

Bioconda: A sustainable and comprehensive software distribution for the life sciences

Ryan Dale¹, Björn Grüning^{*,2}, Andreas Sjödin^{3,4}, Jillian Rowe⁵, Brad A. Chapman⁶, Christopher H. Tomkins-Tinch^{7,8}, Renan Valieris⁹, Bérénice Batut², Adam Caprez¹⁰, Thomas Cokelaer¹¹, Dilmurat Yusuf², Kyle A. Beauchamp¹², Karel Brinda¹³, Thomas Wollmann¹⁴, Gildas Le Corguillé¹⁵, Devon Ryan¹⁶, Anthony Bretraudeau¹⁷, Youri Hoogstrate¹⁸, Brent S Pedersen¹⁹, Simon van Heeringen²⁰, Martin Raden², Sebastian Luna-Valero²¹, Nicola Soranzo²², Matthias De Smet²³, Greg Von Kuster²⁴, Rory Kirchner²⁵, Lorena Pantano⁶, Zachary Charlop-Powers²⁶, Kevin Thornton²⁷, Marcel Martin²⁸, Marius van den Beek²⁹, Daniel Maticzka², Milad Miladi², Sebastian Will³⁰, Kévin Gravouil³¹, Per Unneberg³², Christian Brueffer³³, Clemens Blank², Vitor C. Piro^{34,35}, Joachim Wolff², Tiago Antao³⁶, Simon Gladman³⁷, Ilya Shlyakhter⁸, Mattias de Hollander³⁸, Philip Mabon³⁹, Wei Shen⁴⁰, Jorrit Boekel⁴¹, Manuel Holtgrewe^{42,43}, Dave Bouvier⁴⁴, Julian R. de Ruiter⁴⁵, Jennifer Cabral³⁹, Saket Choudhary⁴⁶, Nicholas Harding⁴⁷, Robert Kleinkauf⁴⁸, Eric Enns³⁹, Florian Eggenhofer², Joseph Brown⁴⁹, Peter J. A. Cock⁵⁰, Henning Timm⁵¹, Cristel Thomas⁵², Xiao-Ou Zhang⁵³, Matt Chambers⁵⁴, Nitesh Turaga⁵⁵, Enrico Seiler³⁴, Colin Brislawn⁵⁶, Elmar Pruesse⁵⁷, Jörg Fallmann⁵⁸, Jerome Kelleher⁴⁷, Hai Nguyen⁵⁹, Lance Parsons⁶⁰, Zhuoqing Fang⁶¹, Endre Bakken Stovner⁶², Nicholas Stoler⁶³, Simon Ye⁶⁴, Inken Wohlers⁶⁵, Rick Farouni⁶⁶, Mallory Freeberg⁶⁷, James E. Johnson⁶⁸, Marcel Bargull⁵¹, Philip Reiner Kensche⁶⁹, Timothy H. Webster⁷⁰, John M Eppley⁷¹, Christoph Stahl⁷², Alexander S Rose⁷³, Alex Reynolds⁷⁴, Liang-Bo Wang^{75,76}, Xavier Garnier^{17,77}, Simon Dirmeier⁷⁸, Michael Knudsen⁷⁹, James Taylor⁸⁰, Avi Srivastava⁸¹, Vivek Rai⁸², Rasmus Agren⁸³, Alexander Junge⁸⁴, Roman Valls Guimera⁸⁵, Aziz Khan⁸⁶, Sebastian Schmeier⁸⁷, Guowei He⁸⁸, Luca Pinello⁶⁶, Emil Hägglund⁸⁹, Alexander S Mikheyev^{90,91}, Jens Preussner⁹², Nicholas R. Waters⁹³, Wei Li⁹⁴, Jordi Capellades⁹⁵, Aroon T. Chande⁹⁶, Yuri Pirola⁹⁷, Saskia Hiltemann⁹⁸, Matthew L. Bendall^{99,100}, Sourav Singh¹⁰¹, W. Augustine Dunn¹⁰², Alexandre Drouin¹⁰³, Tomás Di Domenico¹⁰⁴, Ino de Bruijn¹⁰⁵, David E Larson¹⁰⁶, Davide Chicco¹⁰⁷, Elena Grassi¹⁰⁸, Giorgio Gonnella¹⁰⁹, Jaivarsan B¹¹⁰, Liya Wang¹¹¹, Franck Giacomoni¹¹², Erik Clarke¹¹³, Daniel Blankenberg¹¹⁴, Camy Tran³⁹, Rob Patro⁸¹, Sacha Laurent¹¹⁵, Matthew Gopez³⁹, Bengt Sennblad³², Jasmijn A. Baaijens¹¹⁶, Philip Ewels¹¹⁷, Patrick R. Wright², Oana M. Enache⁸, Pierrick Roger¹¹⁸, Will Dampier¹¹⁹, David Koppstein¹²⁰, Upendra Kumar Devisetty¹²¹, Tobias Rausch¹²², MacIntosh Cornwell¹²³, Adrian Emanuel Salatino¹²⁴, Julien Seiler¹²⁵, Matthieu Jung¹²⁵, Etienne Kornobis¹²⁶, Fabio Cumbo^{127,128,63}, Bianca Katharina Stöcker⁵¹, Oleksandr Moskalenko¹²⁹, Daniel R. Bogema¹³⁰, Matthew L. Workentine¹³¹, Stephen J. Newhouse^{132,133,134}, Felipe da Veiga Leprevost¹³⁵, Kevin Arvai¹³⁶, and Johannes Köster^{†,137,138}

¹Laboratory of Cellular and Developmental Biology, National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda, United States

²Bioinformatics Group, Department of Computer Science, University of Freiburg, Freiburg,

Germany

³Division of CBRN Security and Defence, FOI - Swedish Defence Research Agency, Umeå, Sweden

⁴Department of Chemistry, Computational Life Science Cluster (CLiC), Umeå University, Umeå, Sweden

⁵NYU Abu Dhabi, Abu Dhabi, United Arab Emirates

⁶Harvard T.H. Chan School of Public Health, Boston, United States

⁷Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, United States

⁸Broad Institute of MIT and Harvard, Cambridge, United States

⁹Laboratory of Bioinformatics and Computational Biology, A. C. Camargo Cancer Center, São Paulo, Brazil

¹⁰Holland Computing Center, University of Nebraska, Lincoln, United States

¹¹Institut Pasteur - Bioinformatics and Biostatistics Hub - C3BI, Paris, France

¹²Counsyl, South San Francisco, United States

¹³Center for Communicable Disease Dynamics, Department of Epidemiology, Harvard TH Chan School of Public Health, Boston, United States

¹⁴University of Heidelberg and DKFZ, Heidelberg, Germany

¹⁵UPMC, CNRS, FR2424, ABiMS, Station Biologique, Roscoff, France

¹⁶Bioinformatics core facility, Max Planck Institute for Immunobiology and Epigenetics, Freiburg, Germany

¹⁷INRA, UMR IGEPP, Bioinformatics Platform for Agroecosystems Arthropods (BIPAA), Campus Beaulieu, Rennes, France

¹⁸Erasmus Medical Center, Department of Urology, Rotterdam, The Netherlands

¹⁹Department of Human Genetics, University of Utah, Eccles Institute of Human Genetics, Salt Lake City

²⁰Radboud University, Faculty of Science, Department of Molecular Developmental Biology, Radboud Institute for Molecular Life Sciences, Nijmegen, The Netherlands

²¹MRC Weatherall Institute of Molecular Medicine, University of Oxford, Oxford, United Kingdom

²²Earlham Institute, Norwich Research Park, Norwich, United Kingdom

²³Ghent University Hospital, Ghent University, Belgium

²⁴Institute for CyberScience, Penn State University, University Park, United States

²⁵Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, United States

²⁶The Laboratory for Genetically Encoded Small Molecules, The Rockefeller University, New York, United States

²⁷Department of Ecology and Evolutionary Biology, University of California Irvine, Irvine, United States

²⁸Department of Biochemistry and Biophysics, National Bioinformatics Infrastructure Sweden, Science for Life Laboratory, Stockholm University, Sweden

²⁹Stem Cells and Tissue Homeostasis, Institut Curie, Paris, France

³⁰Institute for Theoretical Chemistry, University of Vienna, Vienna, Austria

³¹Université Clermont Auvergne, INRA, MEDIS, Clermont-Ferrand, France

- ³²Department of Cell and Molecular Biology, National Bioinformatics Infrastructure Sweden, Science for Life Laboratory, Uppsala University, Uppsala, Sweden
- ³³Division of Oncology and Pathology, Department of Clinical Sciences, Lund University, Lund, Sweden
- ³⁴Bioinformatics Unit, Robert Koch Institute, Berlin, Germany
- ³⁵CAPES Foundation, Ministry of Education of Brazil, Brasília, Brazil
- ³⁶Division of Biological Sciences, University of Montana, Missoula, United States of America
- ³⁷Melbourne Bioinformatics, University of Melbourne, Melbourne, Australia
- ³⁸Netherlands Institute of Ecology (NIOO-KNAW), Wageningen, The Netherlands
- ³⁹National Microbiology Laboratory, Public Health Agency of Canada, Winnipeg, Canada
- ⁴⁰Department of Clinical Laboratory, Chengdu Military General Hospital, Chengdu, China
- ⁴¹Department of Oncology-Pathology, National Bioinformatics Infrastructure Sweden, Science for Life Laboratory, Karolinska Institutet, Solna, Sweden
- ⁴²Core Unit Bioinformatics, Berlin Institute of Health, Berlin, Germany
- ⁴³Charité Universitätsmedizin Berlin, Berlin, Germany
- ⁴⁴Department of Biochemistry Molecular Biology, Penn State University, University Park, United States
- ⁴⁵Divisions of Molecular Pathology and Molecular Carcinogenesis, The Netherlands Cancer Institute, Amsterdam, The Netherlands
- ⁴⁶Computational Biology and Bioinformatics, University of Southern California, Los Angeles, United States
- ⁴⁷Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of Oxford, Oxford, United Kingdom
- ⁴⁸_
- ⁴⁹Biological Sciences Division, Pacific Northwest National Laboratory, Richland, United States
- ⁵⁰The James Hutton Institute, Dundee, United Kingdom
- ⁵¹Genome Informatics, Institute of Human Genetics, University Hospital Essen, University of Duisburg-Essen, Essen, Germany
- ⁵²Northrop Grumman Corporation, Technology Services, Rockville, United States
- ⁵³Program in Bioinformatics and Integrative Biology, University of Massachusetts Medical School, Worcester, United States
- ⁵⁴Department of Biochemistry, Molecular Biology and Biophysics (as contractor, not employee), University of Minnesota, Minneapolis, United States
- ⁵⁵Department of Biology, Johns Hopkins University, Baltimore, United States
- ⁵⁶Earth and Biological Sciences Directorate, Pacific Northwest National Laboratory, Richland, United States
- ⁵⁷University of Colorado, Denver, United States
- ⁵⁸Bioinformatics Group, Department of Computer Science and Interdisciplinary Center for Bioinformatics, University of Leipzig, Leipzig, Germany
- ⁵⁹Department of Chemistry Chemical Biology, Rutgers University, Piscataway, United States
- ⁶⁰Lewis-Sigler Institute for Integrative Genomics, Princeton University, Princeton, United States

States

⁶¹Shanghai Institutes for Biological Sciences, University of Chinese Academy of Sciences, Chinese Academy of Sciences, Shanghai China

⁶²Department of Computer Science, Norwegian University of Science and Technology

⁶³Department of Biochemistry and Molecular Biology, Penn State University, University Park, United States

⁶⁴Massachusetts Institute of Technology, Cambridge, United States

⁶⁵Lübeck Interdisciplinary Platform for Genome Analytics (LIGA), Institutes of Neurogenetics and Integrative Experimental Genomics, University of Lübeck, Lübeck, Germany

⁶⁶Massachusetts General Hospital and Harvard Medical School, Boston, United States

⁶⁷EMBL-EBI, Wellcome Genome Campus, Hinxton, Cambridgeshire, United Kingdom

⁶⁸Minnesota Supercomputing Institute, University of Minnesota, Minneapolis, United States

⁶⁹German Cancer Research Center (DKFZ), Foundation under Public Law, Heidelberg, Germany

⁷⁰School of Life Sciences, Arizona State University, Tempe, United States

⁷¹Daniel K. Inouye Center for Microbial Oceanography: Research and Education, Department of Oceanography, University of Hawaii, Honolulu, United States

⁷²Genome Informatics, Institute of Human Genetics, University Hospital Essen, University of Duisburg Essen, Essen, Germany

⁷³RCSB Protein Data Bank, San Diego Supercomputer Center, University of California, San Diego, United States

⁷⁴Altius Institute for Biomedical Sciences, Seattle, United States

⁷⁵Oncology Division, Department of Medicine, Washington University School of Medicine, St. Louis, United States

⁷⁶McDonnell Genome Institute, Washington University School of Medicine, St. Louis, United States

⁷⁷Dyliss - Dynamics, Logics and Inference for biological Systems and Sequences, Inria/IRISA, Campus Beaulieu, Rennes, France

⁷⁸Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland

⁷⁹Department of Molecular Medicine, Aarhus University Hospital, Aarhus, Denmark

⁸⁰Departments of Biology and Computer Science, Johns Hopkins University, Baltimore, United States

⁸¹Department of Computer Science, Stony Brook University, Stony Brook, United States

⁸²Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, United States

⁸³Department of Biology and Biological Engineering, National Bioinformatics Infrastructure Sweden, Science for Life Laboratory, Chalmers University of Technology, Sweden

⁸⁴Disease Systems Biology Program, Novo Nordisk Foundation Center for Protein Research, University of Copenhagen, Copenhagen, Denmark

⁸⁵Center for Cancer Research, University of Melbourne, Melbourne, Australia

⁸⁶Centre for Molecular Medicine Norway (NCMM), Nordic EMBL Partnership, University of Oslo, Oslo, Norway

- ⁸⁷Massey University, Institute of Natural and Mathematical Sciences, North Shore City, New Zealand
- ⁸⁸High Performance Computing, NYU Abu Dhabi, Abu Dhabi, United Arab Emirates
- ⁸⁹Department of Molecular Evolution, Cell and Molecular Biology, Science for Life Laboratory, Biomedical Centre, Uppsala University, Uppsala, Sweden
- ⁹⁰Evolutionary Genomics Lab, Research School of Biology, The Australian National University, Canberra, Australia
- ⁹¹Ecology and Evolution Unit, Okinawa Institute of Science and Technology Graduate University, Onna-son, Kunigami-gun, Okinawa, Japan
- ⁹²ECCPS Bioinformatics Core Unit, Max Planck Institute for Heart and Lung Research, Bad Nauheim, Germany
- ⁹³Department of Microbiology, School of Natural Sciences, National University of Ireland, Galway, Ireland Information and Computational Sciences, James Hutton Institute, Invergowrie, Scotland
- ⁹⁴Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Harvard T.H. Chan School of Public Health, Boston, United States
- ⁹⁵Universitat Rovira i Virgili, Spanish Biomedical Research Center in Diabetes and Associated Metabolic Disorders (CIBERDEM), Reus Spain
- ⁹⁶Applied Bioinformatics Laboratory, 2 Ravinia Drive, Suite 1200 Atlanta, GA 30346, United States
- ⁹⁷Dip. di Informatica Sistemistica e Comunicazione, Univ. degli Studi di Milano-Bicocca, Milan, Italy
- ⁹⁸Erasmus Medical Center, Rotterdam, The Netherlands
- ⁹⁹Computational Biology Institute, Milken Institute School of Public Health, The George Washington University, Washington, D.C., United States
- ¹⁰⁰Department of Microbiology, Immunology Tropical Medicine, The George Washington University School of Medicine and Health Sciences, Washington, D.C., United States
- ¹⁰¹Savitribai Phule Pune University, Pune, Maharashtra, India
- ¹⁰²Boston Children's Hospital, Boston, United States
- ¹⁰³Department of Computer Science and Software Engineering, Université Laval, Québec, Canada
- ¹⁰⁴Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, United Kingdom
- ¹⁰⁵Department of Pathology, Memorial Sloan Kettering Cancer Center, New York, United States
- ¹⁰⁶The McDonnell Genome Institute, Washington University, St. Louis, United States
- ¹⁰⁷Princess Margaret Cancer Centre, Toronto, Canada
- ¹⁰⁸Transcription and Chromatin Lab, Humanitas University, Rozzano, Italy
- ¹⁰⁹ZBH - Center for Bioinformatics, MIN-Fakultät, Universität Hamburg, Hamburg, Germany
- ¹¹⁰Department of Computer Science, Computer science with specialisation in Bioinformatics, VIT University Vellore, India
- ¹¹¹Cold Spring Harbor Laboratory, Cold Spring Harbor, United States
- ¹¹²Clermont Auvergne University, INRA, UNH, Human Nutrition Unit, PFEM, Metabolism Exploration Platform, MetaboHUB-Clermont, Clermont-Ferrand, France

- ¹¹³Department of Microbiology, University of Pennsylvania, United States
¹¹⁴Genomic Medicine Institute, Lerner Research Institute, Cleveland Clinic, Cleveland, United States
¹¹⁵Institute of Microbiology, University Hospital of Lausanne, Switzerland
¹¹⁶Centrum Wiskunde and Informatica, Amsterdam, Netherlands
¹¹⁷Department of Biochemistry and Biophysics, Science for Life Laboratory, Stockholm University, Stockholm, Sweden
¹¹⁸CEA, LIST, Laboratory for data analysis and systems' intelligence, MetaboHUB, France
¹¹⁹Drexel University College of Medicine, Department of Microbiology and Immunology, Philadelphia, United States
¹²⁰The Kirby Institute of Infection and Immunity, University of New South Wales, Sydney, Australia
¹²¹CyVerse, Bio5 institute, University of Arizona, Tucson, United States
¹²²European Molecular Biology Laboratory (EMBL), Genomics Core Facility, Heidelberg, Germany
¹²³New York University School of Medicine, New York City, United States
¹²⁴Department of Molecular Genetics and Biology of Complex Diseases, Institute of Medical Research A Lanari-IDIM, University of Buenos Aires, National Scientific and Technical Research Council (CONICET), Ciudad Autónoma de Buenos Aires, Argentina.
¹²⁵Institut de Génétique et de Biologie Moléculaire et Cellulaire, CNRS, Illkirch, France
¹²⁶Epigenetic Regulation Unit, Pasteur Institute, Paris, France
¹²⁷Department of Engineering, Roma Tre University, Rome, Italy
¹²⁸Institute for Systems Analysis and Computer Science "Antonio Ruberti", National Research Council of Italy, Rome, Italy
¹²⁹UFIT Research Computing, University of Florida, Gainesville, United States
¹³⁰NSW Department of Primary Industries, Elizabeth Macarthur Agricultural Institute, Menangle, Australia
¹³¹Faculty of Veterinary Medicine, University of Calgary, Calgary, Canada
¹³²Department of Biostatistics and Health Informatics, Institute of Psychiatry, Psychology and Neuroscience, King's College London, London, United Kingdom
¹³³NIHR Biomedical Research Centre at South London and Maudsley NHS Foundation Trust and King's College London, United Kingdom
¹³⁴Farr Institute of Health Informatics Research, UCL Institute of Health Informatics, University College London, London, United Kingdom
¹³⁵Pathology Department, University of Michigan, Ann Arbor, United States
¹³⁶GeneDx, Gaithersburg, United States
¹³⁷Algorithms for reproducible bioinformatics, Genome Informatics, Institute of Human Genetics, University Hospital Essen, University of Duisburg-Essen
¹³⁸Dana Farber Cancer Institute, Harvard Medical School, Boston, United States

October 21, 2017

*Co-first author

†To whom correspondence should be addressed.

Abstract

We present Bioconda (<https://bioconda.github.io>), a distribution of bioinformatics software for the lightweight, multi-platform and language-agnostic package manager, Conda. Currently, Bioconda offers a collection of over 2900 software tools, which are continuously maintained, updated, and extended by a growing global community of more than 200 contributors. Bioconda improves analysis reproducibility by allowing users to define isolated environments with defined software versions, all of which are easily installed and managed without administrative privileges.

Introduction

Thousands of new software tools have been released for bioinformatics in recent years, in a variety of programming languages. Accompanying this diversity of construction is an array of installation methods. Software written in C/C++ often has to be compiled manually for different hardware architectures and operating systems, with management left to the user or system administrator. Scripting languages usually deliver their own package management tools for installing, updating, and removing packages, though these are often limited in scope to packages written in the same scripting language, such that external dependencies (e.g., C libraries) have to be installed manually. Published scientific software often consists of simple collections of custom scripts distributed with textual descriptions of the manual steps required to use the software. New analyses often require novel combinations of multiple tools, and the heterogeneity of scientific software makes management of a software stack complicated and error-prone. Moreover, it inhibits reproducible science (Mesirov, 2010; Baker, 2016; Munafò et al., 2017) because it is hard to reproduce a software stack on different machines. System-wide deployment of software has traditionally been handled by administrators, but reproducibility often requires that the researcher (who is often not an expert in administration) is able to maintain full control of the software environment and rapidly modify it without administrative privileges.

The Conda package manager (<https://conda.io>) has become an increasingly popular approach to overcome these challenges. Conda normalizes software installations across language ecosystems by describing each software package with a *recipe* that defines meta-information and dependencies, as well as a *build script* that performs the steps necessary to build and install the software. Conda prepares and builds software packages within an isolated environment, transforming them into relocatable binaries. Conda packages can be built for all three major operating systems: Linux, macOS, and Windows. Importantly, installation and management of packages requires no administrative privileges, such that a researcher can control the available software tools regardless of the underlying infrastructure. Moreover, Conda obviates reliance on system-wide installation by allowing users to generate isolated software environments, within which versions and tools can be managed per-project, without generating conflicts or incompatibilities (see online methods). These environments support reproducibility, as they can be rapidly exchanged via files that describe their installation state. Conda is tightly integrated into popular solutions for reproducible scientific data analysis like Galaxy (Afgan et al., 2016), bcbio-nextgen (<https://github.com/chapmanb/bcbio-nextgen>), and Snakemake (Köster and Rahmann, 2012). Finally, while Conda provides many commonly-used packages by default, it also allows users to optionally include additional repositories (termed *channels*) of packages that can be installed.

Results

In order to unlock the benefits of Conda for the life sciences, the Bioconda project was founded in 2015. The mission of Bioconda is to make bioinformatics software easily installable and manageable via the Conda package manager. Via its channel for the Conda package manager, Bioconda currently provides over 2500 software packages for Linux and macOS. Development is driven by an open community of over 200 international scientists. In the prior two years, package count and the number of contributors have increased

linearly, on average, with no sign of saturation (Fig. 1a,b). The barrier to entry is low, requiring a willingness to participate and adherence to community guidelines. Many software developers contribute recipes for their own tools, and many Bioconda contributors are invested in the project as they are also users of Conda and Bioconda. Bioconda provides packages from various language ecosystems like Python, R (CRAN and Bioconductor), Perl, Haskell, as well as a plethora of C/C++ programs (Fig. 1c). Many of these packages have complex dependency structures that require various manual steps to install when not relying on a package manager like Conda (Fig. 2a, Online Methods). With over 5.9 million downloads, the service has become a backbone of bioinformatics infrastructure (Fig. 1d). Bioconda is complemented by the conda-forge project (<https://conda-forge.github.io>), which hosts software not specifically related to the biological sciences. The two projects collaborate closely, and the Bioconda team maintains over 500 packages hosted by conda-forge. Among all currently available distributions of bioinformatics software, Bioconda is by far the most comprehensive, while being among the youngest (Fig. 2d).

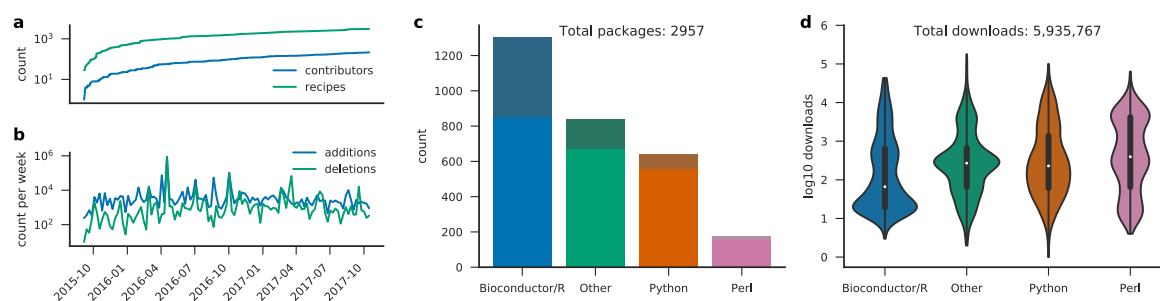


Figure 1: Bioconda development and usage since the beginning of the project (state: October 2017). (a) contributing authors and added recipes over time. (b) code line additions and deletions per week. (c) package count per language ecosystem (saturated colors on bottom represent explicitly life science related packages). (d) total downloads per language ecosystem. The term “other” entails all recipes that do not fall into one of the specific Note that a subset of packages that started in Bioconda have since been migrated to the more appropriate, general-purpose conda-forge channel. Older versions of such packages still reside in the Bioconda channel, and as such are included in the recipe count (a) and download count (d).categories.

To ensure reliable maintenance of such numbers of packages, we use a semi-automatic, agent-assisted development workflow (Fig. 2b). All Bioconda recipes are hosted in a GitHub repository (<https://github.com/bioconda/bioconda-recipes>). Both the addition of new recipes and the update of existing recipes in Bioconda is handled via *pull requests*. Thereby, a modified version of one or more recipes is compared against the current state of Bioconda. Once a pull request arrives, our infrastructure performs several automatic checks. Problems discovered in any step are reported to the contributor and further progress is blocked until they are resolved. First, the modified recipes are checked for syntactic anti-patterns, i.e., formulations that are syntactically correct but bad style (termed *linting*). Second, the modified recipes are built on Linux and macOS, via a cloud based, free-of-charge service (<https://travis-ci.org>). Successfully built recipes are tested (e.g., by running the generated executable). Since Bioconda packages must be able to run on any supported system, it is important to check that the built packages do not rely on particular elements from the build environment. Therefore, testing happens in two stages: (a) test cases are executed in the build environment (b) test cases are executed in a minimal Docker (<https://docker.com>) container which purposefully lacks all non-common system libraries (hence a dependency that is not explicitly defined will lead to a failure). Once the *build* and *test* steps have succeeded, a member of the Bioconda team reviews the proposed changes and, if acceptable, merges the modifications into the official repository. Upon merging, the recipes are built again and uploaded to the hosted Bioconda channel (<https://anaconda.org/bioconda>), where they become available via the Conda package manager. When a Bioconda package is updated to a new version, older builds are generally preserved, and recipes for multiple older versions may be maintained

in the Bioconda repository. The usual turnaround time of above workflow is fast (Fig. 2d). 61% of the pull requests are merged within 5 hours. Of those, 36% are even merged within 1 hour. Only 18% of the pull requests need more than a day. Hence, publishing software in Bioconda or updating already existing packages can be accomplished typically within minutes to a few hours.

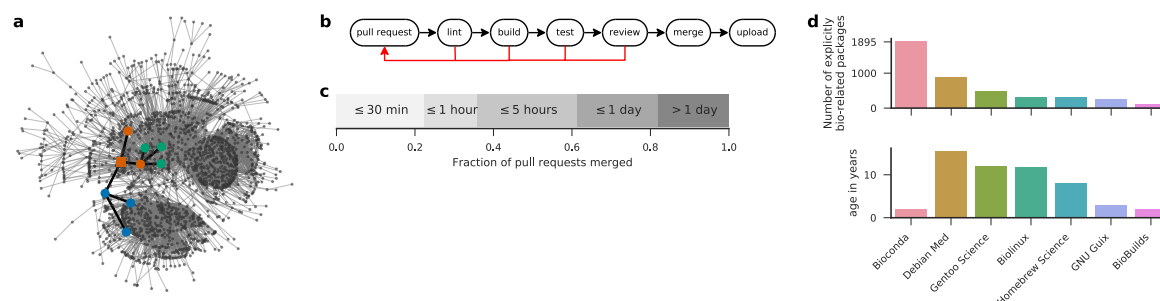


Figure 2: Dependency structure, workflow, comparison with other resources, and turnaround time. (a) largest connected component of directed acyclic graph of Bioconda packages (nodes) and dependencies (edges). Highlighted is the induced subgraph of the CNVkit (Talevich et al., 2016) package and its dependencies (node coloring as defined in Fig. 1, squared node represents CNVkit). (b) Github based development workflow: a contributor provides a pull request that undergoes several build and test steps, followed by a human review. If any of these checks does not succeed, the contributor can update the pull request accordingly. Once all steps have passed, the changes can be merged. (c) Turnaround time from submission to merge of pull requests in Bioconda. (d) Comparison of explicitly life science related packages in Bioconda with Debian Med (<https://www.debian.org/devel/debian-med>), Gentoo Science Overlay (category sci-biology, <https://github.com/gentoo/sci>), Biolinux (Field et al., 2006), Homebrew Science (tag bioinformatics, <https://brew.sh>), GNU Guix (category bioinformatics, <https://www.gnu.org/s/guix>), and BioBuilds (<https://biobuilds.org>). The lower panel shows the age since the first release or commit.

Reproducible software management and distribution is enhanced by other current technologies. Conda integrates itself well with environment modules (<http://modules.sourceforge.net/>), a technology used nearly universally across HPC systems. An administrator can use Conda to easily define software stacks for multiple labs and project-specific configurations. Popularized by Docker, containers provide another way to publish an entire software stack, down to the operating system. They provide greater isolation and control over the environment a software is executed in, at the expense of some customizability. Conda complements container-based approaches. Where flexibility is needed, Conda packages can be used and combined directly. Where the uniformity of containers is required, Conda can be used to build images without having to reproduce the nuanced installation steps that would ordinarily be required to build and install an application within an image. In fact, for each Bioconda package, the build system automatically builds a minimal Docker image containing that package, which is subsequently uploaded and made available via the Biocontainers project (da Veiga Leprevost et al., 2017). As a consequence, every built Bioconda package is available not only for installation via Conda, but also directly available as a container via Docker, Rkt (<https://coreos.com/rkt>), and Singularity (Kurtzer et al., 2017).

Discussion

For reproducible data science, it is crucial that software libraries and tools are provided via an easy to use, unified interface, such that they can be easily deployed and sustainably managed. With its ability to maintain isolated software environments, the integration into major workflow management systems and the fact that no administration privileges are needed, the Conda package manager is the ideal tool to ensure sustainable

and reproducible software management. With Bioconda, we unlock Conda for the life sciences and coordinate closely with other related projects such as conda-forge and Biocontainers. Bioconda offers a comprehensive resource of thousands of software libraries and tools that is maintained by hundreds of international scientists. With almost six million downloads so far, Bioconda packages have been well received by the community. We invite everybody to participate in reaching the goal of a central, comprehensive, and language agnostic collection of easily installable software by maintaining existing or publishing new software in Bioconda.

Funding

The Bioconda project has received support from Anaconda, Inc., Austin, TX, USA, in the form of expanded storage for Bioconda packages on their hosting service (<https://anaconda.org>). Further, the project has been granted extended build times from Travis CI, GmbH (<https://travis-ci.com>). The Bioconda community also would like to thank ELIXIR (<https://www.elixir-europe.org>) for their constant support.

Acknowledgments

We thank the participants of various hackathons (e.g., the GalaxyP contribution fest) for porting numerous packages to Bioconda.

Online Methods

Security Considerations

Using Bioconda as a service to obtain packages for local installation entails trusting that (a) the provided software itself is not harmful and (b) it has not been modified in a harmful way. Ensuring (a) is up to the user. In contrast, (b) is handled by our workflow. First, source code or binary files defined in recipes are checked for integrity via MD5 or SHA256 hash values. Second, all review and testing steps are enforced via the Github interface. This way, it is guaranteed that all packages have been tested automatically and reviewed by a human being. Third, all changes to the repository of recipes are publicly tracked, and all build and test steps are transparently visible to the user. Finally, the automatic parts of the development workflow are implemented in the open-source software *bioconda-utils* (<https://github.com/bioconda/bioconda-utils>). In the future, we will further explore the possibility to sign packages cryptographically.

Software management with Conda

Via the Conda package manager, installing software from Bioconda becomes very simple. In the following, we describe the basic functionality assuming that the user has access to a Linux or macOS terminal. After installing Conda, the first step is to set up the Bioconda channel via:

```
$ conda config --add channels conda-forge
$ conda config --add channels bioconda
```

Now, all Bioconda packages are visible to the Conda package manager. For example, the software CN-Vkit (Talevich et al., 2016), can be searched for with

```
$ conda search cnvkit
```

in order to check if and in which versions it is available. It can be installed with:

```
$ conda install cnvkit
```

CNVkit needs various dependencies from Python and R, which would otherwise have to be installed in separate manual steps (Fig. 2a). Furthermore, Conda enables updating and removing all these dependencies via one unified interface. A key value of Conda is the ability to define isolated, shareable software environments. This can happen ad-hoc, or via YAML (<https://yaml.org>) files. For example, the following defines an environment consisting of Salmon (Patro et al., 2017) and DESeq2 (Love et al., 2014):

```
channels:
  - bioconda
  - conda-forge
  - defaults
dependencies:
  - bioconductor-deseq2 =1.16.1
  - salmon =0.8.2
  - r-base =3.4.1
```

Given that the environment is stored in a file `env.yaml`, it can be created with the name `my-env` via the command:

```
$ conda env create --name my-env --file env.yaml
```

Then, in order to use the environment it can be activated with:

```
$ source activate my-env
```

Within the environment R, Salmon and DESeq2 are available in exactly the defined versions. For example, salmon can be executed with:

```
$ salmon --help
```

It is possible to modify an existing environment by using `conda update`, `conda install` and `conda remove`. For example, we could add a particular version of Kallisto (Bray et al., 2016) and update Salmon to the latest available version with:

```
$ conda install kallisto=0.43.1
$ conda update salmon
```

Finally, the environment can be deactivated again with:

```
$ source deactivate
```

How isolated software environments enable reproducible research

With isolated software environments as shown above, it is possible to define an exact version for each package. This increases reproducibility by eliminating differences due to implementation changes. Note that above we also pin an R version, although the latest compatible one would also be automatically installed without mentioning it. To further increase reproducibility, this pattern can be extended to all dependencies of DESeq2 and Salmon and recursively down to basic system libraries like zlib and boost (<https://www.boost.org>). Environments are isolated from the rest of the system, while still allowing interaction with it: e.g. tools inside the environment are preferred over system tools, while system tools that are not available from within the environment can still be used. Conda also supports the automatic creation of environment definitions from already existing environments. This allows to rapidly explore the needed combination of packages before it is finalized into an environment definition. When used with workflow management systems like Galaxy (Afgan et al., 2016), bcbio-nextgen (<https://github.com/chapmanb/bcbio-nextgen>), and Snakemake (Köster and Rahmann, 2012) that interact directly with Conda, a data analysis can be shipped and deployed in a fully reproducible way, from description and automatic execution of every analysis step down to the description and automatic installation of any required software.

References

- E Afgan, D Baker, den Beek M van, D Blankenberg, D Bouvier, M Čech, J Chilton, D Clements, N Coraor, C Eberhard, B Grüning, A Guerler, J Hillman-Jackson, Kuster G Von, E Rasche, N Soranzo, N Turaga, J Taylor, A Nekrutenko, and J Goecks. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res*, 44:W3–W10, Jul 2016. doi: 10.1093/nar/gkw343. URL <https://doi.org/10.1093/nar/gkw343>.
- Monya Baker. 1,500 scientists lift the lid on reproducibility. *Nature*, 533(7604):452–454, may 2016. doi: 10.1038/533452a. URL <https://doi.org/10.1038%2F533452a>.
- NL Bray, H Pimentel, P Melsted, and L Pachter. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol*, 34:525–7, May 2016. doi: 10.1038/nbt.3519. URL <https://doi.org/10.1038/nbt.3519>.
- F da Veiga Leprevost, BA Grüning, Afitos S Alves, HL Röst, J Uszkoreit, H Barsnes, M Vaudel, P Moreno, L Gatto, J Weber, M Bai, RC Jimenez, T Sachsenberg, J Pfeuffer, Alvarez R Vera, J Griss, AI Nesvizhskii, and Y Perez-Riverol. BioContainers: an open-source and community-driven framework for software standardization. *Bioinformatics*, 33:2580–2582, Aug 2017. doi: 10.1093/bioinformatics/btx192. URL <https://doi.org/10.1093/bioinformatics/btx192>.
- Dawn Field, Bela Tiwari, Tim Booth, Stewart Houten, Dan Swan, Nicolas Bertrand, and Milo Thurston. Open software for biologists: from famine to feast. *Nature Biotechnology*, 24(7):801–803, jul 2006. doi: 10.1038/nbt0706-801. URL <https://doi.org/10.1038%2Fnb0706-801>.
- GM Kurtzer, V Sochat, and MW Bauer. Singularity: Scientific containers for mobility of compute. *PLoS One*, 12:e0177459, 2017. doi: 10.1371/journal.pone.0177459. URL <https://doi.org/10.1371/journal.pone.0177459>.
- J Köster and S Rahmann. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics*, 28:2520–2, Oct 2012. doi: 10.1093/bioinformatics/bts480. URL <https://doi.org/10.1093/bioinformatics/bts480>.
- MI Love, W Huber, and S Anders. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol*, 15:550, 2014. doi: 10.1186/s13059-014-0550-8. URL <https://doi.org/10.1186/s13059-014-0550-8>.
- J. P. Mesirov. Accessible Reproducible Research. *Science*, 327(5964):415–416, jan 2010. doi: 10.1126/science.1179653. URL <https://doi.org/10.1126%2Fscience.1179653>.
- Marcus R. Munafò, Brian A. Nosek, Dorothy V. M. Bishop, Katherine S. Button, Christopher D. Chambers, Nathalie Percie du Sert, Uri Simonsohn, Eric-Jan Wagenmakers, Jennifer J. Ware, and John P. A. Ioannidis. A manifesto for reproducible science. *Nature Human Behaviour*, 1(1):0021, jan 2017. doi: 10.1038/s41562-016-0021. URL <https://doi.org/10.1038%2Fs41562-016-0021>.
- R Patro, G Duggal, MI Love, RA Irizarry, and C Kingsford. Salmon provides fast and bias-aware quantification of transcript expression. *Nat Methods*, 14:417–419, Apr 2017. doi: 10.1038/nmeth.4197. URL <https://doi.org/10.1038/nmeth.4197>.
- E Talevich, AH Shain, T Botton, and BC Bastian. CNVkit: Genome-Wide Copy Number Detection and Visualization from Targeted DNA Sequencing. *PLoS Comput Biol*, 12:e1004873, Apr 2016. doi: 10.1371/journal.pcbi.1004873. URL <https://doi.org/10.1371/journal.pcbi.1004873>.