# bioaRchive: Improving reproducibility of Bioconductor analyses in Galaxy

**Nitesh Turaga**, Eric Rasche, Enis Afgan, Dannon Baker, Galaxy Team

July 7-8th 2015
GCC, Norwich 2015

# As a community we should be focusing on completely reproducible analysis.

PLOS | COMPUTATIO BIOLOGY

**Editorial**

## Ten Simple Rules for Reproducible Computational Research

Geir Kjetil Sandve[1,2]*, Anton Nekrutenko[3], James Taylor[4], Eivind Hovig[1,5,6]

1 Department of Informatics, University of Oslo, Blindern, Oslo, Norway, 2 Centre for Cancer Biomedicine, University of Biochemistry and Molecular Biology and The Huck Institutes for the Life Sciences, Penn State University, University 4 Department of Biology and Department of Mathematics and Computer Science, Emory University, Atlanta, Georgia, Un Biology, Institute for Cancer Research, The Norwegian Radium Hospital, Oslo University Hospital, Montebello, Oslo, N Norwegian Radium Hospital, Oslo University Hospital, Montebello, Oslo, Norway

Replication is the cornerstone of a cumulative science [1]. However, new tools and technologies, massive amounts of data, interdisciplinary approaches, and the complexity of the questions being asked are complicating replication efforts, as are increased pressures on scientists to advance their research [2]. As full replication of studies on independently collected data is often not feasible, there has recently been a call for reproducible research as an attainable minimum standard for assessing the value of scientific claims [3]. This requires that papers in experimental science describe the results and provide a

We further note that reproducibility is just as much about the habits that ensure reproducible research as the technologies that can make these processes efficient and realistic. Each of the following ten rules captures a specific aspect of reproducibility, and discusses what is needed in terms of information handling and tracking of procedures. If you are taking a bare-bones approach to bioinformatics analysis, i.e., running various custom scripts from the command line, you will probably need to handle each rule explicitly. If you are instead performing your analyses through an integrated framework (such as Gene-

### Rule 3: Archive the Exact Versions of All External Programs Used

In order to exactly reproduce a result, it may be necessary to use prog in the exact versions used originally. as both input and output formats change between versions, a newer ve of a program may not even run wi modifying its inputs. Even having which version was used of a program, it is not always trivial t hold of a program in anything bu current version. Archiving the exac sions of programs actually used may save a lot of hassle at later stages. In

Bioconductor is a highly used bioinformatics software toolkit.

BUT ….

ONLY the most recent version of any tool is available to the user.

# What is bioaRchive?

Repository of all versions of all Bioconductor packages and these can be easily obtained from

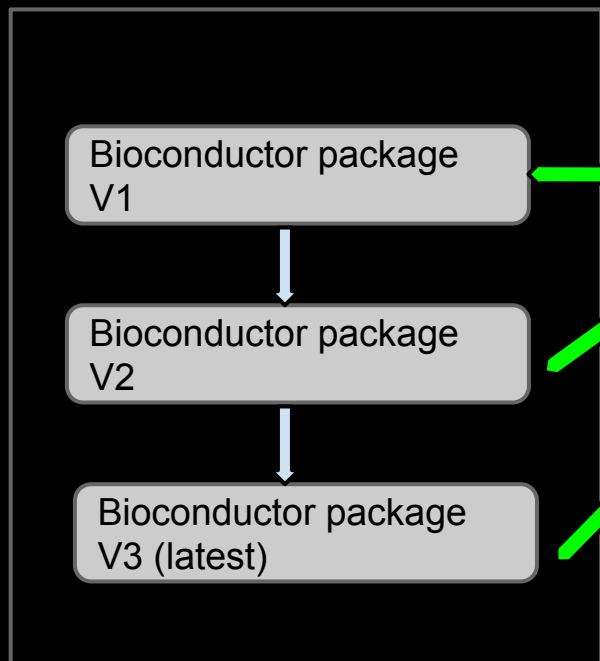bioarchive.galaxyproject.org

Thanks to Eric Rasche for the UI.

# Using bioaRchive

Install version of **Biobase 2.29.0** directly from **bioaRchive**.

```r
install.packages(
  "https://bioarchive.galaxyproject.org/Biobase_2.29.0.tar.gz",
  repos=NULL,method="libcurl")
library("Biobase")
sessionInfo()
```

# Also facilitating Bioconductor and Galaxy interoperability

# Future work

Improve dependency management for bioconductor based analysis

- Missing versions of dependencies
- Packages with multiple dependencies.

# Want to help?

[bioarchive.github.io](bioarchive.github.io)

Contact:

nitesh.turaga@gmail.com