# Secure and practical computational reproducibility in the life sciences

## Bioconda, BioContainers, Galaxy & the de.NBI Galaxy-Docker-Technology

Rolf Backofen, Björn Grüning & The RBC Team

# Portfolio of Tools and S[...]



Since 2013, IF >10

- 9x Nature        3x Science
- 4x Cell          1x Nat. Str. Mol Bio
- 5x Nat. Com.     1x Cell Stem Cell
- 1x Nat. Med.     1x Nat. Micro. Review
- 2x Nat. Meth.    5x Nat. Genetics
- 2x Nat. Biotech. 1x Nat. Rev.Immunol.
- 10x Mol. Cell    16x PNAS,
- 4x Gen. Res.     6x Gen. Bio.
- 1x Circ. Res.    1x J. Clinic. Investig.
- 1x Genes Dev.    6x Embo J.

# Service

- knowledge transfer in RNA bioinformatics



- we can't do all analysis ourselves!
- solutions:
  - **Galaxy**
    - *standard workflows*
    - *FR server: 600+ users*
  - **virtualization**
    - *distribute computation*
  - **training, training, training**
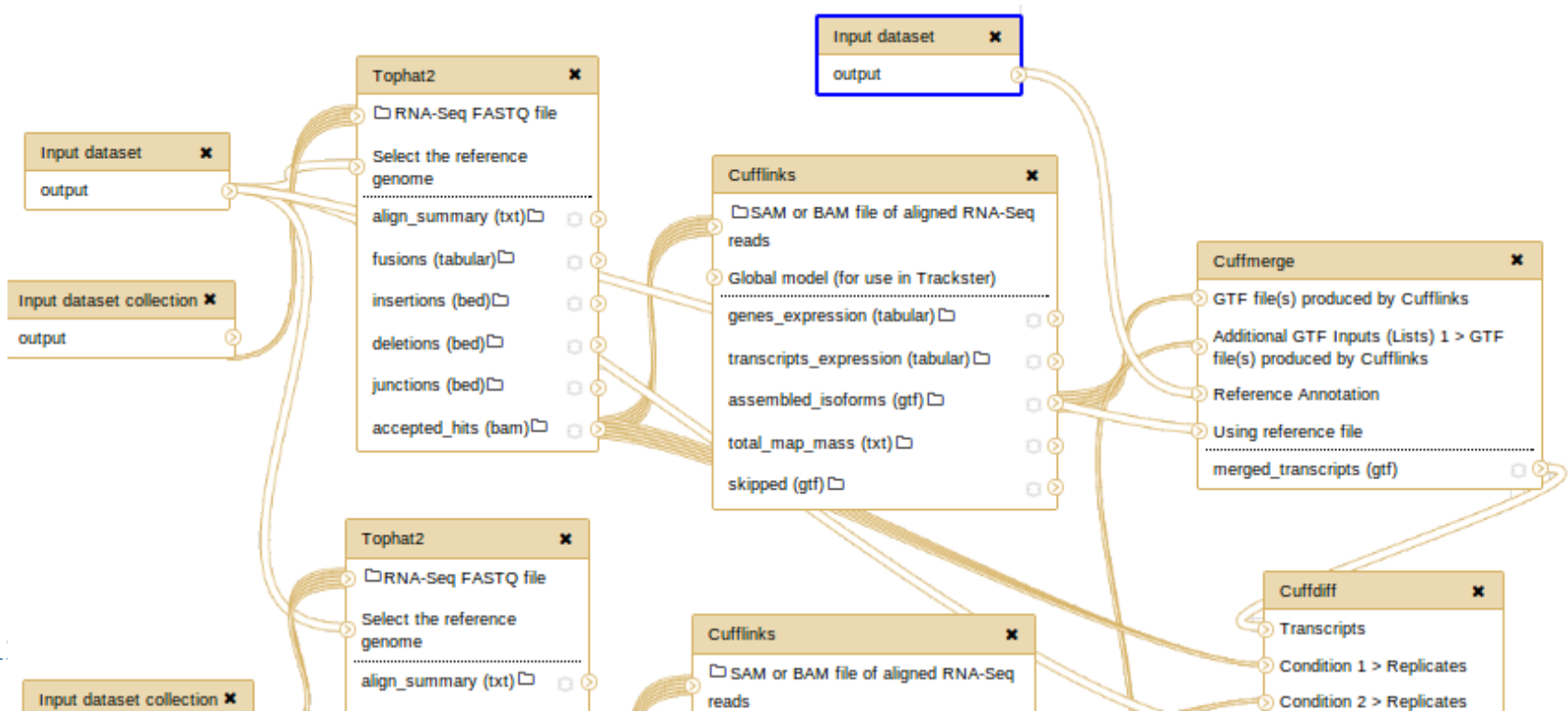
# Galaxy: Accessible Research

☐ reproducibility

- experimental details

- transparency

☐ scalable

☐ easy deployment

- tools available to users

- minimal installation overhead
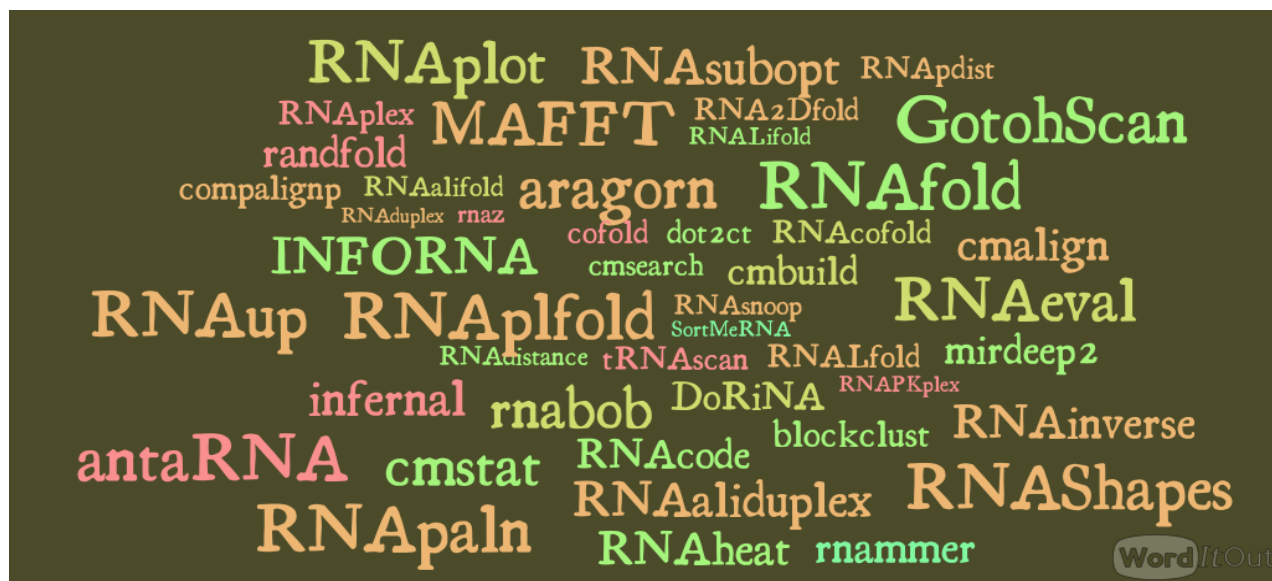
# What is Galaxy?

- Omics data analysis platform
- Accessible via normal web browser

- Every single step is recorded and reproducible
- 5000 citations so far
- 150.000 known user

ocr

# Galaxy RNA Workbench

- already integrated: >4000 tools, ~800 by RBC
  - *We are one of 3 groups worldwide with direct Galaxy commit*

- RNA-specific tools and packages: **60+**



- RBC-specific:
  - *Vienna RNA package, doRiNA, Freiburger RNA Tools*

# Services of Freiburg Galaxy Server

## Workflows

### General
- Build, Test, Use, Share, FAIR

### RNA-seq
- Quality Control
- Differential Expression Analysis

### ChIP-seq, MethylC-seq

### Additionally
- Exome-seq
- Proteomics, Metabolomics
- Imaging, Textmining

# Impact Measurements

- users

- jobs run

# Building a sustainable virtualized infrastructure

# Building a secure virtualized infrastructure



- Isolate every single tool
  - from all other data
  - from other tools
  - from the Operating System
- Isolate the analysis workbench

# adjustable reproducibility and security

**Least reproducible** → **Most reproducible**

**Easy** ← **Hard / Impossible ?**

**Less secure** → **Most secure**

**Easy** ← **Hard / Impossible ?**

# Linux **Containers** for advanced isolation



If there is no standard*, embrace new technologies and make them interconvertible.

* Open Container Initiative (OCI) has released 1.0 of the container and runtime specification 3 days ago.

# Sustainable community-based infrastructure

**BIOCONDA**

provides software for biomedical research.

- 14,400 commits on GitHub
- 287 contributors
- >2,700 packages

**Biocontainers**

provides system-agnostic executable environments for bioinformatics tools

- Uses Docker & rkt
- >2,800 Images
- Automatic builds from BioConda

**Galaxy**

Galaxy is an open web-based platform for data intensive research

- 29,500 commits
- 163 contributors
- RBC is one of the biggest instances

**Our achievement: Cloud-Ready, all integrated into Galaxy**

# Biocontainers

- build out of **conda** package

- all tested automatically

- bioconda-utils / galaxy-lib



By Ryan Dale

# Community



by Johannes Köster

**330 Contributors        6200 merged PR**
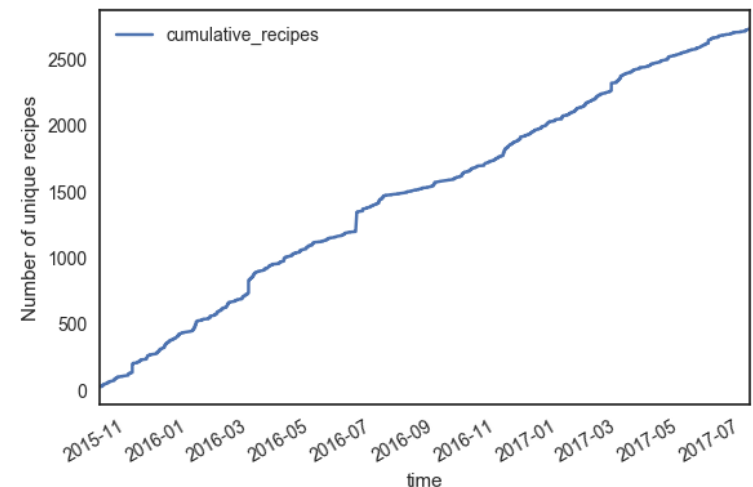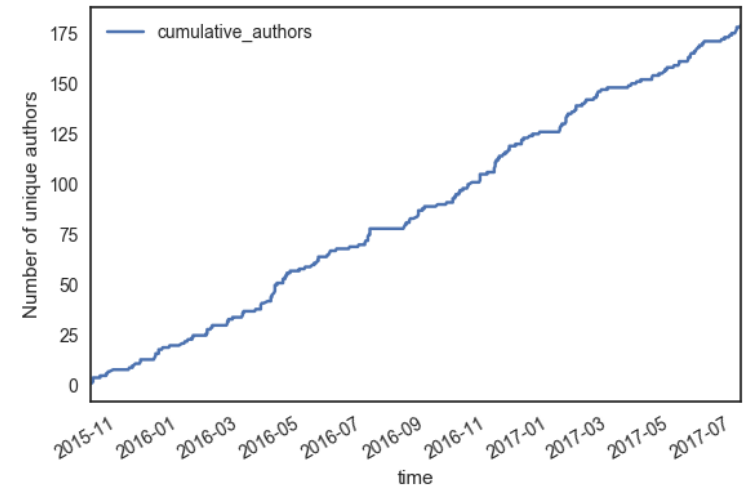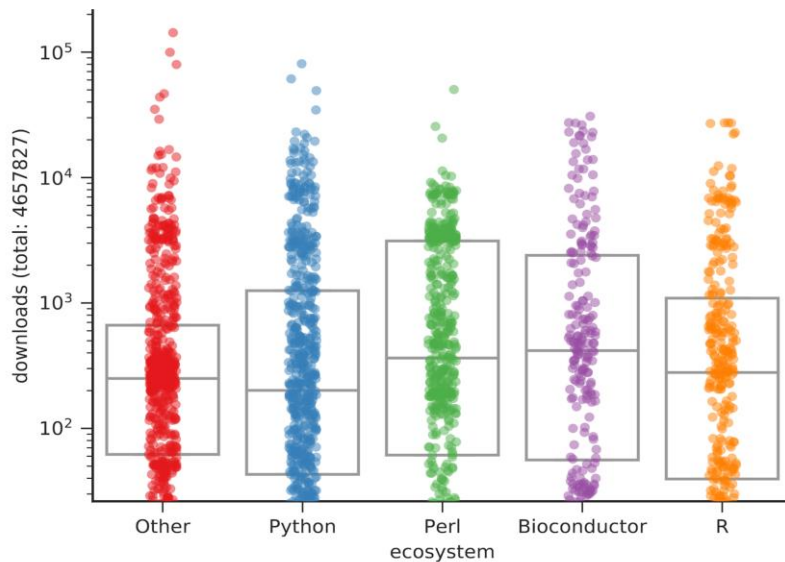
in 2 years

# Reproducibility stack

de.NBI
GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE

## 1. Conda
## 2. Containers
## 3. Virtualization

| OS | OS | OS | OS |
|---|---|---|---|

**Virtual Machine**

| Trimmomatic | Conda package<br>Trimmomatic | Conda package<br>Trimmomatic — Container | Conda package<br>Trimmomatic — Container |
| HISAT2 | Conda package<br>HISAT2 | Conda package<br>HISAT2 — Container | Conda package<br>HISAT2 — Container |
| StringTie | Conda package<br>StringTie | Conda package<br>StringTie — Container | Conda package<br>StringTie — Container |
| Ballgown | Conda package<br>Ballgown | Conda package<br>Ballgown — Container | Conda package<br>Ballgown — Container |
| local environment | local environment | local environment | local environment |

**Least reproducible / secure** → **Most reproducible / secure**

```
> git clone hisat2
> make
> sudo make install
> hisat2 --version
```

```
> conda install hisat2
> hisat2 --version
```

```
> docker run --rm
quay.io/biocontainers/
hisat2 --version
```

# Galaxy in Docker

- Idea: tools to data, not data to tools -> virtualization



- drag & drop based Galaxy flavor generator
- Widely used: >28.000 downloads, ELIXIR (Tjenester for Sensitive Data 2.0 Norway), Cancer Center Amsterdam …

# Galaxy in Docker: Microservices



Isolate the analysis workbench!

- Every single component can be isolated and hardened
- Microservices that communicate with each other
- Can run in an isolated network without Internet
- executes analysis jobs in own isolated containers

# Real world deployments with sensitive data:  Version 1

- central Galaxy server
- job submission to a hardened HTCondor pool
- Galaxy server runs in Docker
- every single job runs in Docker
- every job is isolated from all other data
- one job has only read access to the input data it needs
- entire stack can run in an isolated network

# Real world deployments with sensitive data:  Version 2

- 2 factor Authentication into virtual machines (VM)
- VMs are isolated
- every VM starts its own Galaxy instance
- data is saved encrypted to the local network
- Users can not interact with each other (share workflows, histories etc ...)
- every VM can run Version 1 as well

New Results

### Practical computational reproducibility in the life sciences

Björn Grüning, John Chilton, Johannes Köster, Ryan Dale, Jeremy Goecks, Rolf Backofen, Anton Nekrutenko, James Taylor

doi: https://doi.org/10.1101/200683

New Results

### Bioconda: A sustainable and comprehensive software distribution for the life sciences
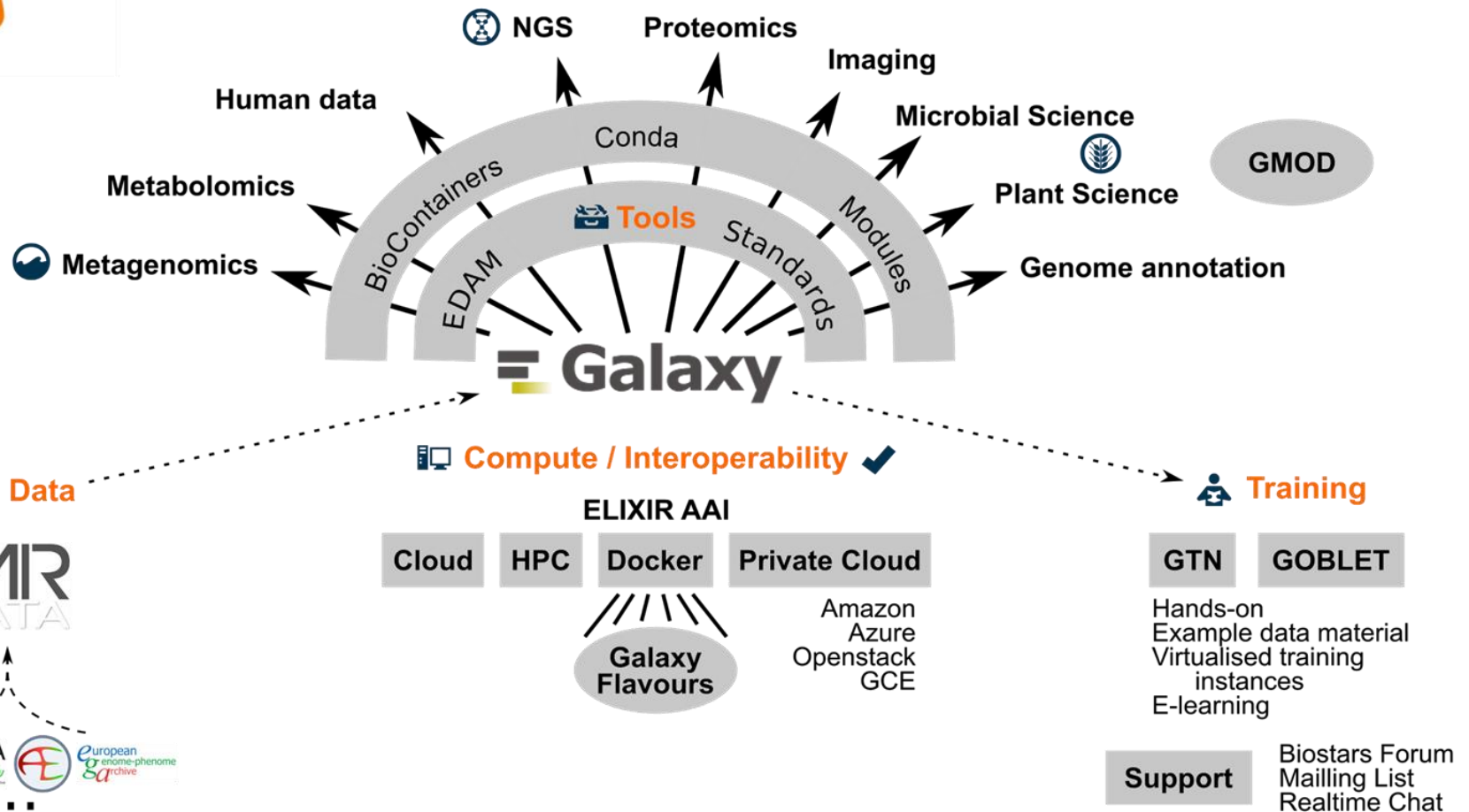
Björn Grüning, Ryan Dale, Andreas Sjödin, Jillian Rowe, Brad A. Chapman, Christopher H. Tomkins-Tinch, Renan Valieris, The Bioconda Team, Johannes Köster

doi: https://doi.org/10.1101/207092

New Results

### Community-driven data analysis training for biology

Bérénice Batut, Saskia Hiltemann, Andrea Bagnacani, Dannon Baker, Vivek Bhardwaj, Clemens Blank, Anthony Bretaudeau, Loraine Guéguen, Martin Čech, John Chilton, Dave Clements, Olivia Doppelt-Azeroual, Anika Erxleben, Mallory Freeberg, Simon Gladman, Youri Hoogstrate, Hans-Rudolf Hotz, Torsten Houwaart, Pratik Jagtap, Delphine Lariviere, Gildas Le Corguillé, Thomas Manke, Fabien Mareuil, Fidel Ramírez, Devon Ryan, Florian Sigloch, Nicola Soranzo, Joachim Wolff, Pavankumar Videm, Markus Wolfien, Aisanjiang Wubuli, Dilmurat Yusuf, Rolf Backofen, Anton Nekrutenko, Björn Grüning

doi: https://doi.org/10.1101/225680

# Conclusion

- **Galaxy**
  - Every single step is recorded and reproducible
  - 5000 citations so far, 150.000 known user
  - Freiburg Server: largest in Europe with 600 users.

- **Deployments with Sensitive data.**
  - Reproduciblity stack:

    *OS -> Packages -> Containers -> VM*
  - Galaxy Docker Flavour Concept
  - Community-based effort:

    *Bioconda, Biocontainers & Galaxy*

Thanks!

Galaxy PROJECT

Yasset Perez-Riverol
Felipe Leprevost ...

Björn Grüning & the RBC team

Johannes Köster
Ryan Dale ...

# Thank you for your attention

# Impact Measurements

- users

- jobs run



| Measurement (from SIG2 report) | Total |
|---|---:|
| Citations/links to website/acknowledgements | 9,589 |
| Supporttickets / GitHub issues | 4,942 |
| Number of downloads | 3,474,974 |
| Web app & database hits | 170,411 |

# Overcoming computational limitations

- Idea: tools to data, not data to tools -> virtualization



- drag & drop based Galaxy flavor generator
- Widely used: >28.000 downloads, ELIXIR (Tjenester for Sensitive Data 2.0 Norway), Cancer Center Amsterdam …

My Appliances

This page is a log of the appliances you have launched. ▾

- ## Usable clouds:
  - Amazon
  - de.NBI-cloud (FR)
  - ....

- ## Our contributions:
  - Virtual container via the flavor generator
  - Can be started world-wide
  - Interface to RBC Galaxy server (user with account can use it)

Cloud: bwcloud

ABOUT        DEVELOPER

# Masterplan



**Training**
- GTN
- Goblet

- Hands-on
- Galaxy Tours
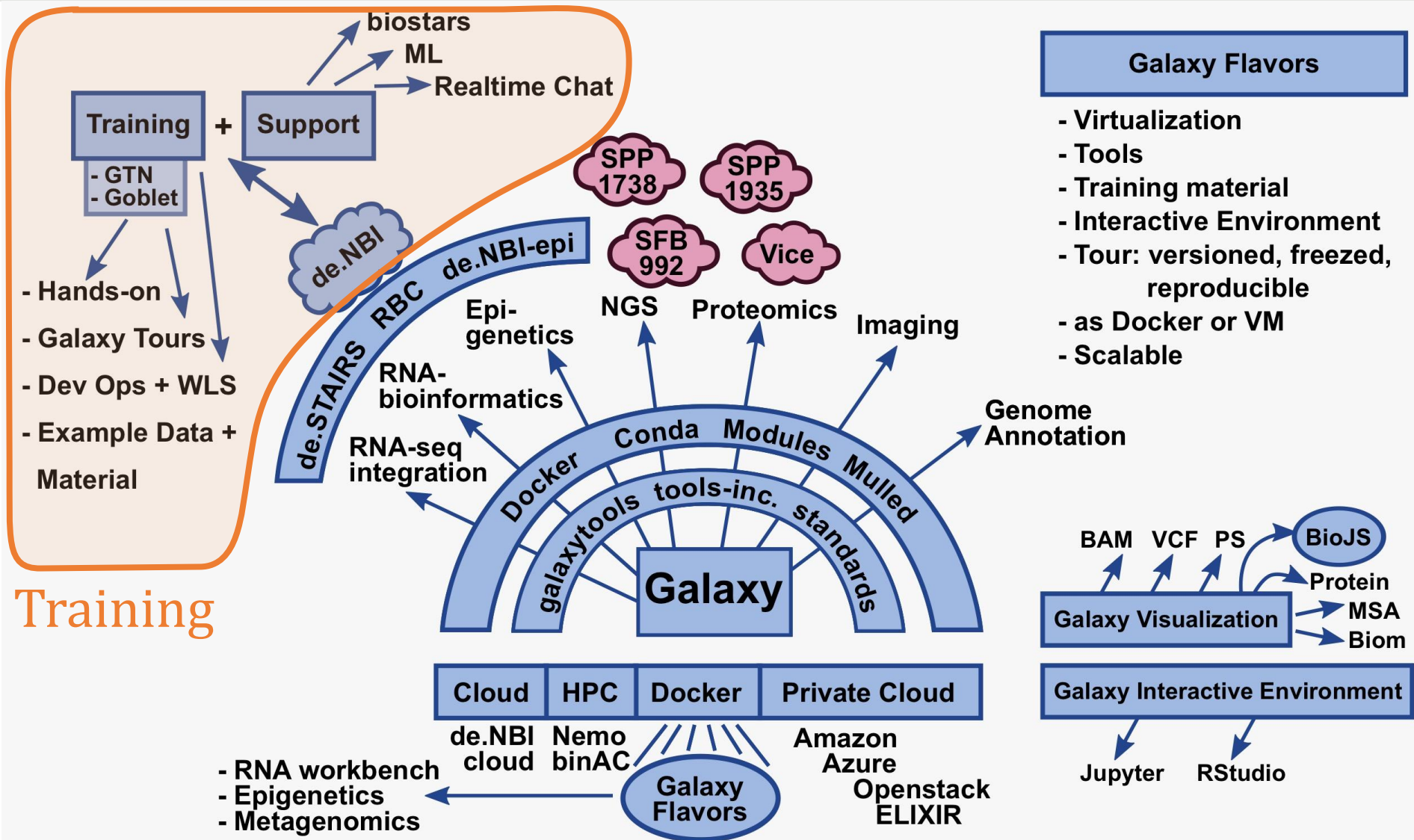- Dev Ops + WLS
- Example Data + Material

**Support** + **Training**

biostars
ML
Realtime Chat

de.NBI

de.STAIRS  RBC  de.NBI-epi

RNA-bioinformatics

RNA-seq integration

Epi-genetics

SPP 1738  SPP 1935

SFB 992  Vice

NGS  Proteomics

Imaging

Docker  Conda  Modules  Mulled

galaxytools  tools-inc.  standards

**Galaxy**

Genome Annotation

Cloud | HPC | Docker | Private Cloud

de.NBI cloud  Nemo binAC  Amazon Azure Openstack ELIXIR

**Galaxy Flavors**

- RNA workbench
- Epigenetics
- Metagenomics

**Galaxy Flavors**

- Virtualization
- Tools
- Training material
- Interactive Environment
- Tour: versioned, freezed, reproducible
- as Docker or VM
- Scalable

BAM  VCF  PS  BioJS
Protein
MSA
Biom

**Galaxy Visualization**

**Galaxy Interactive Environment**

Jupyter  RStudio

# Training by RBC

- ## 28 training courses, ~500 participants

- ## 12 locations across the world

**2016 events**

| | | |
|---|---|---|
| 19th – 20th Jan | ELIXIR EDAM codefest | Freiburg |
| 22nd – 26th Feb | Galaxy HTS data analysis workshop | Freiburg |
| 07th – 08th Mar | RAD-Seq tools and workflows codefest | Online |
| 04th Apr | Conda codefest | Online |
| 06th – 07th Apr | Galaxy DevOps workshop | Heidelberg |
| 27th – 29th Apr | HPC workshop | Oslo, Norway |
| 25th – 29th Jun | Galaxy Community Conference workshop | Indiana |
| 27th Jul | RBC Kick-Off meeting | Freiburg |
| 19th – 23rd Sep | Galaxy HTS data analysis workshop | Freiburg |
| 27th – 28th Sep | GalaxyP codefest | Online |
| 06th – 07th Oct | Galaxy training material codefest | Online |
| 20th – 21st Oct | Swiss German Galaxy workshop | Freiburg |
| 24th – 26th Oct | NETTAB hackathon | Rom, Italy |
| 02nd – 03rd Nov | BioConda codefest | Online |
| 30th Nov – 01st Dec | Galaxy Docker workshop | Barcelona, Spain |
| 01st – 02nd Dec | FAIRDOM/de.NBI Foundry workshop | Frankfurt |
| 15th – 16th Dec | Galaxy RNA-seq data analysis workshop | Freiburg |

**2017 events**

| | | |
|---|---|---|
| 09th – 10th Jan | Galaxy QIIME codefest | Online |
| 16th – 19th Jan | European Galaxy developer workshop | Strasbourg, France |

# Galaxy Tours – Bioinformatics Training 2.0



## Galaxy Tours

- ### Analysis in a box

- ### Real execution of code

- ### User can just use the tour

  ### or change parameters/input

# Thank you for your attention

# Galaxy Tours – Bioinformatics Training 2.0

# De.NBI cloud



- Recently opened: Supercomputer NEMO
    - **15,000** cores, **position 214** in TOP 500
- we got
    - **1,500,000 €** hardware extension for GALAXY
    - Additionally **5% of existing NEMO resources**
    - IMPORTANT: **cloud knowledge of BW-CLOUD**

# Trainings -Material
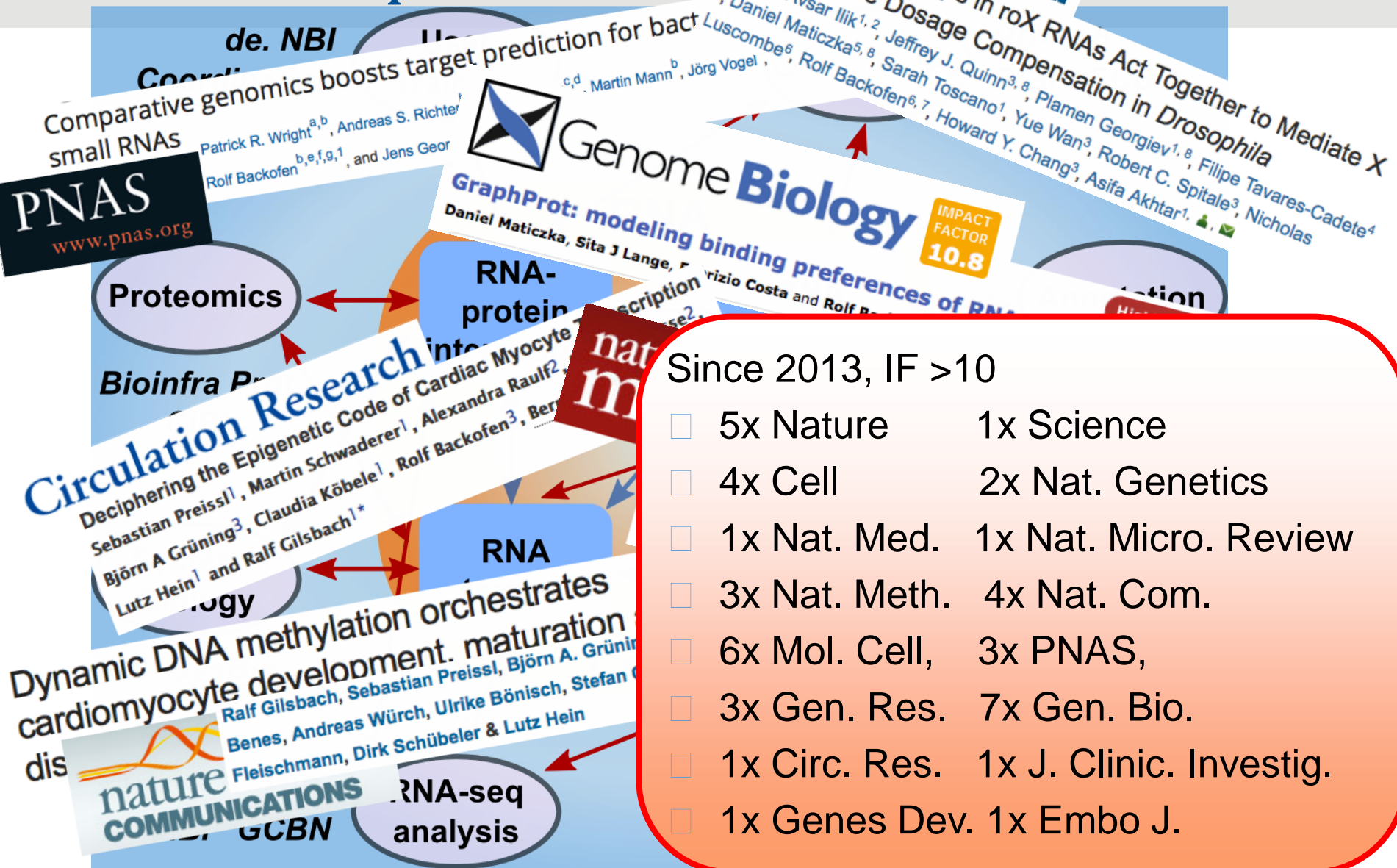
- 11 topics with 68 tutorials for 3 different target audiences

- 51 contributors, 3 contribution feasts
  - Online in October 2016
  - Cambridge in May 2017, organized with ELIXIR and GOBLET
  - Montpellier in June 2017 during the Galaxy Community Conference

- Integration with ELIXIR's training platform (TeSS)

Galaxy Training!  ⊕ Help on Gitter

Epigenetics                                    1

Metagenomics                                   2

# Conclusion

- ## RNA-Bioinformatics Center (RBC):
  - RNA-mediated post-transcriptional regulation
  - integrates all international renowned German RNA bioinformatics groups

- ## Important aspects of our service:
  - RNA workbench based on Galaxy
  - Virtualization for distribution of computational burden
  - Strong interaction with ELIXIR
  - training, training, training (user, user, user)
    *future: eLearning!!*

# Areas of Expertise



Since 2013, IF >10

- 5x Nature          1x Science
- 4x Cell            2x Nat. Genetics
- 1x Nat. Med.    1x Nat. Micro. Review
- 3x Nat. Meth.   4x Nat. Com.
- 6x Mol. Cell,    3x PNAS,
- 3x Gen. Res.    7x Gen. Bio.
- 1x Circ. Res.    1x J. Clinic. Investig.
- 1x Genes Dev. 1x Embo J.