

Running a Local Galaxy

Anushka Brownley
BioTeam Inc.

- Who is BioTeam
- Why Local Galaxy
- Running Galaxy Locally
- Local Galaxy Best Practices

Who is BioTeam

The screenshot shows the BioTeam website homepage. At the top left is the BioTeam logo with the tagline "Enabling Science". A navigation menu includes links for Home, Company, Products, Services, Profiles, News, Articles, Amazon Cloud Training, and Feeds. A search bar is located on the right side of the navigation bar. The main content area features a large featured article titled "INSIDE THE BOX: Solutions for Data Sharing in Life Sciences" with a sub-headline "... modern research produces tons of data and publications are no longer a viable medium for sharing all those data." To the right of this article is a list of other articles, including "Fun and Games with Genomic Analysis At Home", "2013 BioIT Trends Slides", "Intel Managing Big Data Event – Slides", "INSIDE THE BOX: Solutions for Data Sharing in Life Sciences", and "Metrum moves to the cloud". At the bottom of the page, there are four smaller article teasers: "Practical Cloud & Workflow Orchestration", "Mapping Informatics to the Cloud", "Backblaze Storage Pod", and "2011 Trends from the Trenches".

BIOTEAM
Enabling Science

Home Company Products Services Profiles News Articles Amazon Cloud Training Feeds

enter search terms ...

INSIDE THE BOX: Solutions for Data Sharing in Life Sciences
... modern research produces tons of data and publications are no longer a viable medium for sharing all those data.

- Fun and Games with Genomic Analysis At Home
Jul 01, 2013 | 0 comments
- 2013 BioIT Trends Slides
Jun 30, 2013 | 0 comments
- Intel Managing Big Data Event – Slides
Sep 19, 2012 | 0 comments
- INSIDE THE BOX: Solutions for Data Sharing in Life Sciences**
Jul 30, 2012 | 0 comments
- Metrum moves to the cloud
Feb 15, 2012 | 0 comments

Practical Cloud & Workflow Orchestration
Presentation slides from our

Mapping Informatics to the Cloud
Talk slides from a presentation called

Backblaze Storage Pod
102 usable terabytes for \$12,000. Disruptively cheap

2011 Trends from the Trenches
2011 BioITWorld talk slides from Chris Dandician

am The Trenches
2011 Conference & Expo
Software & Hardware

Over a decade of Life Sciences IT consulting

- We are **scientists** forced to learn IT to get research done
- Served over **400** organizations
 - Academic, Non-profit
 - Government, Military
 - Pharm, AgBio, Biotech
 - Cloud & Datacenter Providers



- Active contributors to open-source projects



emboss



GMOD



DIYA



BioPerl

BioPerl



O|B|F



Homebrew
The missing package manager for OS X



GRID ENGINE

•BioTeam

Encapsulate IT best-practices expertise to eliminate redundant effort spent designing and building infrastructure

•Galaxy Project

Decrease the barrier to entry into data analysis by improving accessibility of the Galaxy platform

- BioTeam offers a all-in-one solution to help run a local instance Galaxy



- BioTeam is the the official appliance provider for Galaxy
 - Exclusive partnership with the Galaxy Team
 - Donations back to the Galaxy Project

FEATURES

Powerful Server	Complete analysis tasks quickly
Optimized Galaxy	Configured for stable production use
Analysis Tools	Provides comprehensive set of open-source tools
Automated Updates	Software can be updated automatically
Preinstalled Datasets	Includes 5 model organisms (additional upon request)
Open Platform	Install other software for your own use

Price: \$19,995 (USD)

Why a Local Installation

wiki.galaxyproject.org/BigPicture/Choices

Galaxy Wiki anushkanet Settings Logout | Search:

BigPicture/Choices Edit

Galaxy is available in several different ways.

Which Option to Choose?

Your choices depends upon your needs. Here are the options depending on what you need:

	Main	Local	Cloud	Appliance	Other
Your data sets are moderately sized	Yes	Yes	Yes	Yes	?
Your computational requirements are moderate	Yes	Yes	Yes	Yes	?
You want to share your Galaxy objects with others	Yes	Yes	Yes	Yes	?
All needed Tools are installed on Main .	Yes	?	Yes	Yes	?
Your data sets are very large	No	?	Yes	Yes	?
Your computational requirements are very large	No	?	Yes	Yes	?
You have absolute data security requirements	No	Yes	Yes	Yes	?
No network transfer of data	No	Yes	No	Yes	Yes

- Galaxy Main is a fantastic public resource!!
- Limitations due to popularity
 - Wait times
 - Job and storage quotas
 - Data transfer bottlenecks
 - Pre-defined set of tools and datasets

- Cloud is a scalable option but has its challenges
 - Understand how to manage pricing
 - Data transfer bottlenecks
 - Familiarity with cloud (e.g. Amazon Web Services)
 - Dealing with sensitive data

Running Galaxy locally solves these issues

- IT/informatics expertise
 - Acquire and set up infrastructure
 - Install Galaxy, tools, necessary dependencies
 - Optimize/customize for your use cases
- Define policies
 - Managing usage
 - Data back-up
 - Software updates/upgrades

- Informatics support
 - Handle user questions/requests
 - Gather user feedback
- Ongoing dedicated resource
 - Manage updates
 - Facilitate user support
 - Maintain infrastructure

It's all about control

YOU CONTROL	BENEFIT
Amount of storage	Handle large datasets
Type of hardware	Run compute intensive jobs
What tools to install	Customize to your research
Data access	Granular control of security
Networking architecture	No data transfer bottleneck
Software behavior	Optimize how jobs are run

Local Instance of Galaxy

Galaxy Wiki anushkanet Settings Logout | Search:

Admin/Get Galaxy Edit (C

Get Galaxy: Galaxy Download and Installation

In addition to using the [public Galaxy server](#) (a.k.a. [Main](#)), you can also install your own instance of Galaxy (what this page is about), or create a [cloud-based instance of Galaxy](#). Another option is to use one of the ever-increasing number of [Public Galaxy Servers](#) hosted by other organizations.

See [Big Picture: Choices](#) for help on deciding which of these options may be best for your situation.

Reasons to Install Your Own Galaxy

You only need to download Galaxy if you plan to:

1. [Develop](#) it further
2. [Add](#) new tools
3. [Plug-in](#) new datasources, or
4. [Run](#) a local production server for your site because you have
 1. Sensitive data (e.g., clinical)
 2. Large datasets or processing requirements that are too big to be processed on [Main](#)

Installation Procedure

The installation procedure is simple and is nearly identical for UNIX/Linux and Mac OS X. We are no longer supporting the Windows platform with

Contents

1. [Reasons to Install Your Own Galaxy](#)
2. [Installation Procedure](#)
 1. [Check your Python version](#)
 2. [Get the latest copy from the repository](#)
 3. [Start it up](#)
 4. [Join the Mailing List](#)
 5. [Keep your instance backed up](#)
 6. [Keep your code up to date](#)
3. [Advanced Configuration](#)
4. [Other Help](#)

- Customizing a local installation
 - Customize galaxy itself
 - Install 3rd party/commercial tools
 - Develop your own tools
 - Add shared genome builds
 - Integrate with instruments
 - Manage sensitive or proprietary data

Software	Comment
Python	Version 2.6 or 2.7
Galaxy	Most recent version
Web Server	To enable web hosting
Database	To store histories, users, track datasets, etc.
Analysis Tools	Need to follow install directions for each individual tools

<http://wiki.galaxyproject.org/Admin/GetGalaxy>

- Basic install of Galaxy

```
%
```

```
hg clone https://bitbucket.org/galaxy/galaxy-  
dist
```

```
% sh run.sh
```

- In your browser type: `http://localhost:8080`

- Install tools from toolshed
- Global config file: universe_wsgi.ini
 - Over 100 configurable parameters
 - Add admin users
 - Configure all galaxy settings
- Tool behavior
 - tool_conf.xml
 - Individual tool wrappers (also xml files)

- Galaxy for NGS requires additional tools
- <http://wiki.galaxyproject.org/Admin/Tools/Tool%20Dependencies>
- Set up reference genomes or fetch indexes

- Scalability
 - Handle more users (>5)
 - Run more jobs (>8 concurrent)
 - Support large datasets (>3TB)
- Efficiency
 - Schedule jobs
 - Optimize runs
 - Manage data

Hardware	Recommendation
# of Cores	32-64 processors
RAM	256-512 GB
Storage Amount	>10TB

Best Practices

page was renamed from Develop/Best Practices

Galaxy Software Development Best Practices

Lists software development best practices for the Galaxy Project. These are works in progress and practices vary in how broadly they have been applied.

Contents

1. [Metastandards](#)
2. [Python Standards](#)
3. [JavaScript Standards](#)
 1. [Backbone](#)
4. [Email Threads](#)

Develop

- [Mercurial](#)
- [Source Doc](#)
- [Best Practices](#)
- [Bitbucket](#)
- [Issues & Requests](#)
- [Data Model](#)
- [Search](#)

- Turn off developer settings
 - Speed up the galaxy server
- Switch databases (e.g. PostgreSQL)
 - handle multiple database requests at once
- Switch web servers (e.g. Apache,nginx)
 - Handle file transfers, external authentication

- Minimize redundant storage
 - Galaxy stores everything!
 - Filesystem compression
 - Deleting datasets and histories
- Improve Data Transfer
 - FTP, HTTP, SCP are slow
 - Consider data transfer products (e.g. Aspera, Globus)

- Dev, Test, Production environment
- Automated testing
- Leverage toolshed
- **DEDICATED RESOURCE!!!**
 - ½ FTE to support a large production installation

- Galaxy is a fantastic public resource
- Local Galaxy offers a high level of control over your analysis environment
- There is tremendous value in being aware of the underlying components of Galaxy
- The flexibility makes it complex to manage so **plan accordingly**

Thank You

Contact BioTeam to discuss your Galaxy needs

- Website: www.bioteam.net
- Email: anushka@bioteam.net



EXTRA SLIDES

SlipStream Galaxy Appliance

0010101101010011
SLIPSTREAM GALAXY EDITION
110110100100011
APPLIANCE
Galaxy made easy.



Powerful dedicated
desktop server
pre-configured with a fully
operational production
instance of Galaxy

Hardware Specifications

CPU	2x Intel® Xeon® Processor E5-2690, 8-core (16 cores total)
Memory	24x 16 GB RDIMM (384 GB)
Storage	7x 3TB SAS 6 Gbps HDD (16 TB usable) 1x 100GB SSD
Network	Dual Gigabit network adaptor
Power	Dual redundant power supplies

FEATURES	
Optimized Galaxy	Production configuration, optimized data transfer
Analysis Tools	Comprehensive set of open-source tools
Automated Updates	All software can be updated automatically
Preinstalled Datasets	5 model organisms (additional upon request)
Grid Compute	Grid Engine -based job management
Hardware Maintenance	Warranty includes maintenance
Open Platform	Install other software for your own use
Price: \$19,995 (USD)	

- **EARLY ACCESS PROGRAM (Limited Availability)**
 - Seamless Adoption
 - Dedicated Support
 - Workflow Generation
- **Early Development Partner Feedback**
 - “A device that centralizes functions with respect to data archives, storage, and analysis is a tremendous aid.” – Ed DeLong, MIT

- **Become an Early Access Partner Today!!**

- Web:

www.bioteam.net/slipstream/galaxy-edition