

**High Performance Research Computing**

*A Resource for Research and Discovery*



**TEXAS A&M**  
UNIVERSITY.

# HPRC Maroon Galaxy



# What is the Galaxy Project?

[usegalaxy.org](https://usegalaxy.org)

Galaxy is an open source, web-based platform for data intensive biomedical research. If you are new to Galaxy start here or consult our help resources. You can install your own Galaxy by following the tutorial and choose from thousands of tools from the Tool Shed.

**James Taylor (1979-2020)** believed that scientific progress can best be sustained through the mentoring of students and junior faculty.

To ensure implementation of this vision, the Galaxy community has established a foundation—Junior Training and Educational Connections Hotspot (JTech). JTech's mission is to (1) assist graduate students to participate in computational biology and data science conferences, and (2) organize and host mentoring sessions between senior and junior faculty members at high-profile meetings.

Design by Rebekka Paisner

- public Galaxy instance
- many popular bioinformatic tools are available
- no programming knowledge required
- accessed using a web browser
- reproducible workflow
- shared data and workflows
- try [usegalaxy.org](https://usegalaxy.org) to see if Galaxy is a good fit for your project

The [Galaxy Project](#) is supported in part by [NSF](#), [NHGRI](#), [The Huck Institutes of the Life Sciences](#), [The Institute for CyberScience at Penn State](#), and [Johns Hopkins University](#).

# Galaxy Tutorials

<https://galaxyproject.org/learn>



Use Learn Community Deploy & Develop Support Jobs @jtx

Search Galaxy



Edit

## Learn Galaxy

There are many approaches to learning how to use Galaxy. The most popular is probably to just dive in and use it. Galaxy is simple enough to use that you can do many analyses just by exploring the interface. However, you may miss much of the power this way.

Have you created or know of a resource that is useful for teaching with Galaxy? Then please share it! This will help others and also help get the word out about your resource. Use [this Google form](#) to describe your resource. **Also:** consider joining Galaxy Training Network and contributing your tutorial as described [here!](#)

[Tutorials by Galaxy Training Network](#)

[Tutorials using Galaxy Main](#)

[Interactive Tours](#)

[Tutorials from Lewis-Sigler Institute @ Princeton](#)

## Tutorials by Galaxy Training Network

Thanks to a large [group of wonderful contributors](#) there is a constantly growing [set of tutorials](#) maintained by the [Galaxy Training Network](#). These include:

### Introductory Tutorials

- [Introduction to Galaxy Analyses](#)
- [Data Manipulation](#)
- [User Interface and Features](#)

### Scientific Analyses

- [Assembly](#)
- [Computational chemistry](#)
- [Ecology](#)
- [Epigenetics](#)
- [Genome Annotation](#)
- [Imaging](#)

There are many tutorials available with example input data and step by step analysis for various topics



# Galaxy Shared Workflows

[usegalaxy.org](http://usegalaxy.org)

- can be used to share your workflow used in a publication
- input files are not stored here

## Published Workflows

search name, annotation, owner, and   
Advanced Search

Name	Annotation	Owner	Community Rating	Community Tags	Last Updated
Basic RNA-Seq Analysis - Differential Expression	From the RNA-Seq analysis tutorial during the Functional Genomics Workshop 2012 <a href="https://caps.osu.edu/pfg-workshop">https://caps.osu.edu/pfg-workshop</a>	mejia-guerra	★★★★★		Jun 22, 2012
RNA-seq differential analysis (single-end short reads, 2 conditions, 2 replicates)	Workflow based on Tophat and cuffdiff. Inputs: 4 fastq files (experiments), 1 bam file (pseudoreads), 1 gtf file (annotations). Outputs: bam, bigwig, xls,...	rna-seq-helin-group	★★★★★	illumina rnaseq cuffdiff tophat	Jul 17, 2013
RNA-seq preprocessing	Preprocessing of RNA-seq data by quality trimming of read ends and Illumina adapter removal. Inputs: fastq files. Outputs: fastq files and html files with...	rna-seq-helin-group	★★★★★	illumina rnaseq	Jun 03, 2013
RNA-seq differential expression analysis	RNA-seq differential analysis	chmy	★★★★★	easy	Aug 11, 2017

search results for rna-seq and sorted by Community Rating

# Galaxy Shared Workflow Example

**Galaxy** Workflow Visualize Shared Data Help User Using 0%

**Tools** search tools

**Inputs**

- Get Data
- Collection Operations
- Expression Tools
- GENERAL TEXT TOOLS**
- Text Manipulation
- Filter and Sort
- Join, Subtract and Group
- Datamash
- GENOMIC FILE MANIPULATION**
- FASTA/FASTQ
- FASTQ Quality Control
- SAM/BAM
- BED
- VCF/BCF
- Nanopore
- Convert Formats
- Lift-Over
- COMMON GENOMICS TOOLS**

**Preproc and TH**

- Fastq file (output (input))
- FASTQ Groomer**
  - File to groom
  - output\_file (fastqsanger, fastqcssanger, fastsolexa, fastqillumina, fastqsolexa.gz, fastqillumina.gz, fastqsanger.gz, fastqcssanger.gz, fastqsolexa.bz2, fastqillumina.bz2, fastqsanger.bz2, fastqcssanger.bz2)
- FASTQ Quality Trimmer**
  - FASTQ File
  - output\_file (input)
- Clip**
  - Input file in FASTA or FASTQ format
  - output (input)
- FastQC**
  - Short read data from your current history
  - Contaminant list
  - Adapter list
  - Submodule and Limit specifying file
  - FastQC on input dataset(s): Webpage (html)
  - FastQC on input dataset(s): RawData (txt)

**FastQC Read Quality reports** (Galaxy Version 0.72+galaxy1)

Label

Add a step label.

Step Annotation

Add an annotation or notes to this step. Annotations are available when a workflow is viewed.

**Short read data from your current history**

Data input 'input\_file' (fastq, fastq.gz, fastq.bz2, bam or sam)

**Contaminant list**

Data input 'contaminants' (tabular)

tab delimited file with 2 columns: name and sequence. For example: Illumina Small RNA RT Primer CAAGCAGAAGACGGCATA CGA

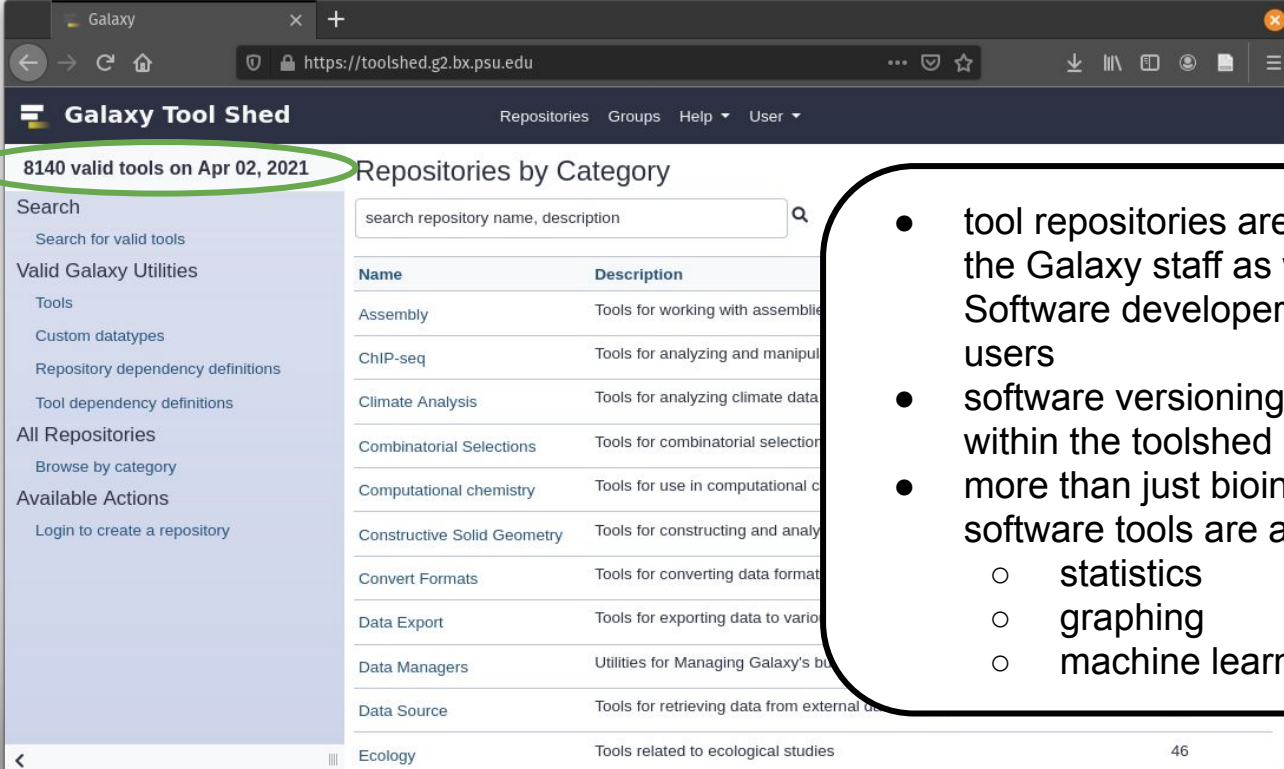
**Adapter list**

Data input 'adapters' (tabular)

list of adapters adapter sequences

# Galaxy Tool Shed

<https://toolshed.g2.bx.psu.edu>



8140 valid tools on Apr 02, 2021

Repositories by Category

Search repository name, description

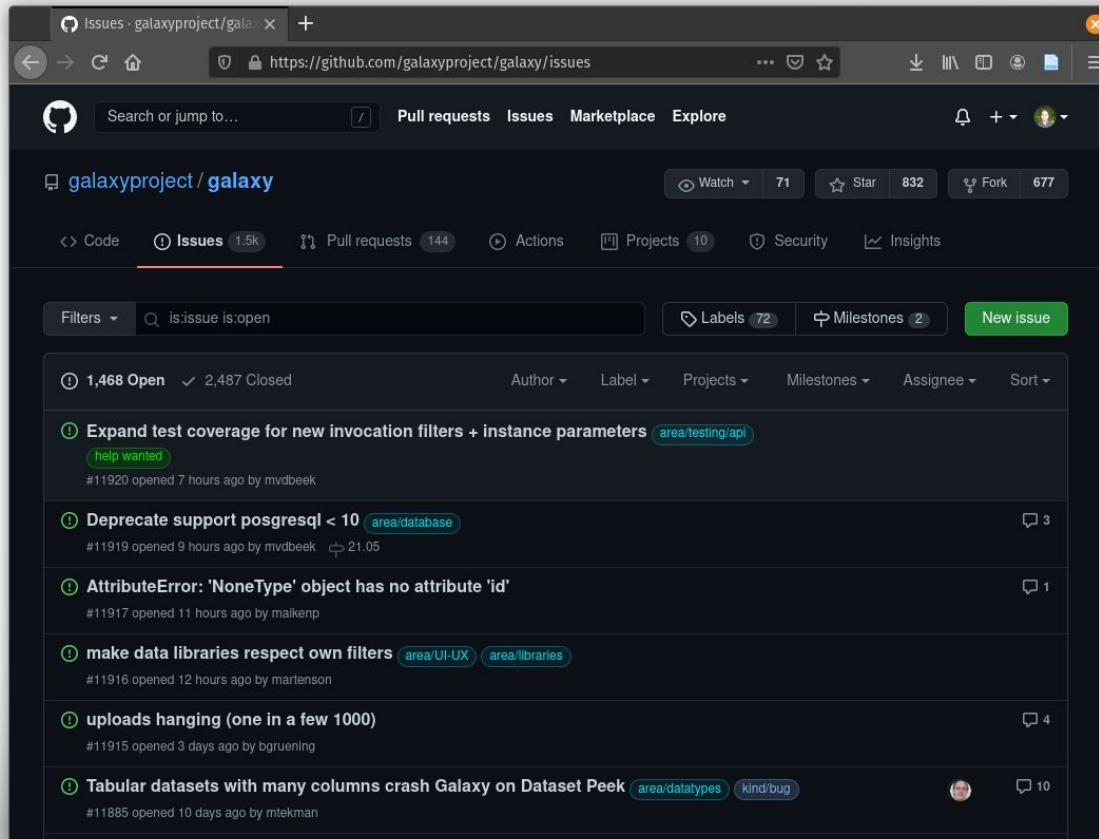
Name	Description
Assembly	Tools for working with assembl...
ChIP-seq	Tools for analyzing and manipu...
Climate Analysis	Tools for analyzing climate data...
Combinatorial Selections	Tools for combinatorial selectio...
Computational chemistry	Tools for use in computational c...
Constructive Solid Geometry	Tools for constructing and analy...
Convert Formats	Tools for converting data format...
Data Export	Tools for exporting data to vario...
Data Managers	Utilities for Managing Galaxy's bu...
Data Source	Tools for retrieving data from external...
Ecology	Tools related to ecological studies

46

- tool repositories are created by the Galaxy staff as well as Software developers and Galaxy users
- software versioning is available within the toolshed
- more than just bioinformatics software tools are available
  - statistics
  - graphing
  - machine learning



# The Galaxy Project has Community Support for Feature Requests and Resolving Issues



# HPRC Maroon Galaxy on the Ada Cluster





# HPRC Maroon Galaxy on the Ada Cluster

The screenshot displays the Galaxy web interface. At the top, there is a navigation bar with tabs for 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Admin', 'Help', and 'User'. The main content area features a dark red banner with the text 'Welcome to the HPRC Maroon Galaxy'. Below this, a yellow warning box states: 'Notice that Maroon Galaxy will be offline on Tuesday Nov 7 from 9:00 am to 5:00 pm for maintenance. All jobs running at that time will be stopped and you will need to restart your jobs after the maintenance is complete.' A blue information box follows, reading: 'Please contact the HPRC helpdesk to request a new tool or indexed genome, report errors or if you just have questions about using Galaxy.' A green box lists 'Newest tools added to Maroon Galaxy' with a list of tools including My HPRC SU balance, Fastq de-interlacer, Salmon 0.7.2, StringTie 1.3.3, PICAPD 2.8.3, deepTools 2.5.1, edgeR 3.14.0, Trinity 2.4.0, MACS2 2.1.1.20160309, GATK 3.6, Exonerate 2.4.0, BUSCO 3.0.2b, and InterProScan 5.25-64.0. The background of the main area is a vibrant image of the Maroon Galaxy. On the left, a sidebar contains a 'Tools' section with a search bar and a list of tool categories such as 'My HPRC SU balance', 'History divider', 'Get Data', 'Text Manipulation', 'Datamash', 'Statistics', 'Filter and Sort', 'Join, Subract and Group', 'Convert Formats', 'Extract Features', 'Fetch Alignments/Sequences', 'Operate on Genomic Intervals', 'Graph/Display Data', 'NCBI SRA Tools', 'Protein tools', 'Sequence Alignment', 'FASTA Tools', 'deepTools', 'BLAST+', 'EMBOSS', 'NGSEP', 'TAMU HPRC NGS TOOLBOX', 'NGS: QC and manipulation', 'NGS: Mapping', 'NGS: SAMtools', 'NGS: Picard Tools', 'NGS: BAMtools', 'NGS: BEDTools', 'NGS: Variant Analysis', 'NGS: de novo assembly', 'NGS: RNA-seq', 'NGS: ChIP-seq', 'NGS: ChV Tools', 'NGS: Population Analysis', and 'NGS: Metagenomics'. On the right, a 'History' panel shows a search bar and a list of datasets, including 'hisat2 picard' and several 'EstimatelibraryComplexity on data 1: Library complexity report' entries.

- TAMU HPRC Galaxy instance
- Accessed using a web browser
- Has been available to TAMU students, staff and faculty since 2015
  - will be retired on June 30, 2021
  - no more jobs submitted after May 31, 2021

# HPRC Maroon Galaxy on the Ada Cluster

- Installed in 2015 (v15.07) as a joint effort by Ping Luo (HPRC) and Dr. Rodolfo Aramayo (Department of Biology) - (currently 75 registered users)
- There are 350+ users across all HPRC Galaxies that are used for teaching and research
- Auto-installation of tools did not consistently work
- No support for compressed files (gzip) other than decompression on upload
- No auto-detection of fastqsanger format on upload
  - user had to run FastQC to evaluate quality score encoding or run Fastq Groomer
- Most newer software tools require v16.01+
- Early support for running tools on a cluster did not work well so tools were hard coded to use a specific amount of resource (cores, memory, time)
  - multiple tool configurations had to be added to support different resource requirements

## ----- 1 DAY JOBS -----

NCBI BLAST+ blastn 480 SUs.  
Search nucleotide database with  
nucleotide query sequence(s)  
(max runtime 1 day, 480 SUs  
required)

## ----- 3 DAY JOBS -----

NCBI BLAST+ blastn 1440 SUs.  
Search nucleotide database with  
nucleotide query sequence(s)  
(max runtime 3 days, 1440 SUs  
required)

## ----- 7 DAY JOBS -----

NCBI BLAST+ blastn 3360 SUs.  
Search nucleotide database with  
nucleotide query sequence(s)  
(max runtime 7 days, 3360 SUs  
required)



# HPRC Maroon Galaxy on the Grace Cluster



# HPRC Maroon Galaxy on the Grace Cluster

**Best Practices for Maroon Galaxy**

- **Contact** the HPRC helpdesk with an email to request a new or updated tool, indexed genome or to report an error. (Maroon Galaxy [docs](#), [slides](#))
- All users begin with a file quota of 1TB. Request an increase if you need more disk space but [permanently delete](#) nonessential files before requesting.
- FTP uploads are removed from ftp directory 48 hours after uploading so import your ftp files into Galaxy the same day as you upload to ftp
- The default job resource parameters for all tools is 1 core with 9GB memory for 24 hours (24 SUs).
  - Configuring a job to use 360GB memory for 1 hour requires 48 SUs. (360GB memory for 168 hours = 8,064 SUs).
- Only tools that support multi-core processing have the Job Resources Parameters option which allow you to select cores, memory and time.
  - Configuring a job to use 48 cores and 360GB memory for 1 hour requires 48 SUs. (48 cores for 168 hours = 8,064 SUs).
  - Configuring a job to use 80 cores and 2.93TB memory for 1 hour requires 80 SUs. (80 cores for 168 hours = 13,440 SUs).

COVID-19 related research on Galaxy: [training](#), [tutorials](#), [documents](#)

**Current known issues:**  
Some tools cannot be removed from favorites. Choose your favorites wisely.

**Uploading data from your computer**  
from Galaxy Project


**Local file upload**

<https://galaxy-grace.hprc.tamu.edu/maroon>

- TAMU HPRC Galaxy instance
- Available now to TAMU students, staff and faculty



# HPRC Grace Maroon Galaxy Features

- Installed in Spring 2021 (v21.01) on the HPRC Grace cluster.
- This is a significant deployment for the HPRC bio community for research and teaching.
- Improved integration of automatically installing tools from the Galaxy Toolshed.
- New support added to maintain file compression after uploading
  - Some software doesn't support compressed file format as input but file decompression is supported
- Multiple tool versions are available within a single tool 
- New auto-detection of fastqsanger format on upload
  - no need to use Fastq Groomer unless working with very early Illumina sequence data
- Integrated support for selecting cluster resources for each tool (cores, memory, time)

tools supporting **single-core** processing

Job Resource Parameters

Specify job resource parameters

Memory (GB)

9

Maximum Job Memory

Time (hours)

24

Maximum job time

tools supporting **multi-core** processing

Job Resource Parameters

Specify job resource parameters

Cores & Memory

48 cores & 360GB memory

1 core & 7GB memory

6 cores & 45GB memory

12 cores & 90GB memory

24 cores & 180GB memory

48 cores & 360GB memory

tools supporting **multi-core big memory** processing

Job Resource Parameters

Specify job resource parameters

Cores & Big Memory

48 cores & 360GB memory

20 cores & 730GB memory

40 cores & 1.44TB memory

60 cores & 2.19TB memory

80 cores & 2.93TB memory





# Galaxy Tutorials and Chat Features

HPRC Maroon Galaxy Analyze Data Workflow Visualize Shared Data Admin Help User Using 0%

Tools  
gffread

Galaxy Training! Contributors Help Extras Search

## Welcome to Galaxy Training!

Collection of tutorials developed and maintained by the worldwide Galaxy community

### Galaxy for Scientists

Topic	Tutorials
<a href="#">Introduction to Galaxy Analyses</a>	10
<a href="#">Assembly</a>	9
<a href="#">Climate</a>	4
<a href="#">Computational chemistry</a>	6
<a href="#">Ecology</a>	6

### Galaxy Tips & Tricks

Topic	Tutorials
<a href="#">Using Galaxy and Managing your Data</a>	17

### Galaxy for Developers and Admins

Topic	Tutorials
<a href="#">Galaxy Server administration</a>	

Local file upload Click to run unavailable.

OPEN CHAT

The same tutorials found at galaxyproject.org/learn can be accessed here and followed without leaving Maroon Galaxy

More than 300 Galaxy users and admins discussing anything Galaxy related





# Circos Tutorial Demonstration

The screenshot shows the HPRC Maroon Galaxy interface. The main content area is titled 'Galaxy Training!' and features a 'Visualisation' section. Below this, there are 'Requirements' and 'Material' sections. The 'Material' section contains a table with columns for Lesson, Slides, Hands-on, Input dataset, Workflows, Galaxy tour, and Galaxy instances. The row for 'Visualisation with Circos' has its 'Hands-on' icon circled in orange, and an orange arrow points to it from the left. The interface also includes a search bar, a chat window on the right, and a footer with navigation links.

Lesson	Slides	Hands-on	Input dataset	Workflows	Galaxy tour	Galaxy instances
Visualisations in Galaxy						
Genomic Data Visualisation with JBrowse						
Visualisation with Circos						

# Check Your HPRC SUs Balance

The screenshot shows the HPRC Maroon Galaxy interface. The left sidebar contains a list of tools, with 'My HPRC SU balance' highlighted by an orange arrow. The main panel displays the tool's configuration, including a search bar, a dropdown menu set to 'Show my current SU balance', and an 'Execute' button. Below the button, there is an information icon and a message: 'If you are unable to run this tool then you either need to renew your HPRC account or you are out of SUs.' A description follows: 'This tool retrieves a summary of your HPRC SU balance or allows the user to set the default account.' Below that, it says: 'Run this tool selecting the option 'Show my current SU balance' to get a list of your project account numbers.'

**You can also change your default HPRC project account and view SUs charged for each job run during the current fiscal year**

Account	Default	Allocation	Used & Pending SUs	Balance
10000000000001	Y	5000.00	-4990.00	10.00

# Shared Data Libraries

The screenshot shows the HPRC Maroon Galaxy interface. The top navigation bar includes 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data' (highlighted with an orange circle), 'Admin', 'Help', and 'User'. A dropdown menu for 'Shared Data' is open, showing options: 'Data Libraries', 'Histories', 'Workflows', 'Visualizations', and 'Pages'. Below the navigation, there is a search bar and buttons for '+ Folder', '+ Datasets', 'Export to His...', 'Delete', 'Details', and 'include deleted'. The breadcrumb path is 'Libraries / C\_dubliniensis\_CD36 data'. A table lists data libraries with columns: Name, Description, Date Updated (UTC), and State. One entry is visible: 'DR34\_R1.fastq.gz' with description 'uploaded fastqsanger.gz file', size '162.4 MB', and date '2020-11-09 08:28 PM'. A 'Manage' button is next to it. At the bottom, there is a pagination control showing '1' of 1 page, 15 items per page, and '1 total'.

Files can be added to a 'Data Library' which you can share with your Galaxy colleagues. Send a request to the HPRC helpdesk if you would like a Data Library for your group

# NCBI Blast Databases Available

NCBI BLAST+ blastn Search nucleotide database with nucleotide query sequence(s) (Galaxy Version 2.10.1+galaxy0) ☆ Favorite ▼ Options

**Nucleotide query sequence(s)**

3: M.tuberculosis\_genome.fasta

(-query)

**Subject database/sequences**

Locally installed BLAST database

**Nucleotide BLAST database**

Select/Unselect all

NCBI nt (partially non-redundant) 12 Apr 2021

**Type of BLAST**

megablast - Traditional megablast used to find very similar (e.g., intraspecies or closely related species) sequences

blastn - Traditional BLASTN requiring an exact match of 11, 1

blastn-short - BLASTN program optimized for sequences sho

dc-megablast - Discontiguous megablast used to find more di

Additional NCBI BLAST and custom databases can be added upon request

# HPRC Maroon Galaxy Access

- Try Galaxy at [usegalaxy.org](https://usegalaxy.org) to see if it appropriate for your project
- How to access HPRC Maroon Galaxy on the Grace Cluster
  - Available to Texas A&M students, staff and faculty with a NetID and an HPRC account
  - Apply for an HPRC account first
    - <https://hprc.tamu.edu/apply>
  - Then send an email request for a Maroon Galaxy account
    - [help@hprc.tamu.edu](mailto:help@hprc.tamu.edu)
  - Need to use TAMU [VPN](#) when connecting to Galaxy from off campus
  - Login to Maroon Galaxy using your TAMU NetID and password
- Read the Galaxy Usage Notes
  - <https://hprc.tamu.edu/wiki/SW:Galaxy>





**HIGH PERFORMANCE  
RESEARCH COMPUTING**  
TEXAS A&M UNIVERSITY

**Thank you.**

**Any questions?**

