

Application Note

BioContainers: An open-source and community-driven framework for software standardization

Felipe da Veiga Leprevost^{2,θ}, Björn A. Grüning^{4,5,θ}, Saulo Alves Aflitos⁶, Hannes L. Röst⁷, Julian Uszkoreit⁸, Harald Barsnes^{9,10}, Marc Vaudel^{11, 12}, Pablo Moreno¹, Laurent Gatto¹³, Jonas Weber⁴, Mingze Bai¹, Rafael C Jimenez¹, Timo Sachsenberg¹⁴, Julianus Pfeuffer¹⁵, Roberto Vera Alvarez¹⁶, Johannes Griss^{1, 17}, Alexey I. Nesvizhskii^{2, 3}, Yasset Perez-Riverol^{1,*}

¹EMBL Outstation, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, U.K.

²Department of Pathology, University of Michigan, Ann Arbor, 48109, USA.

³Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, 48109, USA

⁴Bioinformatics Group, Department of Computer Science, Albert-Ludwigs-University Freiburg, Freiburg, Germany.

⁵Center for Biological Systems Analysis (ZBSA), Albert-Ludwigs-University Freiburg, Freiburg, Germany.

⁶Applied Bioinformatics, Wageningen University, Wageningen, The Netherlands.

⁷Department of Genetics, Stanford University, USA.

⁸Medizinisches Proteom-Center, Ruhr-Universität Bochum, Bochum, Germany.

⁹Proteomics Unit (PROBE), Department of Biomedicine, University of Bergen, Bergen, Norway.

¹⁰Computational Biology Unit (CBU), Department of Informatics, University of Bergen, Bergen, Norway.

¹¹KG Jebsen Center for Diabetes Research, Department of Clinical Science, University of Bergen, Norway.

¹²Center for Medical Genetics and Molecular Medicine, Haukeland University Hospital, Bergen, Norway.

¹³Computational Proteomics Unit and Cambridge Centre for Proteomics, Department of Biochemistry, University of Cambridge, Cambridge, UK.

¹⁴Center for Bioinformatics, University of Tübingen, Tübingen, Germany.

¹⁵Algorithmic Bioinformatics group, Freie Universität Berlin, Berlin, Germany.

¹⁶Computational Biology Branch, National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland, USA.

¹⁷Division of Immunology, Allergy and Infectious Diseases, Department of Dermatology, Medical University of Vienna, Austria.

* To whom correspondence should be addressed. θ These authors contribute equally to this work.

Associate Editor: Prof. Alfonso Valencia

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: BioContainers (biocontainers.pro) is an open-source and community-driven framework which provides platform independent executable environments for bioinformatics software. BioContainers allows labs of all sizes to easily install bioinformatics software, maintain multiple versions of the same software and combine tools into powerful analysis pipelines. BioContainers is based on popular open-source projects *Docker* and *rkt* frameworks, that allow software to be installed and executed under an isolated and controlled environment. Also, it provides infrastructure and basic guidelines to create, manage and distribute bioinformatics containers with a special focus on omics technologies. These containers can be integrated into more comprehensive bioinformatics pipelines and different architectures (local desktop, cloud environments or HPC clusters).

Availability: The software is freely available at github.com/BioContainers/.

Contact: yperez@ebi.ac.uk, European Molecular Biology Laboratory, EMBL-EBI, Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10 1SD, UK, Tel: +44-1223-492686, Fax: +44-1223-494468.

Introduction

Bioinformatics have emerged as a crucial contributor to our understanding of the function and behavior of systems biology with the development of novel algorithms, the connection of various tools into complex pipelines (Perez-Riverol *et al.*, 2014) and their deposition and dissemination. These developments have been moved from single and individual tools to complex and integrated workflow systems such as OpenMS (Röst *et al.*, 2016), Taverna (Wolstencroft *et al.*, 2013) and Galaxy (Afgan *et al.*, 2016), creating two major challenges for software developers and the bioinformatics community: (i) software availability and (ii) reproducible experiments. Several algorithms software and pipelines in bioinformatics require substantial effort for correct installation and configuration (e.g. conflicting system dependencies). A good starting point for the replicability and reproducibility of the original results should be well-documented (software parameters, dependencies, etc.) and easily installable software (Leprevost *et al.*, 2014). Container based technologies such as Docker (docker.com) or rkt (https://coreos.com/rkt) have emerged to overcome these challenges by automating the deployment of applications inside so-called software *containers*. A software container provides an isolated environment for the installation and execution of a specific software, without affecting other parts of the system. Different groups have proposed the use of Docker containers to solve bioinformatics problems (Moreews *et al.*, 2015) (Belmann *et al.*, 2015). However, most of these projects have been limited to individual efforts and only explore the potential of Docker technology in bioinformatics.

In this manuscript, we present BioContainers (biocontainers.pro), a community-driven project that provides the infrastructure and guidelines to create, manage and distribute bioinformatics containers. The BioContainers architecture facilitates the requests and maintenance of bioinformatics containers, and the interaction between the users and the community. With more than 30 contributors, the community-driven approach guarantees the sustainability and scalability of the project. In addition, BioContainers has been integrated with the BioConda (https://bioconda.github.io/) project enabling the automatic generation of containers for each BioConda recipe. At the time of writing, BioContainers provides more than 2076 containers that can be searched, tagged and accessed through a common web registry (biocontainers.pro/registry/). Finally, we discuss the integration of BioContainers as a container provider with other open-source projects such Galaxy (https://galaxyproject.org/) and PhenoMeNal H2020 (http://phenomenal-h2020.eu/home/).

BioContainers Architecture

The BioContainers architecture is built on two main components: (i) a GitHub organization (github.com/BioContainers/) including all Dockerfiles (for the Dockerfile-based containers), the specification, and tools to create/manage containers; (ii) the BioContainers registries and Registry-UI (biocontainers.pro/registry/) where the available containers are built by an automatic system and made available for download, ready-to-use, by the Docker or rkt (see example, http://biocontainers.pro/docs/101/running-example/). Fig. 1 shows the BioContainers infrastructure from the user request to the final deployment of the container.

Users of BioContainers can request a software container by opening an issue in the container's repository containing information about software (name, URL or binary to be packaged). A member of the BioContainers community will pick up the issue and generate the specific container. An automated build system is configured/deployed making the new container available within hours. To integrate both registries we developed

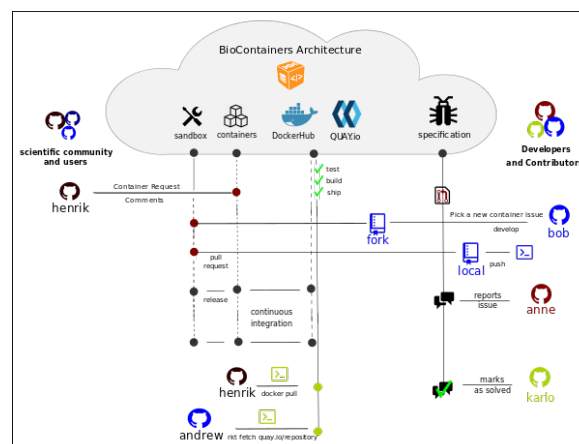


Fig. 1. Overview of the BioContainers architecture: Users and developers can use the BioContainers infrastructure by interacting via GitHub account page. All container Dockerfiles are freely available and people are encouraged to participate submitting pull requests or asking for new containerized software. Containers can be acquired via Docker command line interface, or by downloading the Dockerfile directly from the GitHub organization.

a Registry-UI (biocontainers.pro/registry/) that allows users to search, tag and find BioContainers independently of where they have been deployed. The user can then use docker or rkt to pull or fetch the corresponding container:

```
$> docker pull biocontainers/blast
$> docker run -v /home/user/workplace:/data/
    biocontainers/blast blastp -query seq.fa -db
    zebrafish.fa
```

Dockerfile-based and Mulled Containers

In order to create and build a new container, the BioContainer developers can follow two approaches: (i) create a BioConda recipe for the software or (ii) create a Dockerfile recipe in the container's repository (http://github.com/BioContainers/containers). In the first approach the developer should create a BioConda recipe following the BioConda guidelines (https://bioconda.github.io/guidelines.html). A container generation tool (https://github.com/BioContainers/automulled/) automatically creates a container for the BioConda package and pushes it into BioContainers quay.io registry. These "mulled containers" are generated using the *invulcro* tool (https://github.com/invulcro/invulcro) which enables the generation of containers without any Dockerfile definition, reusing already existing recipes from other package managers, like Conda or Alpine. In summary, *invulcro* will install the given (Conda) package into a build-time container which has the the preferred package manager already installed and copies the resulting new image layer on top of a runtime environment defined by BioContainers (busybox). (ii) In the second approach, a recipe file must be named Dockerfile which holds all the instructions necessary for creating the complete container. As part of the project specifications, we are providing a template for developers to "containerize" their own applications (https://github.com/BioContainers/specs/blob/master/container-specs.md). For each BioContainers the developer should provide metadata about the software such as the name, version, license, web-page and the maintainer. Both strategies are already aligned and the metadata needed to create a BioConda recipe in the YAML file is the same we recommended for the

Dockerfile. This metadata enables BioContainers to find, describe and maintain each containers following best practices (Leprevost *et al.*, 2014).

Tools and Future Directions

At the time of writing, BioContainers provides more than 2076 containers ready to be used. The integration with the BioConda project (bioconda.github.io) has enabled us to create a new type of containers without any Dockerfile, reusing already existing BioConda recipes. The Galaxy Project has recently proposed Docker containers as a new way to solve workflow dependencies (biocontainers.pro/docs). Also, the PhenoMeNal H2020 project has adopted and implemented BioContainers guidelines and deploying their containers into the BioContainers architecture. The BioContainers community is now working on new ways for testing containers and for workflow/pipelines integration.

Funding

F.V.L. and A.I.N. are supported by NIH grant numbers R01-GM-094231 and U24-CA0-210967 (to A.I.N.). H.L.R. is supported by the Swiss National Science Foundation (SNSF grant P2EZP3 162268) and EMBO (ALTF 854-2015). H.B. is supported by the Bergen Research Foundation and the Research Council of Norway. T.S., J.P. and J.U. acknowledge funding from BMBF (de.NBI, grant nos. FKZ 031 A 535A and FKZ 031 A 534A). Y.P.-R. is supported by US NIH BD2K grant [U54 GM114833]. LG is supported by the BBSRC Strategic Longer and Larger grant (Award BB/L002817/1). P.M. is supported by EC Horizon 2020 grant agreement 654241.

References

- Afgan, E., Baker, D., van den Beek, M., Blankenberg, D., Bouvier, D., Čech, M., Chilton, J., Clements, D., Coraor, N., Eberhard, C., Grüning, B., Guerler, A., Hillman-Jackson, J., Von Kuster, G., Rasche, E., Soranzo, N., Turaga, N., Taylor, J., Nekrutenko, A., and Goecks, J. (2016). The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Research*, **44**(W1), W3–W10.
- Belmann, P., Dröge, J., Bremges, A., McHardy, A. C., Sczyrba, A., and Barton, M. D. (2015). Bioboxes: standardised containers for interchangeable bioinformatics software. *GigaScience*, **4**.
- Leprevost, F. d. V., Barbosa, V. C., Francisco, E. L., Perez-Riverol, Y., and Carvalho, P. C. (2014). On best practices in the development of bioinformatics software. *Bioinformatics and Computational Biology*, **5**, 199.
- Moreews, F., Sallou, O., Ménager, H., Le bras, Y., Monjeaud, C., Blanchet, C., and Collin, O. (2015). BioShaDock: a community driven bioinformatics shared Docker-based tools registry. *F1000Research*, **4**.
- Perez-Riverol, Y., Wang, R., Hermjakob, H., Müller, M., Vesada, V., and Vizcaíno, J. A. (2014). Open source libraries and frameworks for mass spectrometry based proteomics: A developer's perspective. *Biochimica et Biophysica Acta (BBA) - Proteins and Proteomics*, **1844**(1, Part A), 63–76.
- Röst, H. L., Sachsenberg, T., Aiche, S., Bielow, C., Weissner, H., Aicheler, F., Andreotti, S., Ehrlich, H.-C., Gutenbrunner, P., Kenar, E., Liang, X., Nahnsen, S., Nilse, L., Pfeuffer, J., Rosenberger, G., Rurik, M., Schmitt, U., Veit, J., Walzer, M., Wojnar, D., Wolski, W. E., Schilling, O., Choudhary, J. S., Malmström, L., Aebersold, R., Reinert, K., and Kohlbacher, O. (2016). OpenMS: a flexible open-source software platform for mass spectrometry data analysis. *Nature Methods*, **13**(9), 741–748.
- Wolstencroft, K., Haines, R., Fellows, D., Williams, A., Withers, D., Owen, S., Soiland-Reyes, S., Dunlop, I., Nenadic, A., Fisher, P., Bhagat, J., Belhajjame, K., Bacall, F., Hardisty, A., Hidalgo, A. N. d. l., Vargas, M. P. B., Sufi, S., and Goble, C. (2013). The Taverna workflow suite: designing and executing workflows of Web Services on the desktop, web or in the cloud. *Nucleic Acids Research*, **41**(W1), W557–W561.