

An introduction to Docker for reproducible research

Carl Boettiger

Center for Stock Assessment Research,
110 Shaffer Rd, Santa Cruz, CA 95050, USA
cboettig(at)gmail.com

ABSTRACT

As computational work becomes more and more integral to many aspects of scientific research, computational reproducibility has become an issue of increasing importance to computer systems researchers and domain scientists alike. Though computational reproducibility seems more straightforward than replicating physical experiments, the complex and rapidly changing nature of computer environments makes being able to reproduce and extend such work a serious challenge. In this paper, I explore common reasons that code developed for one research project cannot be successfully executed or extended by subsequent researchers. I review current approaches to these issues, including virtual machines and workflow systems, and their limitations. I then examine how the popular emerging technology Docker combines several areas from systems research - such as operating system virtualization, cross-platform portability, modular re-usable elements, versioning, and a ‘DevOps’ philosophy, to address these challenges. I illustrate this with several examples of Docker use with a focus on the R statistical environment.

INTRODUCTION

Reproducible research has received an increasing level of attention throughout the scientific community [19, 22] and the public at large [25]. All steps of the scientific process, from data collection and processing, to analyses, visualizations and conclusions depend ever more on computation and algorithms, *computational reproducibility* has received particular attention [18]. Though in principle this algorithmic dependence should make such research easier to reproduce – computer codes being both more portable and potentially more precise to exchange and run than experimental methods – in practice this has led to an ever larger and more complex black box that stands between what was actually done and what is described in the literature. Crucial scientific processes such as replicating the results, extending the approach or testing the conclusions in other contexts, or even merely installing the software used by the original researchers can become immensely time-consuming if not impossible.

Systems research & reproducibility

Systems research has long concerned itself with the issues of computational reproducibility and the technologies that can facilitate those objectives [6]. Docker is a new but already very popular open source tool that combines many of these approaches in a user friendly implementation, including: (1) performing Linux container (LXC) based operating system (OS) level virtualization, (2) portable deployment of containers across platforms, (3) component reuse, (4) sharing, (5) archiving, and (6) versioning of container images. While Docker’s market success has largely focused on the needs of businesses in deploying web applications and the potential for a lightweight alternative to full virtualization, these features have potentially important implications for systems research in the area of scientific reproducibility.

In this paper, I seek to address two audiences. First, that of the domain scientist, those conducting research in ecology, bioinformatics, economics, psychology and so many other disciplines in which computation plays an ever-increasing role. I seek to help this audience become more aware of the concerns and challenges in making these computations more reproducible, extensible, and portable to other researchers, and highlight Docker as platform to address these challenges. The second audience is that of the computer systems researchers: readers familiar with both the challenges to reproducibility (versioning, software dependencies, etc.) and the technical components underlying Docker (virtualization, LXC containers, jails, hashes, etc.), but who may be unfamiliar with Docker software. I hope the domain scientist will be motivated to try Docker as a tool to address reproducibility. Meanwhile, I expect the systems researcher to see Docker technology as an area ripe for further systems research (e.g. to what extent can differences between Linux kernels frustrate reproducibility in a Docker image?) rather than as a tool for their own study (because Docker applies only to software at the level above the Linux kernel, it will be of little use in making studies of differences in kernel behavior, hardware performance, execution times, and so forth reproducible.)

A cultural problem

It is worth observing from the outset that the primary barrier to computational reproducibility in many domain sciences has nothing to do with the technological approaches discussed here, but stems rather from a reluctance to publish the code used in generating the results in the first place [2]. Despite extensive evidence to the contrary [14], many researchers and journals continue to assume that summary descriptions or

pseudo-code provide a sufficient description of computational methods used in data gathering, processing, simulation, visualization, or analysis. Until such code is available in the first place, one cannot even begin to encounter the problems that the approaches discussed here set out to solve. As a result, few domain researchers may be fully aware of the challenges involved in effectively re-using published code.

A lack of requirements or incentives no doubt plays a crucial role in discouraging sharing [2, 24]. Nevertheless, it is easy to underestimate the significant barriers raised by a lack of familiar, intuitive, and widely adopted tools for addressing the challenges of computational reproducibility. Surveys and case studies find that a lack of time, more than innate opposition to sharing, discourages researchers from providing code [7, 23].

Four technical challenges

By restricting ourselves to studies of where code has been made available, I will sidestep for the moment the cultural challenges to reproducibility so that I may focus on the technical ones; in particular, those challenges for which improved tools and techniques rather than merely norms of behavior can contribute substantially to improved reproducibility.

Studies focusing on code that has been made available with scientific publications regularly find the same common issues that pose substantial barriers to reproducing the original results or building on that code [4, 8, 10, 16, 27], which I attempt to summarize as follows.

1. “Dependency Hell”

A recent study by researchers at the University of Arizona found that less than 50% of software could even be successfully built or installed [4]. Though sufficiently knowledgeable users may be able to overcome the issues in at least some of these cases,¹ similar results are seen in an ongoing effort by other researchers to replicate that study [27], and has also been observed in other independent studies [8, 16]. Installing or building software necessary to run the code in question assumes the ability to recreate the computational environment of the original researchers.

Differences in numerical evaluation, such as arise in floating point arithmetic or even ambiguities in standardized programming languages (“order-of-evaluation” problems) can be responsible for differing results between or even within the same computational platform [14]. Such issues make it difficult to restrict the true dependencies of the code to higher level environments such as that of a given scripting language, independent of the underlying OS or even hardware itself.

2. Imprecise documentation

Documentation on how to install and run code associated with published research is another frequent barrier to replication. A study by Lapp [16] found this impairs a researcher’s ability to install and build the software necessary, as even small holes in the documentation were found to be major barriers, particularly for “novices” [8] – where novices may be experts in nearby languages but unfamiliar with the package

¹See discussion at <https://gist.github.com/pbailis/9629050> and <https://gist.github.com/samth/9641364>

managers and other tools of the language involved. This same problem is discussed in [3]. Imprecise documentation goes well beyond issues of the software environment itself: incomplete documentation of parameters involved meant as few as 30% of analyses ($n = 34$) using the popular software STRUCTURE could be reproduced in the study of [10].

3. Code rot

Software dependencies are not static elements, but receive regular updates that may fix bugs, add new features or deprecate old features (or even entire dependencies themselves). Any of these changes can potentially change the results generated by the code. As some of these changes may indeed resolve valid bugs or earlier problems with underlying code, it will often be insufficient to demonstrate that results can be reproduced when using the original versions, a problem sometimes known as “code rot.” Researchers will want to know if the results are robust to the changes. The case studies in [16] provide examples of these problems.

4. Barriers to adoption and reuse in existing solutions

Technological solutions such as workflow software, virtual machines, continuous integration services, and best practices from software development would address many of the issues frequently frustrating reproducibility. However, researchers face significant barriers to entry in learning these tools and approaches which are not part of their typical curriculum, or lack incentives commensurate with the effort required [7, 15].

Though a wide variety of approaches exists to work around these challenges, few operate on a low enough level to provide a general solution. Clark *et al.* [3] provide an excellent description of this situation:

In scientific computing the environment was commonly managed via Makefiles & Unix-y hacks, or alternatively with monolithic software like Matlab. More recently, centralized package management has provided curated tools that work well together. But as more and more essential functionality is built out across a variety of systems and languages, the value – and also the difficulty – of coordinating multiple tools continues to increase. Whether we are producing research results or web services, it is becoming increasingly essential to set up new languages, libraries, databases, and more.

There are two dominant approaches to this issue of coordinating multiple tools: Workflows and Virtual Machines (VMs).

CURRENT APPROACHES

Two dominant paradigms have emerged to address these issues so far: workflow software [1, 13] and virtual machines [5, 12]. Workflow software provides very elegant technical solutions to the challenges of communication between diverse software tools, capturing provenance in graphically

driven interfaces, and handling issues from versioning dependencies to data access. Workflow solutions are often built by well-funded collaborations between domain scientists and computer scientists, and can be very successful in the communities within which they receive substantial adoption. Nonetheless, most workflow systems struggle with relatively low total adoption overall [5, 9].

Dudley & Butte [5] give several reasons that such comprehensive workflow systems have not been more successful:

- (i) efforts are not rewarded by the current academic research and funding environment;
- (ii) commercial software vendors tend to protect their markets through proprietary formats and interfaces;
- (iii) investigators naturally tend to want to ‘own’ and control their research tools;
- (iv) even the most generalized software will not be able to meet the needs of every researcher in a field; and finally
- (v) the need to derive and publish results as quickly as possible precludes the often slower standards-based development path.

In short, workflow software expects a new approach to computational research. In contrast, virtual machines (VMs) offer a more direct approach. Since the computer Operating System (OS) already provides the software layer responsible for coordinating all the different elements running on the computer, the VM approach captures the OS and everything running on it whole-cloth. To make this practical, Dudley & Butte [5] and Howe [12] both propose using virtual machine images that will run on the cloud, such as Amazon’s EC2 system, which is already based upon this kind of virtualization.

Critics of the use of VMs to support reproducibility highlight that the approach is too much of a black box and thus ill suited for reproducibility [28]. While the approach sidesteps the need to either install or even document the dependencies, this also makes it more difficult for other researchers to understand, evaluate, or alter those dependencies. Moreover, other research cannot easily build on the virtual machine in a consistent and scalable way. If each study provided its own virtual machine, any pipeline combining the tools of multiple studies would quickly become impractical or impossible to implement.

A “DevOps” approach

The problems highlighted here are not unique to *academic* software, but impact software development in general. While the academic research literature has frequently focused on the development of workflow software dedicated to particular domains, or otherwise to the use of virtual machines, the software development community has recently emphasized a philosophy (rather than a particular tool), known as *Development* and *Systems Operation*, or more frequently just “DevOps.” The approach is characterized by scripting, rather than documenting, a description of the necessary dependencies for software to run, usually from the Operating System (OS) on up. Clark *et al.* [3] describe the DevOps approach along with both its relevance to reproducible research and

examples of its use in the academic research context. They identify the difficulties I have discussed so far in terms of effective documentation:

Documentation for complex software environments is stuck between two opposing demands. To make things easier on novice users, documentation must explain details relevant to factors like different operating systems. Alternatively, to save time writing and updating documentation, developers like to abstract over such details.

The authors contrast this to the DevOps approach, where dependency documentation is *scripted*:

A DevOps approach to “documenting” an application might consist of providing brief descriptions of various install paths, along with scripts or “recipes” that automate setup.

This elegantly addresses both the demand for simplicity of use (one executes a script instead of manually managing the environmental setup) and comprehensiveness of implementation. Clark *et al.* [3] are careful to note that this is not so much a technological shift as a philosophical one:

The primary shift that’s required is not one of new tooling, as most developers already have the basic tooling they need. Rather, the needed shift is one of philosophy.

Nevertheless, a growing suite of tools designed explicitly for this purpose have rapidly replaced the use of general purpose tools (such as Makefiles, bash scripts) to become synonymous with the DevOps philosophy. Clark *et al.* [3] reviews many of these DevOps tools, their different roles, and their application in reproducible research.

I focus the remainder of this paper on one of the most recent and rapidly growing among these, called Docker, and the role it can play in reproducible research. Docker offers several promising features for reproducibility that go beyond the tools highlighted in [3]. Nevertheless, my goal in focusing on this technology is not to promote a particular solution, but to anchor the discussion of technical solutions to reproducibility challenges in concrete examples.

DOCKER

Docker is an open source project that builds on many long-familiar technologies from operating systems research: LXC containers, virtualization of the OS, and a hash-based or git-like versioning and differencing system, among others (see docs.docker.com/faq for an excellent overview of what Docker adds to plain LXC). The official documentation docs.docker.com already provides a thorough introduction in how to use Docker software; here my focus is on describing the implications this has for reproducible research. Readers may also find it helpful to see more detailed examples of

using Docker to capture, share, and interact with a specific computational environments. Some such examples, along with more detailed documentation on use can be found at github.com/rocker-org.

I introduce the most relevant concepts from Docker through the context of the four challenges for reproducible research I have discussed above. In brief, these challenges can be addressed by distributing a Dockerfile capable of reconstructing the researcher's development environment, ideally along with depositing the corresponding binary Docker image in an appropriate repository.

1. Docker images: resolving 'Dependency Hell'

A Docker based approach works similarly to a virtual machine image in addressing the dependency problem by providing other researchers with a binary image in which all the software has already been installed, configured and tested. (A machine image can also include all data files necessary for the research, which may simplify the distribution of data.)

A key difference between Docker images and other virtual machines is that the Docker images share the Linux kernel with the host machine. For the end user the primary consequence of this is that any Docker image must be based on a Linux system with Linux-compatible software, which includes R, Python, Matlab, and most other scientific programming needs.²

Sharing the Linux kernel makes Docker more light-weight and higher performing than complete virtual machines – a typical desktop computer could run no more than a few virtual machines at once but would have no trouble running 100's of Docker containers (a container is simply the term for running instance of an image). This feature has made Docker particularly attractive to industry and is largely responsible for the immense popularity of Docker. For our purposes this is a nice bonus, but the chief value to reproducible research lies in other aspects.

2. Dockerfiles: Resolving imprecise documentation

Though Docker images can be created interactively, this leaves little transparent record³ of what software has been installed and how. Dockerfiles provide a simple script (similar to a Makefile) that defines exactly how to build up the image, consistent with the DevOps approach I mentioned previously.

With a syntax that is simpler than other provisioning tools (*e.g.* Chef, Puppet, Ansible) or Continuous Integration (CI) platforms (*e.g.* Travis CI, Shippable CI); users need little more than a basic familiarity with shell scripts and a Linux distribution software environment (*e.g.* Debian-based `apt-get`) to get started writing Dockerfiles.

This approach has many advantages:

²Note that re-distribution of an image in which proprietary software has been installed will be subject to any relevant licensing agreement.

³The situation is in fact slightly better than the virtual machine approach because these changes are versioned. Docker provides tools to inspect differences (diffs) between the images, and I can also roll back changes to earlier versions.

- While machine images can be very large (many gigabytes), a Dockerfile is just a small plain text file that can be easily stored and shared.
- Small plain text files are ideally suited for use with a version management system such as `subversion` or `git`, which can track any changes made to the `Dockerfile`
- the `Dockerfile` provides a human readable summary of the necessary software dependencies, environmental variables and so forth needed to execute the code. There is little possibility of the kind of holes or imprecision in such a script that so frequently cause difficulty in manually implemented documentation of dependencies. This approach also avoids the burden of having to tediously document dependencies at the end of a project, since they are instead documented as they are installed by writing the `Dockerfile`.
- Unlike a `Makefile` or other script, the `Dockerfile` includes all software dependencies down to the level of the OS, and is built by the Docker `build` tool, making it very unlikely that the resulting build will differ when being built on different machines. This is not to say that all builds of a `Dockerfile` are bitwise identical. In particular, builds executed later will install more recent versions of the same software, if available, unless the package managers used are explicitly configured otherwise. I address this issue in the next section.
- It is possible to add checks and tests following the commands for installing the software environment, which will verify that the setup has been successful. This can be important in addressing the issue of code-rot which I discuss next.
- It is straightforward for other users to extend or customize the resulting image by editing the script directly.

3. Tackling code-rot with image versions

As I have discussed above, changes to the dependencies, whether they are the result of security fixes, new features, or deprecation of old software, can break otherwise functioning code. These challenges can be significantly reduced because Docker defines the software environment to a particular operating system and suite of libraries, such as the Ubuntu or Debian distribution. Such distributions use a staged release model with `stable`, `testing` and `unstable` phases subjected to extensive testing to catch such potential problems [20], while also providing regular security updates to software within each stage. Nonetheless, this cannot completely avoid the challenge of code-rot, particularly when it is necessary to install software that is not (yet) available for a given distribution.

To address this concern, one must archive a binary copy of the image used at the time the research was first performed. Docker provides a simple utility to save an image as a portable `tarball` file that can be read in by any other Docker installation, providing a robust way to run the exact versions of all software involved. By testing both the `tarball` archive and the image generated by the latest `Dockerfile`, Docker provides a simple way to confirm whether or not code rot has effected the function of a particular piece of code.

Binary Docker images can be efficiently shared through the Docker Hub as well, as I describe later.

4. *Barriers to adoption and re-use*

A technical solution, no matter how elegant, will be of little practical use for reproducible research unless it is both easy to use and adapt to the existing workflow patterns of practicing domain researchers.

Though most of the concerns I have discussed so far can be addressed through well-designed workflow software or the use of a DevOps approach to provisioning virtual machines by scripts, neither approach has seen widespread adoption by domain researchers, who work primarily in a local rather than cloud-based environment using development tools native to their personal operating system. To gain more widespread adoption, reproducible research technologies must make it easier, not harder, for a researcher to perform the tasks they are already doing (before considering any additional added benefits).

These issues are reflected both during the original research or development phase and in any subsequent reuse. Another researcher may be less likely to build on existing work if it can only be done by using a particular workflow system or monolithic software platform with which they are unfamiliar. Likewise, a user is more likely to make their own computational environment available for reuse if it does not involve a significant added effort in packaging and documenting [23].

Though Docker is not immune to these challenges, it offers an interesting example of a way forward in addressing these fundamental concerns. Here I highlight these features in turn:

- Integrating into local development environments
- Modular reuse
- Portable environments
- Public repository for sharing
- Versioning

Integrating into local development environments

Due in part to the difficulty of moving large VM images around a network, proponents of virtual machines for reproducible research often propose that these machines would be available exclusively as cloud computing environments [5] rather than downloaded to a user's personal computer. Though cloud computing offers some advantages such as scalable resources for computationally intensive tasks, many researchers prefer to work with tools installed locally on their personal computer, at least during testing and development. This can reduce costs, issues of network latency, the ability to work off-line, and will be most familiar to students and new developers.

It is possible to run virtual machines locally on most common laptop and desktop computers, as demonstrated in the approach now being pioneered at UC Berkeley [3]. This approach provides users with a pixel-identical environment whether working locally or on the cloud, which has particular advantages for student instruction [3]. However, it

remains to be seen if existing researchers will be willing to forgo native applications for such tasks as file browsing, text editing, or version management and rely exclusively on this standardized virtual environment.

In contrast, Docker's approach is optimized for a more integrated workflow. Rather than replace a user's existing toolchain with a standardized virtual environment, editors and all, Docker is optimized at the level of single applications. A developer can thus rely on familiar tools while still ensuring that the execution of their code always occurs on the standardized container environment, thus ensuring its portability and reproducibility. This approach can be accomplished either by linking volumes or directories between the container and the host, or simply by instructing the Dockerfile to copy the code to the container (As we will see, modular reuse makes this latter strategy efficient).

Rather than put the burden on the researcher to adopt a very different workflow, the researcher can thus use their familiar editors, etc., while immediately benefiting from the fact their computations can be easily deployed on the cloud or the machines of collaborators without any further effort than installing Docker software. The docker approach is particularly well suited for moving between local and cloud platforms when a web-based integrated development environment is available, such as RStudio Server, or (to lesser extent) an iPython notebook.

On systems not already based on the Linux kernel (such as Mac or Windows platforms), Docker is installed (see docs.docker.com/installation) through means of a very small (about 24 MB) VM platform called `boot2docker` (github.com/docker/boot2docker). While this poses an additional challenge for tight integration with desktop tools, recent advances in Docker are rapidly bridging this gap as well. For example, Docker 1.3, released between drafts of this manuscript, supports shared volumes on Macs (blog.docker.com).

Portable computation & sharing

A particular advantage of this approach is that the resulting computational environment is immediately *portable*. LXC containers by themselves are unlikely to run in the same way, if at all, across different machines, due to differences in networking, storage, logging and so forth. Docker handles the packaging and execution of a container so that it works identically across different machines, while exposing the necessary interfaces for networking ports, volumes, and so forth. This is useful not only for the purposes of reproducible research, where other users may seek to reconstruct the computational environment necessary to run the code, but is also of immediate value to the researcher themselves. For instance, a researcher might want to execute their code on a cloud server which has more memory or processing power than their local machine, or would want a co-author to help debug a particular problem. In either case, the researcher can export a snapshot of their running container:

```
docker export container-name > container.tar
```

and then run this identical environment on the cloud or collaborators' machine.

Sharing these images is further facilitated by the Docker Hub technology (hub.docker.com). While Docker images tend to be much smaller than equivalent virtual machines, moving around even 100's of gigabytes can be a challenge. The Docker Hub provides a convenient distribution service, freely storing the pre-built images, along with their metadata, for download and reuse by others. The Docker Hub is a free service and an open source software product so that users can run their own private versions of the Hub on their own servers, for instance, if security of the data or the longevity of the public platform is a concern. Docker also supports *Automated Builds* through the Docker Hub. This acts as a kind of Continuous Integration (CI) service that verifies the image builds correctly whenever the Dockerfile is updated, particularly if the Dockerfile includes checks for the environment.

One can share a public copy of the image just created by using the `docker push` command, followed by the name of the image using the command:

```
docker push username/r-recommended
```

If a Dockerfile is made available on a public code repository such as Github or Bitbucket, the Hub can automatically build the image whenever a change is made to the Dockerfile, making the `push` command unnecessary. A user can update their local image using the `docker pull <imagename>`, which downloads any changes that have since been made to the copy of the image on the Hub.

Re-usable modules

The approach of Docker offers a technical solution to what is frequently seen as the primary weakness of the standard VM approach to reproducibility - reusing and remixing elements. To some extent this is already addressed by the DevOps approach of Dockerfiles, providing a scripted description of the environment that can be tweaked and altered, but also includes something much more fundamental to Docker.

The challenge to reusing VMs can be summarized as “you can’t install an image for every pipeline you want...” [28]. While providing a VM may make it easy for other researchers to run a particular piece of software, it becomes very difficult to combine multiple software components in future research if each must run inside it’s own VM. The lack of reusable, scalable modules in the VM model poses a major barrier to future reuse. In contrast, (as the analogy to shipping containers in the name might imply) Docker containers are optimized for this kind of stacking and modular reuse. There are at least three ways in which Docker supports this kind of extensibility.

Most primitively, because Dockerfile itself provides a script for creating the computational environment, future researchers can extend or modify the resulting machine image by simply editing the Dockerfile. This same approach is also possible for VM approaches that rely on DevOps tools [3]. In contrast to VMs, however, Docker is also modular by design: both in how Dockerfiles are defined and how Docker containers can be linked.

First, Docker facilitates modular reuse by building one container on top of another through the use of `FROM` directive in Dockerfiles. This acts like a software dependency; but unlike other software, a Dockerfile must have exactly one dependency (one `FROM` line). Particular version of the dependency can be specified using the `:` notation, or omitted to default to the latest version. To install software, Dockerfiles leverage existing package managers on common Linux distributions such as `apt` on Debian. These also permit installing either specific or only the latest version.

While a Dockerfile can have only a single `FROM` line, sometimes it may be necessary to build on the computational environment provided by more than one container. To address this, Docker defines a syntax for linking multiple containers together. This allows each container to act as a building block providing just what is necessary to run one particular service or element, and exposing just what is needed to link it together with other blocks. For instance, one could have one container running a PostgreSQL database which serves data to another container running a python environment to analyze the data:

```
docker run -d --name db training/postgres
docker run -d -P --link db:db training/webapp \
python app.py
```

This separates the task for running the database (first line) from the application used to analyze the data (second line), allowing a more modular approach to reuse: another researcher could use the same database container while connecting it to different scripts.

Unlike the much more heavyweight virtual machine approach, a single computer can easily run 100s of such services each in their own container. A rapidly expanding ecosystem of software around Docker, such as `fig` (github.com/docker/fig) facilitates the complexity of running multiple containers. This feature making it easy to break computational elements down into logically reusable chunks that come, batteries included, with everything they need to run reproducibly.

Versioning

In addition to version managing the Dockerfile, the images themselves are versioned using a git-like hash system (*e.g.* see `docker commit`, `docker push/docker pull`, `docker history`, `docker diff`). Docker images and containers have dedicated metadata specifying the date, author, parent image, and other details (see `docker inspect`). One can roll back an image through the layers of history of its construction, then build off an earlier layer, or roll back changes made interactively in a container. For instance, here I inspect recent changes made to the `ubuntu:14.04` image:

```
docker history ubuntu:14.04
```

One can identify an earlier version, and roll back to that version just by adjusting the Docker tag to match the hash of that version. For instance:

```
docker tag 25f ubuntu:14.04
```

If one now inspects the history, which shows that it now begins from this earlier point:

```
docker history ubuntu:14.04
```

This same feature also means that Docker can perform incremental uploads and downloads that send only the differences between images, (just like `git push` or `git pull` for git repositories), rather than transfer the full image each time.

CONCLUSIONS

Best Practices

The effectiveness of this approach for supporting reproducible research nonetheless depends on how each of these features are adopted and implemented by individual researchers. I summarize a few of these practices here:

- *Use Docker containers during development.* A key feature of the Docker approach is the ability to mimic as closely as possible the current workflow and development practices of the user. Code executing inside a container on a local machine can appear identical to code running natively, but with the added benefit that one can simply recreate or snapshot and share the entire computational environment with a few simple commands. This works best if researchers set up their computational environment in a container from the outset of the project.
- *Write Dockerfiles instead of installing interactive sessions.* As we have noted already, Docker can be used in a purely interactive manner to record and distribute changes to a computational environment. However, the approach is most useful for reproducible research when researchers begin by defining their environment explicitly in the DevOps fashion by writing a Dockerfile.
- *Adding tests or checks to the Dockerfile.* Dockerfile commands need not be limited to installing software, but can also include execution. This can help verify that an image has build successfully with all the software necessary to run the research code of interest.
- *Use and provide appropriate base images.* Though Docker supports modular design, it remains up to the researchers to take advantage of it. An appropriate workflow might involve one Dockerfile that includes all the software dependencies a researcher usually uses in the course of their development, which can then be extended by separate Docker images for particular projects. Re-using existing images reduces the effort required to set up an environment, contributes to the standardization of computational environments within a field, and best leverages the ability of Docker's distribution system to download only differences.
- *Share Docker images and Dockerfiles.* The Docker Hub significantly reduces the barriers for making even large images (which can exceed the file size limits of journals common scientific data repositories such as Dryad and Figshare) readily available to other researchers.

- *Archive tarball snapshots.* Despite similar semantics to git, Docker's versioning system works rather differently than version management of code. Docker can roll back layers⁴ that have been added to an image, but not revert to the earlier state of a particular layer. In consequence, to preserve a bitwise identical snapshot of a container used to generate a given set of results, it is necessary to archive the image tarball itself – one can not simply rely on the Docker history to recover an earlier state.

Limitations and future developments

Docker has the potential to address shortcomings of certain existing approaches to reproducible research challenges that stem from recreating complex computational environments. Docker also provides a promising case study in other issues. Its versioning, modular design, portable containers, and simple interface have proven successful in industry and could have promising implications for reproducible research in scientific communities. Nonetheless, these advances raise questions and challenges of their own.

- Docker does not provide complete virtualization but relies on the Linux kernel provided by the host. Systems research can provide insight on what limitations to reproducibility this introduces [11].
- Docker is limited to 64 bit host machines, making it impossible to run on older hardware (at this time).
- On Mac and Windows machines Docker must still be run in a fully virtualized environment. Though the boot2docker tool streamlines this process, it remains to be seen if the performance and integration with the host machine's OS is sufficiently seamless or creates a barrier to adoption by users on of these systems.
- Potential computer security issues may still need to be evaluated. Among other changes, future support for digitally signing Docker images may make it easier to build off of only trusted binaries.
- Most importantly, it remains to be seen if Docker will be significantly adopted by any scientific research or teaching community.

Further considerations

Combining virtualization with other reproducible-research tools

Using Docker containers to distribute reproducible research should be seen as an approach that is synergistic with, rather than an alternative to, other technical tools for ensuring computational reproducibility. Existing tools for managing dependencies for a particular language [21] can easily be employed within a Docker-based approach, allowing the operating-systems level virtualization to sidestep potential issues such as external library dependencies or conflicts with existing user libraries. Other approaches that facilitate reproducible research also introduce additional software dependencies and possible points of failure [7]. One example includes

⁴Technically AUFS (advanced multi layered unification filesystem) layers, see wikipedia.org/wiki/aufs

dynamic documents [17, 22, 26] which embed the code required to re-generate the results within the manuscript. As a result, it is necessary to package the appropriate typesetting libraries (e.g. L^AT_EX) along with the code libraries such that the document executes successfully for different researchers and platforms.

Impacting cultural norms?

I noted at the outset that cultural expectations responsible for a lack of code sharing practices in many fields are a far more extensive primary barrier to reproducibility than the technical barriers discussed here. Nevertheless, it may be worth considering how solutions to these technical barriers can influence the cultural landscape as well. Many researchers may be reluctant to publish code today because they fear a it will be primarily a one-way street: more technical savvy researchers than themselves can benefit from their hard work, while they may not benefit from the work produced by others. Lowering the technical barriers to reuse provides immediate practical benefits that make this exchange into a more balanced, two-way street. Another concern is that the difficulty imposed in preparing code to be shared, such as providing even semi-adequate documentation or support for other users to be able to install and run it in the first place is too high [23]. Thus, lowering these barriers to re-use through the appropriate infrastructure may also reduce certain cultural barriers to sharing.

ACKNOWLEDGEMENTS

CB acknowledges support from NSF grant DBI-1306697, and also from the Sloan Foundation support through the rOpenSci project. CB also wishes to thank Scott Chamberlain, Dirk Eddelbuettel, Rich FitzJohn, Yihui Xie, Titus Brown, John Stanton-Geddes and many others for helpful discussions about reproducibility, virtualization, and Docker that have helped shape this manuscript.

REFERENCES

[1]Altintas, I. et al. 2004. Kepler: an extensible system for design and execution of scientific workflows. *Proceedings. 16th international conference on scientific and statistical database management, 2004.* (2004).

[2]Barnes, N. 2010. Publish your computer code: it is good enough. *Nature.* 467, 7317 (Oct. 2010), 753–753.

[3]Clark, D. et al. 2014. BCE: Berkeley’s Common Scientific Compute Environment for Research and Education. *Proceedings of the 13th Python in Science Conference (SciPy 2014).* (2014).

[4]Collberg, C. et al. 2014. *Measuring Reproducibility in Computer Systems Research.*

[5]Dudley, J.T. and Butte, A.J. 2010. In silico research in the era of cloud computing. *Nat Biotechnol.* 28, 11 (Nov. 2010), 1181–1185.

[6]Eide, E. 2010. Toward Replayable Research in Networking and Systems. *Archive ’10, the nSF workshop on archiving experiments to raise scientific standards* (2010).

[7]FitzJohn, R. et al. 2014. Reproducible research is still a challenge. <http://ropensci.org/blog/2014/06/09/reproducibility/>.

[8]Garijo, D. et al. 2013. Quantifying reproducibility in computational biology: The case of the tuberculosis drugome. *{PLoS} {ONE}*. 8, 11 (Nov. 2013), e80278.

[9]Gil, Y. et al. 2007. Examining the challenges of scientific workflows. *Computer.* 40, 12 (2007), 24–32.

[10]Gilbert, K.J. et al. 2012. Recommendations for utilizing and reporting population genetic analyses: the reproducibility of genetic clustering using the program structure. *Mol Ecol.* 21, 20 (Sep. 2012), 4925–4930.

[11]Harji, A.S. et al. 2013. Our Troubles with Linux Kernel Upgrades and Why You Should Care. *ACM SIGOPS Operating Systems Review.* 47, 2 (2013), 66–72.

[12]Howe, B. 2012. Virtual appliances, cloud computing, and reproducible research. *Computing in Science & Engineering.* 14, 4 (Jul. 2012), 36–41.

[13]Hull, D. et al. 2006. Taverna: a tool for building and running workflows of services. *Nucleic Acids Research.* 34, Web Server (Jul. 2006), W729–W732.

[14]Ince, D.C. et al. 2012. The case for open computer programs. *Nature.* 482, 7386 (Feb. 2012), 485–488.

[15]Joppa, L.N. et al. 2013. Troubling Trends in Scientific Software Use. *Science (New York, N.Y.).* 340, 6134 (May 2013), 814–815.

[16]Lapp, Hilmar 2014. Reproducibility / repeatability big-Think (with tweets) @hlapp. *Storify.* <http://storify.com/hlapp/reproducibility-repeatability-bigthink>.

[17]Leisch, F. 2002. Sweave: Dynamic Generation of Statistical Reports Using Literate Data Analysis. *Compmat.* W. Härdle and B. Rönz, eds. Physica-Verlag HD.

[18]Merali, Z. 2010. Computational science: ...Error. *Nature.* 467, 7317 (Oct. 2010), 775–777.

[19]Nature Editors 2012. Must try harder. *Nature.* 483, 7391 (Mar. 2012), 509–509.

[20]Ooms, J. 2013. Possible directions for improving dependency versioning in r. *arXiv.org.* <http://arxiv.org/abs/1303.2140v2>.

[21]Ooms, J. 2014. The openCPU system: Towards a universal interface for scientific computing through separation of concerns. *arXiv.org.* <http://arxiv.org/abs/1406.4806>.

[22]Peng, R.D. 2011. Reproducible research in computational science. *Science.* 334, 6060 (Dec. 2011), 1226–1227.

[23]Stodden, V. 2010. The scientific method in practice: Reproducibility in the computational sciences. *SSRN Journal.* (2010).

[24]Stodden, V. et al. 2013. Setting the Default to Reproducible. (2013), 1–19.

[25]The Economist 2013. How science goes wrong. *The Economist*. <http://www.economist.com/news/leaders/21588069-scientific-research-has-changed-world-now-it-needs-change-itself-how-science-goes-wrong>.

[26]Xie, Y. 2013. *Dynamic documents with R and knitr*. Chapman; Hall/CRC.

[27]2014. Examining reproducibility in computer science. <http://cs.brown.edu/~sk/Memos/Examining-Reproducibility/>.

[28]2012. Mick Watson on Twitter: @ewanbirney @pathogenomenick @ctitusbrown you can't install an image for every pipeline you want... <https://twitter.com/BioMickWatson/status/265037994526928896>.