludes files and directories from the build context, i.e., the set of files considered by **docker build**. Ignoring data files also speeds up the build in case of very large files or many small files, the latter of which makes it common to ignore for example the **.git** folder of the Git VCS. You should include dummy or test data into the image to be able to ensure that a container is functional without a larger custom dataset, e.g., for automated tests or instructions in the user manual. For all these cases you should provide clear instructions for users in the **README** how to obtain actual or dummy data. When publishing your workspace, e.g., on Zenodo, having data and script contents as

regular files outside of the container makes them more accessible to others, for example for reuse or analysis.

A mount can also be used to access *output data* from a container. This can be an extra mount, or the same **data** directory. Alternatively, you can use the **docker cp** command to access files from a running or stopped container, but this requires a specific handling, e.g., naming the container when starting it or using multiple shells, that require very detailed instructions for users.

You can use the **-v/--volume** or **--mount** flags to **docker run** to configure bind mounts of directories or files [44], including options, as shown in the following examples. If the target path exists within the image, the bind mount will replace it for the started container.

```
# mount directory
docker run --volume /home/user/project:/project mycontainer

# mount directory as read-only
docker run --volume /home/user/project:/project:ro mycontainer

# mount multiple directories, one with write access relative to current path
docker run --volume /home/user/Download/article-x-supplement/data:/data:ro --v
```

How your container expects external resources to be mounted into the container should be included in the example commands (see Rule 3). In these commands you can also make sure to avoid issues with file permissions by using Docker’s **--user** option. For example, by default, writing a new file from inside the container will be owned by user **root** on your host, because that is the default user within the container.

6. Capture structured environment metadata

Labels and build arguments can be very helpful to both provide metadata and allow for customization of a build. This structured metadata is machine-readable may help tools more than humans (see Rule 3 for user documentation within a **Dockerfile**). In addition to the metadata creation by users, Docker captures useful information in the image metadata automatically, such as the version of Docker used for building the image, which you can access with the **docker inspect** command.

Labels

Labels serve as structured metadata that can be exposed by APIs, e.g., <https://microbadger.com/labels>, along with tools to inspect the container binaries, e.g., **docker inspect**. For example, software versions, maintainer contact information, along with vendor specific metadata are commonly seen. The OCI Image Format Specification provides some common label keys (see the “Annotations” section in [45]) to help standardize field names across container tools, as shown below. These labels match the **org.label-schema.*** specification, which has been deprecated in favour of the new namespace but are still found a lot in existing containers.

[SJE: this was new to me. On the one hand, I like the idea of LABEL, but have they really been adopted by the community? The document from 2017 cited above hasn’t moved on since Jun 2017. Is this really a thing? And is the Dockerfile the right place for all this metadata? It feels like a bit too fine-grained for my liking. But if this is good practice, then I shall start adopting it!]

```

LABEL org.opencontainers.image.created='2019-12-10' \
      org.opencontainers.image.authors='author@example.org' \
      org.opencontainers.image.url='https://github.com/nuest/ten-simple-rules' \
      org.opencontainers.image.documentation='https://github.com/...' \
      org.opencontainers.image.version='0.0.1'

LABEL org.opencontainers.image.vendor='nuest' \
      org.opencontainers.image.title='Demo title' \
      org.opencontainers.image.description='Demo description'

```

You can add multiple fields within the same instruction. Important metadata attributes to include as labels would include any of the following, ideally with globally unique identifiers:

- Author and contact (e.g., email, project website, or ORCID; you will often see the deprecated MAINTAINER instruction - use a label instead)
- Research organizations (identified with <https://ror.org/>)
- Funding agency/grant number
- A repository link where the `Dockerfile` is published, e.g., a GitHub project or a repository record with a DOI, e.g., Zenodo, where you can pre-register a DOI and add it to your `Dockerfile` before publishing the record
- License

Proper software citation is still a work in progress [46] and you should follow current recommendations of projects such as CodeMeta (<https://codemeta.github.io/>) for describing software and the Citation File Format (<https://citation-file-format.github.io/>) to enable citations of software.

Build arguments

Build arguments can provide more dynamic metadata and also allow for customization of a build. As an example, the following build argument would default to 1.0.0 but allow you to change it with `--build-arg MYVERSION=2.0.0` when building the image:

```
ARG MYVERSION=1.0.0
```

Along with specifying versions, e.g., a git commit hash, or adding a date and timestamp, build arguments can be useful to provide the context of the build, e.g., building user, production versus development environment, and automated or not. Examples of build arguments that are useful to include to describe a container are:

7. Enable interactive usage and one-click execution

Containers are very well suited for day-to-day development tasks (see also Rule 10), because they support common interactive environments for data science and software development. But they are also useful for a “headless” execution of a full workflows, e.g. as demonstrated in [47]. A workflow that does not support headless execution may even be seen as irreproducible. These two usages can be configured by the `Dockerfile`’s author and exposed to the user based on the `Dockerfile`’s `CMD` and `ENTRYPOINT` instructions (see Rule 8 for having multi-container environments). An images main purpose is reflected by the default process and configuration, though the `CMD` and `ENTRYPOINT` can also be changed at runtime. It is considered good practice to have a combination of default command and entrypoint that meets reasonable user

expectations. For example, a container known to be a workflow should execute the workflow, or provide instructions for how to do so. Since you are likely the primary user of the `Dockerfile` and image, so you should choose a selection or combination that works for you, only then accomodate other user's needs. A possible weakness with using containers is the limitation on only providing one default command and entrypoint. However tools, e.g., The Scientific Filesystem [48], have been developed to expose multiple entrypoints, environments, help messages, labels, and even install sequences. With plain Docker, you can override the defaults as part of the `docker run` command or in an extra `Dockerfile` using the primary image as a base. In any case you should document different variants, if you choose to provide them, in a `Makefile` [49]. To support both one click execution and interactive interfaces and even allow for custom configuration, it's helpful to expose settings via a configuration file which can be bound from the host, via environment variables [50], or special Docker-based wrappers such as Kliko [51].

[SJE: I think I get the two style of working – batch and interactive might be my way of thinking about them. But it feels a bit clunky above and certainly too detailed. might be better served with two example dockerfiles: once doing a batch analysis, and the other whose goal is to allow some interaction.]

Interactive graphical interfaces, such as RStudio, Jupyter, or Visual Studio Code, can run in a container to be used across platforms via a regular web browser. The HTML-based user interface is exposed over HTTP. Use the `EXPOSE` instruction to document the ports of interest for both humans and tools, because they need to be bound to the host to be accessible to the user using the `docker run` option `-p/--publish <host port>:<container port>`. The container should also print to the screen the used ports along with any login credentials needed. For example, as done in the last few lines of the output of running a Jupyter Notebook server locally (lines abbreviated).

```
docker run -p 8888:8888 jupyter/datascience-notebook:7a0c7325e470
```

```
[...]
[I 15:44:31.323 NotebookApp] The Jupyter Notebook is running at:
[I 15:44:31.323 NotebookApp] http://9027563c6465:8888/?token=6a92d [...]
[I 15:44:31.323 NotebookApp] or http://127.0.0.1:8888/?token=6a92 [...]
[I 15:44:31.323 NotebookApp] Use Control-C to stop this server and [...]
```

[SJE: that's good, does rstudio do similar? need to check.]

A person who is unfamiliar with Docker but wants to use your image may rely on graphical tools like Kitematic [52] or ContainDS for assistance in managing containers on their machine without using the Docker CLI.

Interactive usage of a command-line interfaces is quite straightforward to access from containers, if users are familiar with this style of user interface. Running the container will provide a shell where a tool can be used and help or error messages can assist the user. For example, complex workflows in any programming language can, with suitable pre-configuration, be triggered by running a specific script file. If your workflow can be executed via a CLI you may use that to validate correct functionality of an image in automated builds, e.g. using a small toy example and checking the output, by checking successful responses from HTTP endpoints provided by the container, e.g. via an HTTP response code of 200, or by using a controller such as Selenium [53].

The followig example runs a simple R command counting the lines in this articles source file. The file path is passed as an environment variable.

```

docker run \
  --env CONFIG_PARAM="/data/ten-simple-rules-dockerfiles.Rmd" \
  --volume $(pwd):/data \
  jupyter/datascience-notebook:7a0c7325e470 \
  R --quiet -e "
l = length(readLines(Sys.getenv('CONFIG_PARAM')));
print(paste('Number of lines: ', l))
"

> l = length(readLines(Sys.getenv('CONFIG_PARAM')));
> print(paste('Number of lines: ', l))
[1] "Number of lines: 568"

```

If there is only a regular desktop application, the hosts window manager can be connected to the container. This has notable security implications, which are reduced by using the “X11 forwarding” natively supported by Singularity [54], which can execute Docker containers, or by leveraging supporting tools such as **x11docker** [55]. Bidge containers [56] and exposing a regular desktop via the browser (e.g., for Jupyter Hub [57]) are further alternatives. This variety of approaches render seemingly more convenient uncontainerised environments, i.e. just using the local machine, unnecessary in favour of reproducibility and portability.

8. Establish templates for new projects

[SJE: this seemed a bit lacking in meat, compared to the other rules. Do you have an example Dockerfile template that you are thinking of?]

It is likely going to be the case that over time you will develop workflows that are similar in nature to one another. In this case, you should consider adopting a standard workflow that will give you a clean slate for a new project. If you decide to build your own standard, collaborate with your community during development of the standard to ensure it will be useful to others. Part of your project template should be a protocol for publishing the container image to a suitable container registry, and taking into consideration of how the code can be given a DOI or proper publication (e.g., Zenodo, Journal of Open Source Software).

As an example, cookie cutter templates [58], project starter kits [REF], or community templates (e.g., [59]) can provide files, directory organization, and build instructions that include basic steps for getting started. A good project template should get you started with a template for documentation, setting up testing via continuous integration (CI), building a container, and even choosing a license. In the case of using a common **Dockerfile** or base, Docker’s build caching will take shared lines into account and speed up build time even between projects. When developing or working on projects with containers you can switch between isolated project environments by stopping the container and restarting it when you are ready to work again, even on another machine or in a cloud environment. You can even run projects in parallel without interference. At most, if a port is shared by two projects to expose a user interface in the browser, you would need to configure non-conflicting ports.

In the case of more complex applications that require web servers, databases, and workers or messaging, the entire infrastructure can be brought up or down with a combination of templates and orchestration tools like **docker-compose** [60]. **docker-compose** also allows definition of services using multiple containers via its own **docker-compose.yml** file. This file can help to template options including mounted volumes, permissions, environment variables, and exposed ports. A **docker-compose**

configuration also helps to keep one *purpose* per container, which often means one process running in the container, and to combine existing stable building blocks instead of bespoke massive containers for specific purposes.

9. Publish one Dockerfile per project in a code repository with version control

Because a **Dockerfile** is a plain text-based format, it works well with VCS. Including a **Dockerfile** alongside your code and data is an effective way to consistently build your software, to show visitors to the repository how it is built and used, to solicit feedback and collaborate with your peers, and to increase the impact and sustainability of your work (cf. [61]). Online collaboration platforms (e.g., GitHub, GitLab) also make it easy to use CI services to test building your image in an independent build environment. Continuous integration increases stability and trust, and gives the ability to publish images automatically. If your **Dockerfile** includes an interactive user interface, you can also adapt it so that it is ready-to-use as a Binder instance [17], providing an online work environment to any user with a simple click of a link. Furthermore, the commit messages in your version controlled repository preserve a record of all changes to the **Dockerfile**.

While there are exceptions to the rule (cf. [62]), it's generally a simple and clear approach to provide one **Dockerfile** per project. Alternatively, you could build an *image stack* to serve different use cases with a common base image, but this reduces understandability and scatters information across multiple places. We recommend to avoid this at the price of a longer **Dockerfile** and mitigate with consistent form (see Rule 3). If you find that you need to use more than one container to run your workflow, use **docker-compose** and consider if it's possible to use build arguments to flip between states (e.g., development vs. production, see Rule 7) or to separate tools into different repositories. Yet most importantly, you should publish *all* files **COPY**ied into the container, e.g., test data or files for software installation from source (see Rule 5), in the same repository.

10. Use the container daily, rebuild the image weekly, clean up and preserve if need be

Using containers for research workflows does not only require technical understanding, but also an awareness of risks that can be managed efficiently by following a number of good *habits*, which we outline below. While there is no firm rule, if you use a container daily, is good practice to rebuild that container every once or two weeks. At the time of publication of research results it is good practice to save a copy of the image in a public data repository so that readers of the publication can access the resources that produced the published results.

First, use your container every time you work on a project and not just as a final step during publication. If the container is the only platform you use, the confidence in proper documentation of the computing environment can be very high [63]. You should prioritize this usage over others, e.g., non-interactive execution of a full workflow, because it gives you personally the highest value and does not limit your use or others' use of your data and code at all (see Rule 7).

Second, for reproducibility, we can treat containers as transient and disposable, and even intentionally rebuild an image at regular intervals. Ideally, containers that we built years ago should rebuild seamlessly, but this is not necessarily the case, especially with

rapidly changing technology relevant to machine learning and data science. It can almost be guaranteed future changes will interfere with the image build process or even the functionality of the software within the image. The longer that you wait to rebuild the image, the more change happens, the higher the risk that your build will not work, and the harder it becomes to identify and fix the problem. If you are using an interactive container and find that you need to manually install a package or change a parameter, it is best practice to add this dependency to the container and rebuild it right away, but one tends to take shortcuts. Therefore, a habitual deletion of a container and cache-less rebuild of the image (a) increases security due to updating underlying software, (b) helps to reveals issues requiring manual interference, i.e., changes to code or configuration not documented in the **Dockerfile** but perhaps should be, and (c) allows you to more incrementally debug issues. This habit can be supported by using continuous deployment or CI strategies. If a VCS repository with a **Dockerfile** is linked to an automated build, then a CI pipeline can build images and execute a validating run the new image, e.g., using contained demo data (see Rule 5).

In the case of needing setup or configuration for the first two habits, it is good practice to provide a **Makefile** alongside your container, which can capture the specific commands. The effective use of a **Makefile** can help avoiding undocumented steps or manual, extra commands to be run on the local machine. A fully scripted configuration makes it easier for both you and future users, and can increase trust in your workflow.

Third, from time to time you can reduce the system resources occupied by Docker images and their layers or unused containers, volumes and networks by running **docker system prune --all**. After a prune is performed, it follows naturally to rebuild a container for local usage, or to pull it again from a newly built registry image. This habit can be automated with a cron job [64].

[SJE: “sure, okay, I’d like to export an image. Give me a recipe to show me how to do it, e.g. for this paper!”] Fourth, you can export the image to file and deposit it in a public data repository, where it not only becomes citable but also provides a snapshot of the *actual* environment you used at a specific point in time. You should include instructions how to import and run the workflow based on the image archive and add your own image tags for clarity. Depositing the image next to other project files, i.e., data, code, and the used **Dockerfile**, in a public repository makes them likely to be preserved, but is highly unlikely that over time you will be able to recreate it precisely from the accompanying **Dockerfile**. Publishing the image and the contained metadata therein (e.g., the Docker version used) may even allow future science historians to emulate the Docker runtime environment. Applying proper preservation strategies (cf. [61]) can be highly complex, but simply running an image “as-is”, i.e. with the default command and entrypoint (see Rule 7), and observing the output is quite likely to work for many years into the future. If the image does not work anymore, a user can still extract the image contents and explore the files of each layer manually, or if an import still works, with exploration tools like dive [37]. However, if you want to ensure usability and extendability, then you could run import, run, and export an image regularly to make sure the export format still works with the then current version of Docker.

The exported image and a version controlled **Dockerfile** together allow you to freely experiment and continue development of your workflow and keeping the image up to date, e.g., updating versions of pinned dependencies (see Rule 4) and regular image building (see above).

Finally, for a sanity check and to foster even higher trust in the stability and documentation of your project, you can ask a colleague or community member to be your code copilot (see https://twitter.com/Code_Copilot) to interact with your workflow container on a machine of their own. You can do this shortly before submitting your reproducible workflow for peer-review, so you are well positioned for

the future of scholarly communication and open science where these may be standard practices required for publication [22,65–67].

Example Dockerfiles

To demonstrate the ten rules, we maintain a GitHub repository with example `Dockerfiles`, some of which we took from public repositories and updated to adhere to the rules (see `Dockerfile.original`):
<https://github.com/nuest/ten-simple-rules-dockerfiles/>
[SJE: I would like to see these annotated. Perhaps this is something I can meaningfully do?]

Conclusion

Reproducibility in research is an endeavor of incremental improvement and best efforts, not about achieving the perfect solution, which may be not achievable for many researchers with limited resources, and the definition of which may change over time. Containerisation can support researchers. Containers can be optimized for very different use cases, such as supporting many hosts, being small in size for quick deployments in cloud infrastructures, but in this article we have provided guidance for using `Dockerfiles` to create containers for use and communication in computational research. Our goal is to help the researcher to work towards creating a “time capsule” [68] in form of a research compendium. Given some expertise and the right tools, such a compendium can come as close as needed to the original workflow with reasonably little effort for better science. Even if such a capsule decays over time, the effort to create and document it provides incredibly useful and valuable transparency for the project. We encourage researchers to value these steps taken by their peers to use `Dockerfiles` to create “time capsules”, and promote change in the way scholars communicate towards an open and friendly “preproducibility” [69]. So please, make a best effort with your current knowledge, and strive to write readable `Dockerfiles` for functional containers, that are realistic about what might break and what is unlikely to break. In a similar vein, we accept that you should freely break these rules if another way makes more sense *for your use case*. Most importantly, share and exchange your `Dockerfile` freely and collaborate in your community to spread the knowledge about containers as a tool for research and scholarly collaboration and communication. Together we can develop common practices and shared materials for better transparency, higher efficiency, and faster innovation.

Acknowledgements

DN is supported by the project Opening Reproducible Research II (<https://o2r.info/>; <https://www.uni-muenster.de/forschungaz/project/12343>) funded by the German Research Foundation (DFG) under project number PE 1632/17-1. DN and SJE are supported by a Mozilla mini science grant.

Contributions

DN conceived the idea, outlined the first rules, and contributed to all rules. VS wrote the first draft and contributed to all rules. BM revised the text and contributed to all rules. SJE contributed to the overall structure and selected rules. THE contributed to

the rule structure and particularly Rule 1. THi gave extensive feedback on early drafts and contributed to the discussion. This articles was written collaboratively on GitHub, where all contributions in form of text or discussions comments are documented: <https://github.com/nuest/ten-simple-rules-dockerfiles/>.

References

1. Wikipedia contributors. Version control [Internet]. Wikipedia. 2019. Available: https://en.wikipedia.org/w/index.php?title=Version_control&oldid=926593231
2. Marwick B. How computers broke science – and what we can do to fix it [Internet]. The Conversation. 2015. Available: <https://theconversation.com/how-computers-broke-science-and-what-we-can-do-to-fix-it-49938>
3. Donoho DL. An invitation to reproducible computational research. *Biostatistics*. 2010;11: 385–388. doi:10.1093/biostatistics/kxq028
4. Wilson G, Aruliah DA, Brown CT, Hong NPC, Davis M, Guy RT, et al. Best Practices for Scientific Computing. *PLOS Biology*. 2014;12: e1001745. doi:10.1371/journal.pbio.1001745
5. Wilson G, Bryan J, Cranston K, Kitzes J, Nederbragt L, Teal TK. Good enough practices in scientific computing. *PLOS Computational Biology*. 2017;13: e1005510. doi:10.1371/journal.pcbi.1005510
6. Rule A, Birmingham A, Zuniga C, Altintas I, Huang S-C, Knight R, et al. Ten simple rules for writing and sharing computational analyses in Jupyter Notebooks. *PLOS Computational Biology*. 2019;15: e1007007. doi:10.1371/journal.pcbi.1007007
7. Sandve GK, Nekrutenko A, Taylor J, Hovig E. Ten Simple Rules for Reproducible Computational Research. *PLoS Comput Biol*. 2013;9: e1003285. doi:10.1371/journal.pcbi.1003285
8. Nüst D. Author Carpentry : Docker for reproducible research [Internet]. Author Carpentry : Docker for reproducible research. 2017. Available: <https://nuest.github.io/docker-reproducible-research/>
9. Chapman P. Reproducible data science environments with Docker Phil Chapman’s Blog [Internet]. 2018. Available: <https://chapmandu2.github.io/post/2018/05/26/reproducible-data-science-environments-with-docker/>
10. rOpenSci Labs. R Docker tutorial [Internet]. 2015. Available: <https://ropenscilabs.github.io/r-docker-tutorial/>
11. Udemy, Zhibanko V. Docker Containers for Data Science and Reproducible Research [Internet]. Udemy. 2019. Available: <https://www.udemy.com/course/docker-containers-data-science-reproducible-research/>
12. Psomopoulos FE. Lesson ”Docker and Reproducibility” in Workshop ”Reproducible analysis and Research Transparency” [Internet]. Reproducible analysis and Research Transparency. 2017. Available: <https://reproducible-analysis-workshop.readthedocs.io/en/latest/8.Intro-Docker.html>
13. Docker Inc. Best practices for writing Dockerfiles [Internet]. Docker Documentation. 2020. Available: https://docs.docker.com/develop/develop-images/dockerfile_best-practices/
14. Brinckman A, Chard K, Gaffney N, Hategan M, Jones MB, Kowalik K, et al. Computing environments for reproducibility: Capturing the “Whole Tale”. *Future Generation Computer Systems*. 2018; doi:10.1016/j.future.2017.12.029
15. Code Ocean [Internet]. 2019. Available: <https://codeocean.com/>
16. Šimko T, Heinrich L, Hirvonsalo H, Kousidis D, Rodríguez D. REANA: A System for Reusable Research Data Analyses. *EPJ Web of Conferences*. 2019;214:

06034. doi:10.1051/epjconf/201921406034 824

17. Jupyter P, Bussonnier M, Forde J, Freeman J, Granger B, Head T, et al. Binder 825
2.0 - Reproducible, interactive, sharable environments for science at scale. Proceedings 826
of the 17th Python in Science Conference. 2018; 113–120. 827
doi:10.25080/Majora-4af1f417-011 828

18. Nüst D, Konkol M, Pebesma E, Kray C, Schutzzeichel M, Przibytzin H, et al. 829
Opening the Publication Process with Executable Research Compendia. D-Lib 830
Magazine. 2017;23. doi:10.1045/january2017-nuest 831

19. Cohen J, Katz DS, Barker M, Chue Hong NP, Haines R, Jay C. The Four Pillars 832
of Research Software Engineering. IEEE Software. 2020; doi:10.1109/MS.2020.2973362 833

20. Wikipedia contributors. Docker (software) [Internet]. Wikipedia. 2019. 834
Available: [https://en.wikipedia.org/w/index.php?title=Docker_\(software\)&oldid=928441083](https://en.wikipedia.org/w/index.php?title=Docker_(software)&oldid=928441083) 835

21. Boettiger C, Eddelbuettel D. An Introduction to Rocker: Docker Containers for 836
R. The R Journal. 2017;9: 527–536. doi:10.32614/RJ-2017-065 837

22. Chen X, Dallmeier-Tiessen S, Dasler R, Feger S, Fokianos P, Gonzalez JB, et al. 838
Open is not enough. Nature Physics. 2019;15: 113. doi:10.1038/s41567-018-0342-2 839

23. Docker Inc. Dockerfile reference [Internet]. Docker Documentation. 2019. 840
Available: <https://docs.docker.com/engine/reference/builder/> 841

24. Kurtzer GM, Sochat V, Bauer MW. Singularity: Scientific containers for 842
mobility of compute. PLOS ONE. 2017;12: e0177459. doi:10.1371/journal.pone.0177459 843

25. Boettiger C. An Introduction to Docker for Reproducible Research. SIGOPS 844
Oper Syst Rev. 2015;49: 71–79. doi:10.1145/2723872.2723882 845

26. Ben Marwick. 1989-excavation-report-Madjebebe. 2015; 846
doi:10.6084/m9.figshare.1297059 847

27. Docker Inc. Official Images on Docker Hub [Internet]. Docker Documentation. 848
2019. Available: https://docs.docker.com/docker-hub/official_images/ 849

28. Nüst D, Hinz M. Containerit: Generating Dockerfiles for reproducible research 850
with R. Journal of Open Source Software. 2019;4: 1603. doi:10.21105/joss.01603 851

29. Stencila. Stencila/dockta [Internet]. Stencila; 2019. Available: 852
<https://github.com/stencila/dockta> 853

30. Docker Inc. Official Images on Docker Hub [Internet]. Docker Documentation. 854
2020. Available: https://docs.docker.com/docker-hub/official_images/ 855

31. Husain H, Silkaitis R. Machine-learning-apps/repo2docker-action [Internet]. ML 856
Apps; 2019. Available: 857
<https://github.com/machine-learning-apps/repo2docker-action> 858

32. Henderson S. ScottyHQ/repo2docker-githubci [Internet]. 2019. Available: 859
<https://github.com/scottyhq/repo2docker-githubci> 860

33. Preston-Werner T. Semantic Versioning 2.0.0 [Internet]. Semantic Versioning. 861
2013. Available: <https://semver.org/> 862

34. Docker Inc. Use multi-stage builds [Internet]. Docker Documentation. 2020. 863
Available: 864
<https://docs.docker.com/develop/develop-images/multistage-build/> 865

35. Halchenko YO, Hanke M. Open is Not Enough. Let's Take the Next Step: An 866
Integrated, Community-Driven Computing Platform for Neuroscience. Frontiers in 867
Neuroinformatics. 2012;6. doi:10.3389/fninf.2012.00022 868

36. Wikipedia contributors. Lint (software) [Internet]. Wikipedia. 2019. Available: 869
[https://en.wikipedia.org/w/index.php?title=Lint_\(software\)&oldid=907589761](https://en.wikipedia.org/w/index.php?title=Lint_(software)&oldid=907589761) 870

37. Goodman A. Wagoodman/dive [Internet]. 2019. Available: 871
<https://github.com/wagoodman/dive> 872

38. The Python Software Foundation. Requirements Files — pip User Guide [Internet]. 2019. Available: https://pip.pypa.io/en/stable/user_guide/#requirements-files
39. Continuum Analytics. Managing environments — conda documentation [Internet]. 2017. Available: <https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html>
40. R Core Team. The DESCRIPTION file in "writing r extensions" [Internet]. 1999. Available: <https://cran.r-project.org/doc/manuals/r-release/R-exts.html#The-DESCRIPTION-file>
41. Eddelbuettel D, Horner J. Littler: R at the command-line via 'r' [Internet]. 2019. Available: <https://CRAN.R-project.org/package=littler>
42. npm. Creating a package.json file npm Documentation [Internet]. 2019. Available: <https://docs.npmjs.com/creating-a-package-json-file>
43. The Julia Language Contributors. 10. Project.Toml and Manifest.Toml · Pkg.Jl [Internet]. 2019. Available: <https://julialang.github.io/Pkg.jl/v1/toml-files/>
44. Docker Inc. Use bind mounts [Internet]. Docker Documentation. 2019. Available: <https://docs.docker.com/storage/bind-mounts/>
45. Opencontainers. Opencontainers/image-spec v1.0.1 - Annotations [Internet]. GitHub. 2017. Available: <https://github.com/opencontainers/image-spec/blob/v1.0.1/annotations.md>
46. Katz DS, Chue Hong NP. Software Citation in Theory and Practice. In: Davenport JH, Kauers M, Labahn G, Urban J, editors. Mathematical Software – ICMS 2018. Springer International Publishing; 2018. pp. 289–296. doi:10.1007/978-3-319-96418-8_34
47. Verstegen JA. JudithVerstegen/PLUC_Mozambique: First release of PLUC for Mozambique [Internet]. Zenodo; 2019. doi:10.5281/zenodo.3519987
48. Sochat V. The Scientific Filesystem. GigaScience. 2018;7. doi:10.1093/gigascience/giy023
49. Wikipedia contributors. Make (software) [Internet]. Wikipedia. 2019. Available: [https://en.wikipedia.org/w/index.php?title=Make_\(software\)&oldid=929976465](https://en.wikipedia.org/w/index.php?title=Make_(software)&oldid=929976465)
50. Knoth C, Nüst D. Reproducibility and Practical Adoption of GEOBIA with Open-Source Software in Docker Containers. Remote Sensing. 2017;9: 290. doi:10.3390/rs9030290
51. Molenaar G, Makhathini S, Girard JN, Smirnov O. Klico—The scientific compute container format. Astronomy and Computing. 2018;25: 1–9. doi:10.1016/j.ascom.2018.08.003
52. Docker Inc. Docker/kitematic [Internet]. Docker; 2019. Available: <https://github.com/docker/kitematic>
53. Selenium contributors. SeleniumHQ/selenium [Internet]. Selenium; 2019. Available: <https://github.com/SeleniumHQ/selenium>
54. Singularity. Frequently Asked Questions Singularity [Internet]. 2019. Available: <http://singularity.lbl.gov/archive/docs/v2-2/faq#can-i-run-x11-apps-through-singularity>
55. Viereck M. X11docker: Run GUI applications in Docker containers. Journal of Open Source Software. 2019;4: 1349. doi:10.21105/joss.01349
56. Yaremenko E. JAreMko/docker-x11-bridge [Internet]. 2019. Available: <https://github.com/JAreMko/docker-x11-bridge>
57. Panda Y. Yuvipanda/jupyter-desktop-server [Internet]. 2019. Available: <https://github.com/yuvipanda/jupyter-desktop-server>
58. {Cookiecutter contributors}. Cookiecutter/cookiecutter [Internet]. cookiecutter; 2019. Available: <https://github.com/cookiecutter/cookiecutter>

59. Marwick B. Benmarwick/rrtools [Internet]. 2019. Available: <https://github.com/benmarwick/rrtools>
60. Docker Inc. Overview of Docker Compose [Internet]. Docker Documentation. 2019. Available: <https://docs.docker.com/compose/>
61. Emsley I, De Roure D. A Framework for the Preservation of a Docker Container International Journal of Digital Curation. International Journal of Digital Curation. 2018;12. doi:10.2218/ijdc.v12i2.509
62. Kim B, Ali TA, Lijeron C, Afgan E, Krampis K. Bio-Docklets: Virtualization Containers for Single-Step Execution of NGS Pipelines. bioRxiv. 2017; 116962. doi:10.1101/116962
63. Marwick B. README of 1989-excavation-report-Madjebebe. 2015; doi:10.6084/m9.figshare.1297059
64. Wikipedia contributors. Cron [Internet]. Wikipedia. 2019. Available: <https://en.wikipedia.org/w/index.php?title=Cron&oldid=929379536>
65. Eglen S, Nüst D. CODECHECK: An open-science initiative to facilitate sharing of computer programs and results presented in scientific publications. Septentrio Conference Series. 2019; doi:10.7557/5.4910
66. Schönbrodt F. Training students for the Open Science future. Nature Human Behaviour. 2019;3: 1031–1031. doi:10.1038/s41562-019-0726-z
67. Eglen SJ, Mounce R, Gatto L, Currie AM, Nobis Y. Recent developments in scholarly publishing to improve research practices in the life sciences. Emerging Topics in Life Sciences. 2018;2: 775–778. doi:10.1042/ETLS20180172
68. Blank D, Twitter. Twitter thread on reproducibility time capsules on Twitter [Internet]. Twitter. 2019. Available: <https://twitter.com/dougblank/status/1135904909663068165>
69. Stark PB. Before reproducibility must come preproducibility [Internet]. Nature. 2018. doi:10.1038/d41586-018-05256-0