

# Singularity Containers in Bioinformatics

**Vipin T Sreedharan**

Computing Infrastructure @NEXUS

Computational Reproducibility Seminar

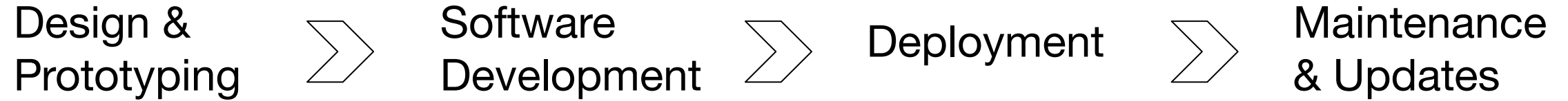
18 October 2023



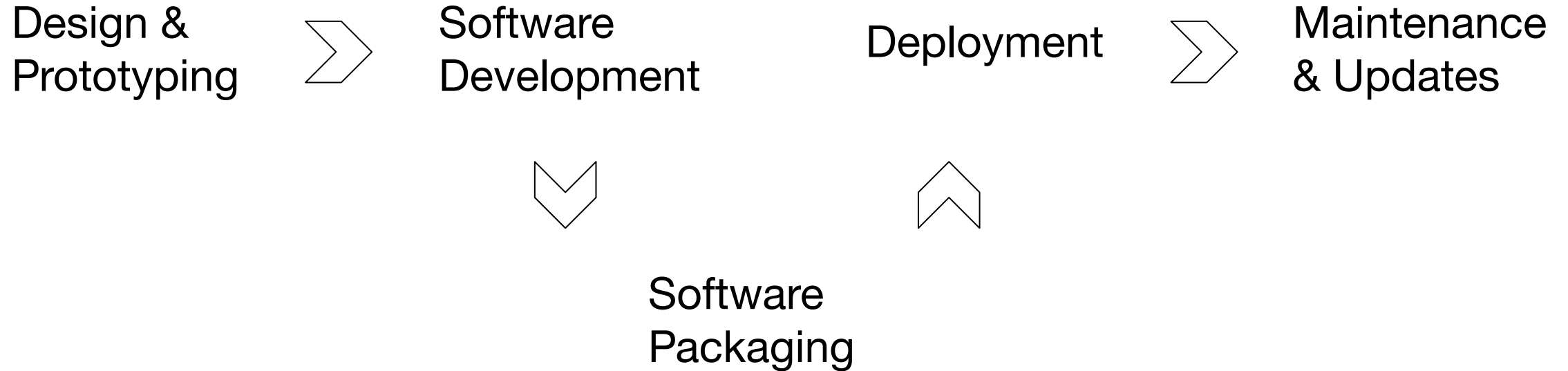
# Objectives

- Software packaging
- Galaxy platform
- Container images for bioinformatics tools

# Software Deployment Process

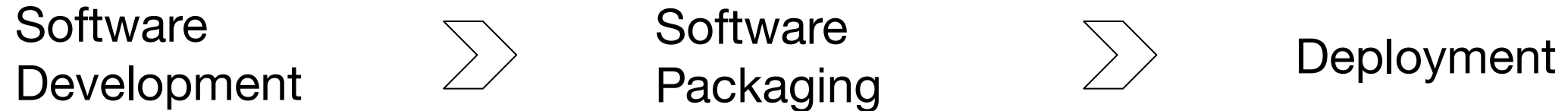


# Software Deployment Process





# Software Deployment Process



## Package Manager

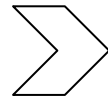
- Package manager performs dependency resolution checks and install
- In **1993**, the earliest form of package manager began to appear from Linux
- Some of these early package managers (**dpkg**, **rpm**) live on today

# Software Deployment before Package Managers

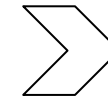
- Historically, software was provided either via FTP or basic websites
- The **configure** process starts using a C compiler and checks your system for application dependencies
- If the configure script completed successfully, a **Makefile** would be created
- Once a **Makefile** existed, you would then proceed to run the **make** command
- Finally, after the **make** process has been completed, you would need to run **make install** in order to actually install the software

# Software Deployment Process

Software  
Development



Software  
Packaging



Deployment

## 32 Package Managers



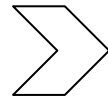
Language package managers



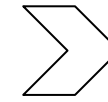
System package managers

# Software Deployment Process

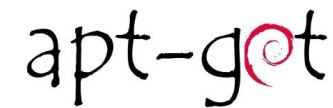
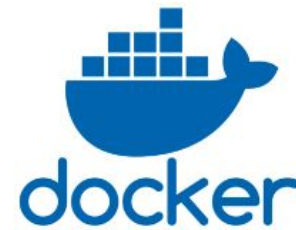
Software  
Development



Software  
Packaging



Deployment





Galaxy is an **open-source** platform for FAIR **data analysis** that enables users to:

- Use **tools** from various domains (that can be plugged into **workflows**) through its graphical web interface
- Run code in **interactive environments** (RStudio, Jupyter...) along with other tools or workflows
- **Manage data** by sharing and publishing results, workflows, and visualizations
- **Ensure reproducibility** by capturing the necessary information to repeat and understand data analyses

# Galaxy User Interface



The screenshot shows the Galaxy web interface. On the left is a 'Tools' sidebar with a search bar and a list of tool categories including 'Get Data', 'GENERAL TEXT TOOLS', 'GENOMIC FILE MANIPULATION', 'COMMON GENOMICS TOOLS', and 'Assembly'. The main content area features an announcement for the 'James P. Taylor (JXTX) Foundation for Open Science' with a video thumbnail and a 'Learn More' button. Below the announcement is a banner for SARS-CoV-2 data analysis. On the right is a 'History' panel showing a search bar and a list of datasets, including '2: SNPs' and '1: Exons'. The browser address bar shows 'usegalaxy.org'.

## Public servers



**Galaxy Tool Shed** Repositories Groups Help User

9588 valid tools on Sep 24, 2023

**Search**  
Search for valid tools

**Valid Galaxy Utilities**  
Tools  
Custom datatypes  
Repository dependency definitions  
Tool dependency definitions

**All Repositories**  
Browse by category

**Available Actions**  
Login to create a repository

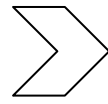
### Repositories by Category

**Name** **Description** **Repositories**

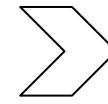
Assembly	Tools for working with assemblies	198
Astronomy	Tools for astronomy	8
ChIP-seq	Tools for analyzing and manipulating ChIP-seq data.	78
Climate Analysis	Tools for analyzing climate data	12
CLIP-seq	Tools for CLIP-seq	4
Combinatorial Selections	Tools for combinatorial selection	9
Computational chemistry	Tools for use in computational chemistry	180
Constructive Solid Geometry	Tools for constructing and analyzing 3-dimensional shapes and their properties	11
Convert Formats	Tools for converting data formats	140
Data Export	Tools for exporting data to various destinations	17
Data Managers	Utilities for Managing Galaxy's built-in data cache	106
Data Source	Tools for retrieving data from external data sources	107
Ecology	Tools related to ecological studies	74
Entomology	Tools that involve insect studies	4
Epigenetics	Tools for analyzing Epigenetic/Epigenomic datasets	48
Fasta Manipulation	Tools for manipulating fasta data	121
Fastq Manipulation	Tools for manipulating fastq data	104
Flow Cytometry Analysis	Tools for manipulating and analyzing FCS files	45

# Software Sustainability in Science

Software  
Development



Software  
Packaging



Deployment



apt-get



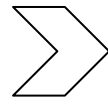
# Software Sustainability in Science

## Standards needed for packaging bioinformatics software

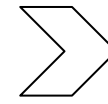
- Programming language
- OS independent
- Multiple versions of the software
- HPC and Cloud compatible
- easy to maintain

# Software Sustainability in Science

Software  
Development



Software  
Packaging



Deployment



BIOCONDA



apt-get



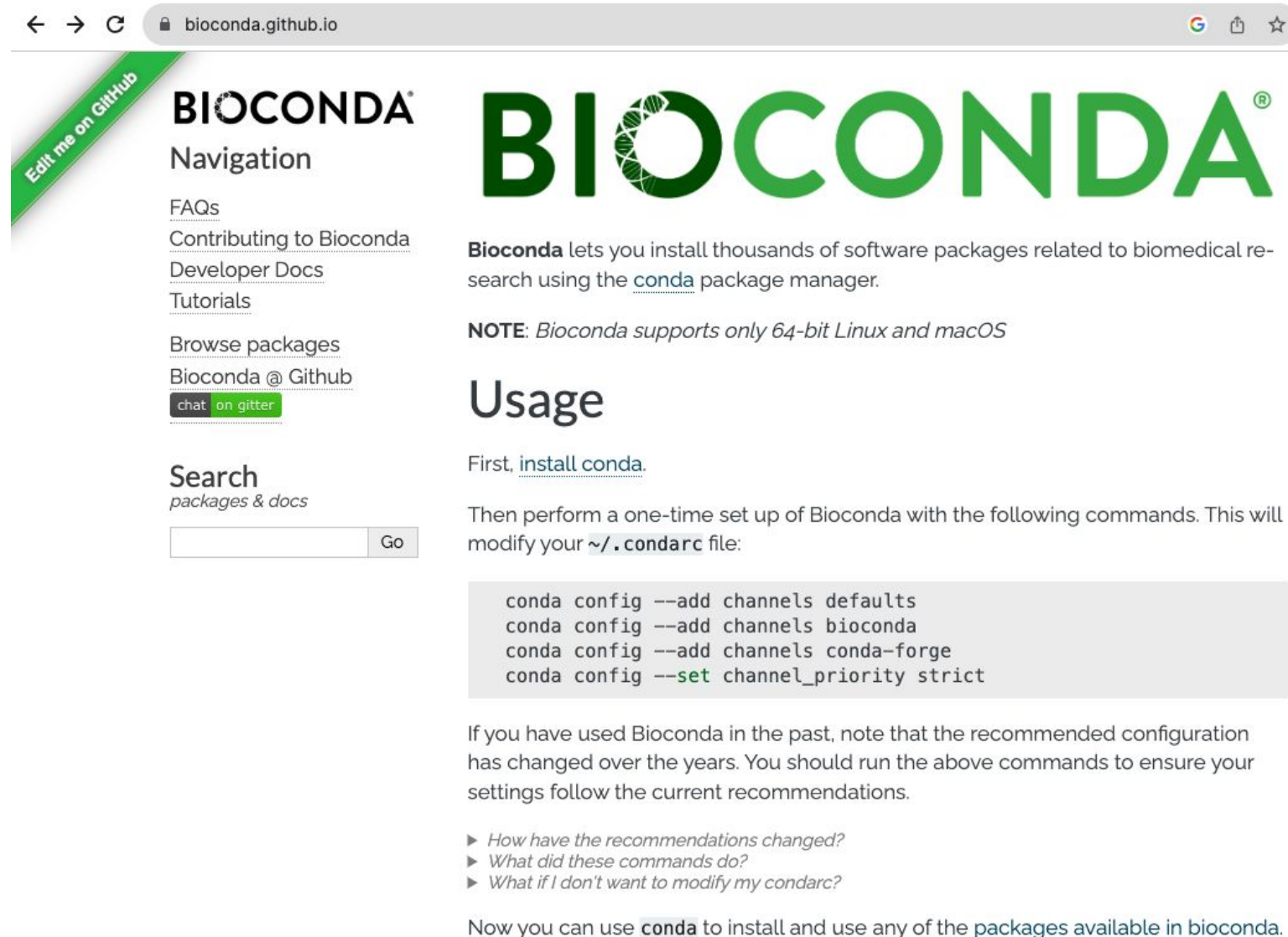


# Software Sustainability in Science

Code Blame 35 lines (29 loc) · 713 Bytes

```
1  {% set version = "0.2.2" %}
2
3  package:
4    name: ribodiff
5    version: {{ version }}
6
7  source:
8    url: https://github.com/ratschlab/RiboDiff/archive/v0.2.2.tar.gz
9    md5: b02833b4412959746032f0cf23a944d8
10
11 build:
12   noarch: python
13   number: 2
14
15 requirements:
16   host:
17     - python <3
18     - pip
19   run:
20     - python <3
21     - numpy >=1.8.0
22     - scipy >=0.13.3
23     - matplotlib >=1.3.0
24     - statsmodels >=0.5.0
25
26 test:
27   imports:
28     - ribodiff
29   commands:
30     - TE.py --help
31
32 about:
33   home: http://public.bmi.inf.ethz.ch/user/zhongy/RiboDiff/index.html
34   license: GPL 3
35   summary: 'RiboDiff is a statistical tool that detects the protein translational efficiency change from Ribo-Seq (ribosome footprinting) and RNA-Seq data.'
```

# Software Sustainability in Science



The screenshot shows the Bioconda website homepage. At the top, there is a browser address bar with the URL 'bioconda.github.io'. Below the address bar, on the left, is a navigation menu with links for 'FAQs', 'Contributing to Bioconda', 'Developer Docs', 'Tutorials', 'Browse packages', and 'Bioconda @ Github'. A search bar is located below the navigation menu. On the right side of the page, the Bioconda logo is prominently displayed. Below the logo, there is a brief description of Bioconda as a package manager for biomedical research. A note specifies that Bioconda supports only 64-bit Linux and macOS. The 'Usage' section provides instructions on how to install conda and configure it to use Bioconda channels. A code block shows the specific conda commands for this configuration. Finally, there are three bullet points with links to more information about the configuration changes and how to modify the conda file.

← → ↻ bioconda.github.io

**BIOCONDA**  
Navigation

FAQs  
Contributing to Bioconda  
Developer Docs  
Tutorials  
Browse packages  
Bioconda @ Github  
chat on gitter

**BIOCONDA**  
®

**Bioconda** lets you install thousands of software packages related to biomedical research using the `conda` package manager.

**NOTE:** *Bioconda supports only 64-bit Linux and macOS*

## Usage

First, [install conda](#).

Then perform a one-time set up of Bioconda with the following commands. This will modify your `~/.condarc` file:

```
conda config --add channels defaults
conda config --add channels bioconda
conda config --add channels conda-forge
conda config --set channel_priority strict
```

If you have used Bioconda in the past, note that the recommended configuration has changed over the years. You should run the above commands to ensure your settings follow the current recommendations.

- ▶ [How have the recommendations changed?](#)
- ▶ [What did these commands do?](#)
- ▶ [What if I don't want to modify my condarc?](#)

Now you can use `conda` to install and use any of the [packages available in bioconda](#).

# Software Sustainability in Science

## Large scale bioinformatics data analysis

- Bioconda standard package management
- Different workflow managers
- Success story of the bioconda community
- easy to maintain



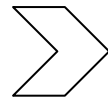
**nextflow**

# Software Sustainability in Science

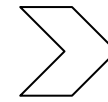


# Software Sustainability in Science

Software  
Development



Software  
Packaging



Deployment



conda  
18.8K Packages

BIOCONDA



docker



BioContainers

Galaxy Term - **Mulled**

- without Dockerfile
- layer donning approach for building containers

# Software Sustainability in Science

- 9874 bioinformatics packages
- 95404 singularity images



# Software Sustainability in Science

```
profiles {
```

nextflow

```
  samtools {
```

```
    process.container = 'https://depot.nexus.ethz.ch/singularity/samtools:1.2'
```

```
    singularity.enabled = true
```

```
    singularity.cacheDir = "$PWD"
```

```
  }
```

```
}
```

```
rule samtools:
```

```
  input:
```

```
    "inputs/{unmapped}.bam"
```

```
  output:
```

```
    "results/{unmapped}.bam"
```

```
  threads: 1
```

```
  singularity: "https://depot.nexus.ethz.ch/singularity/samtools:1.2"
```

```
  shell:
```

```
    """
```

```
    cat {input} > {output}
```

```
    """
```

```
    "samtools view {input} > {output}"
```

snakemake

# Software Sustainability in Science



# Software Sustainability in Science

The screenshot shows a web browser window with the address bar containing 'depot.nexus.ethz.ch'. The page title is 'Image Depot Search Interface (IDSI)'. The search bar contains the text 'samtools:1.2'. Below the search bar, it indicates 'Type to search' and 'Number of results: 9'. The search results are listed as follows:

- bioconductor-rsamtools:1.22.0-r3.2.2\_1
- bioconductor-rsamtools:1.26.1-r3.3.1\_0
- bioconductor-rsamtools:1.26.1-r3.3.2\_0
- bioconductor-rsamtools:1.26.1-r3.4.1\_0
- bioconductor-rsamtools:1.28.0-r3.4.1\_0
- samtools:1.2
- samtools:1.2-0
- samtools:1.2-0
- samtools:1.2.rglab-0

# Acknowledgments



- NEXUS Software Engineering
- HPC Storage Service
- S4D Support Service
- ETH Gitlab Service