Establishing a National Genomics Virtual Laboratory with Galaxy CloudMan

Enis Afgan, Clare Sloggett, Andrew Lonie, Michael Pheasant GCC 2012, Chicago

3 things

- New CloudMan features
- Blend library
- Creating a National Genomics Virtual Lab

CloudMan Features

- Start/launch through a web browser or the command line
- Choose between four cluster types
- Multiple cloud support
- Terminate/restart
- Scale
- Auto-scale
- Spot support
- Persist changes
- Share-an-instance (customized one too)
- Expand the file system
- Customize (via CLI or the Tool Shed): tools, data, references
- Access through ssh
- File system access for any S3 bucket -> data library
- Control the Galaxy process
- API (via Blend)
- Build your own image (via mi-deployment)

CloudMan Features

- Start/launch through a web browser or the command line
- Choose between four cluster types
- Multiple cloud support
- Terminate/restart
- Scale
- Auto-scale
- Spot support
- Persist changes
- Share-an-instance (customized one too)
- Expand the file system
- Customize (via CLI or the Tool Shed): tools, data, references
- Access through ssh
- File system access for any S3 bucket -> data library
- Control the Galaxy process
- API (via Blend)
- Build your own image (via mi-deployment)

Multiple Cloud Support

- Support for AWS, OpenNebula (thanks to Mattias de Hollander!), and OpenStack cloud middleware
- Run the same setup on a private or public cloud
 - Infrastructure setup is still automated

Spot Support

- Support for AWS Spot instances
 - Generally at ~30% cost of a full instance
 - Depend on demand and can terminate at any moment

Use S3 Buckets as File Systems

- Mount any S3 bucket as a local file system
- Gain fast access to AWS Public Datasets (e.g., 1000genomes)
- Also works on non-AWS clouds

APIs & Blend

CloudMan has an API

• Blend

A Python library for interacting with Galaxy's and CloudMan's APIs

\$ [sudo] pip install blend-lib

Requesting the infrastructure

Request an instance of CloudMan

from blend.cloudman.launch import CloudManLaunch
cml = CloudManLaunch('access_key', 'secret_key',
cloud)

cml.launch('cluster_name', 'img_id', 'inst_type', 'pwd')
cml.get_status()

Manipulating the infrastructure

Manage an instance of CloudMan from blend.cloudman import CloudMan cm = CloudMan('http://115.146.92.145', 'pwd') cm.initialize(type="Galaxy") cm.get_status() cm.add_nodes(3) cm.enable_autoscaling(0, 5) cm.get_galaxy_state()

Working with Galaxy

Import data into a Data Library and execute a workflow with it from blend.galaxy import GalaxyInstance

- gi = GalaxyInstance('http://usegalaxy.org', 'your API key')
- lib = gi.libraries.create_library('Set 13')
- d1 = gi.libraries.upload_file_from_url(lib['id'], 'http://
 tinyurl.com/gcc-exons')
- d2 = gi.libraries.upload_file_from_url(lib['id'], 'http://
 tinyurl.com/gcc-snps')
- gi.histories.create_history('Run 13')
- w = gc.workflows.get_workflows()[0]
- gi.workflows.show_workflow(w['id'])
- dataset_map = {'25': {'id':d1['id'], 'src':'ld'},

'27': {'id':d2['id'], 'src':'ls'}}

gi.workflows.run(lib['id'], dataset_map, 'F'+lib['id'])

Project Versions latest

RTD Search

```
Go
Full-text doc search.
```

Table Of Contents

Blend

- About
- Installation
- Usage
- Development
- API Documentation
- BioCloudCentral.org AP
 - CloudMan API
 - Galaxy API
- Testing Getting help
- Indices and tables

Next topic

API documentation for interacti

This Page

Show Source

Blend

Docs and examples included http://blend.readthedocs.org/

About

Blend is a Python (2.6 or higher) library for interacting with BioCloudCentral.org, CloudMan, and Galaxy's API. Conceptually, it makes it possible to script and automate the process of cloud infrastrucutre provisioning and scaling, as well as running of analyses within Galaxy. In reality, it makes it possible to do things like this:

Create a CloudMan compute cluster, via an API and directly from your local machine:

```
from blend.cloudman.launch import CloudManLaunch
cml = CloudManLaunch('<your cloud access key>', '<your cloud secret key')
cml.launch('Blend CloudMan', 'ami-<ID>', 'm1.small', 'password')
cml.get_status()
```

Manipulate your CloudMan instance and react to the current needs:

```
from blend.cloudman import CloudMan
cm = CloudMan("instance IP", "password")
cm.initialize(type="Galaxy")
cm.add_nodes(3)
cluster_status = cm.get_status()
cm.remove_nodes(2)
```

· Interact with Galaxy via a straighforward API:

```
from blend.galaxy import GalaxyInstance
gi = GalaxyInstance('<Galaxy IP>', key='your API key')
libs = gi.libraries.get_libraries()
gi.workflows.show_workflow('workflow ID')
gi.workflows.run_workflow('workflow ID', input_dataset_map)
```

Note

Although this library allows you to blend these three services into a cohesive unit, the library itself can be used with

CREATING A NATIONAL GENOMICS VIRTUAL LABORATORY

Univ. of Queensland and Victorian Life Sciences Compute Initiative (VLSCI)

Australian National Research Cloud

Provide computational infrastructure to support research

Compute and Storage (~25,000 cores + multiple PBs)



shell vs. IDE

Genomics Virtual Lab



Workflow platforms



Researcher activities

Cluster-on-the-cloud



Data catalogs



Researcher activities

| | | 🗐 CloudMan f | rom Galaxy | Admin Report bugs <u>Wiki</u> <u>Screencast</u> | | | | | Nutition readression Nutition readression Control of the second second Nutition Alignments Metagenomic analyses FASTA manipulation | | | Ch B: Cufflie transcript or Ch 7: Cufflie anne supers | inks on data 2 |
|---|--|--|--|--|--|----------------|-------------------------------------|---|---|--|-------------------------------------|--|--------------------------------|
| | | Cloud Man Welcome to <u>Cloud</u> is your first time | Initial Cluster Configuration | hin. If this ta store is | | | | | NS.BE (8451+) NGS: OC and manipulation NGS: Instal NGS: Assembly NGS: Mapping NGS: Indel Analysis | | | assembled) (% S: Cafflie transcript e (% 4: Cafflie | transcripts inks on data 1: |
| | | configured, defaul on which jobs are Termina | Welcome to CloudMan. This application will allow you to man the services provided within. To get started, choose the type to work with and provide the associated value, if any. | age this cluster and of cluster you'd like | | | | | NGS: RNA Analysis NGS: SAN Tools NGS: GATY Tools NGS: Peak Calling SNP/WGK-Color, UP Pops SNP/WGK-Color, UP Pops | | | 2:2-th 1:0-2h | asion a annotation |
| | 1 | Status Cluster name | Galaxy Cluster: Galaxy application, available tools, refer job manager, and a data volume. Specify the initial storage s G8 | rence datasets, SGE size (in Gigabytes): | | \sum | | 7 | | X NX - ubuntu@ip-1 | 0-101-30-41:2000 - CloudMan (GPL Ed | tion) | 10 Dubuntu |
| | | Disk status: Worker statu Service statu | Share-an-Instance Cluster: derive your duster form : duster. Specify the provided duster share-string (for example cm-0011923649e9271f17c4f83ba6846db0/shared/2011-0 | someone else's on? le, 8-1921-00): | | | | | Accessories | QTLCart | | | |
| C C Antochoud | ntral × \biocloudcentral.herokuapp.com/launch | Cluster stat | Ouster sh O Data Cluster: a persistent data volume and SGE. Specifisize (in Ggabytes): | CloudMan from (| Galaxy | Adm | n Report bugs Wiki Screencast | | Internet Other Programming | CLC Sequence Viewer | STANK. | | |
| oCloudCe | dCentral Coulter Doutlicitus and Gates plattere on Coul Computing in Coulter Counter States plattere on Coul Computing in Counter States Plattere States on Coulter States and Counter States on | | | CloudMan Console Welcome to <u>CloudMan</u> . This app within. Your previous data store manage services provided by the | CloudMan Console Welcome to <u>CouldMan</u> . This application allows you to manage this instance cloud cluster and the services provided within. Your previous data store has been reconnected. Once the cluster has initialized, use the controls below to | | | | | dendroscope | 3 | NO STATIS | townet. |
| ster name C | rent Protocols Demo | Name of your clu any name you ch Your choice of pa and accessing th | Hide extra options | Terminate cluster | Add nodes ¥ | Remove nodes v | Access Galaxy | | Ubuntu Software Center | X Exchanger S fastDNAml Forester ATV | 1AS | 11/00 | A AN |
| ud A | azon (AWS EC2) · | Choose from the provide below m | Start CloudMan Cluster | Status | | | | | | gap4 | C.C.A. | | |
| ess key 🛛 | AJKMSM6GLSW7V2CPA | Your Access Key the <u>security credentials page</u> . | | Cluster name: ghem | | | | | | a maxdLoad2 | | 10 A | |
| key E | Mvij9SLV.uxvb9OgeD58qiUXLkEZLa1 | Your Secret Access Key. For the Amazon cloud, a available from the security credentials page. | leo | Disk status: 0 / 0 (0%) | lable: 7 Beguest - 1 7 | | Autoscaling is off. Turn on? | | | X mesquite X Mr Bayes Multi | 35 8 | 19 1 J | 1 |
| ce type | e (4 ECUs / 7.5GB RAM) - | Type (je, virtual hardware configuration) of the inst start. | ance to | Service status: 1dle: 4 Avi Service status: Applications | inable: 2 Requested: 5 | | | | | 💥 Mr Bayes Multi 🎇 oligoarray | A. | Star Di | 10 |
| Jdvanced start. Lan instance | options | | | Cluster status log | | | 0 | | | K omegamap K pfaat | Crain Vent | | AL |
| website is an op it started doing a g the Amazon clos | service developed by the <u>CloudBioLinux</u> a lable biological analysis on cloud resource The <u>open source code</u> is available on GitHe | nd <u>CloudMan</u> communities. The goal is to make it is. See this guide for a detailed usage example ib allowing you to also run this service locally. | t easy when | | | | | | Welcome to Ubuntu 11.1 | 8 (GNU/Linux 3.0.0-14-v | irtual x86_64) | se galaxy | |

This site can be used for any of the available clouds. Note that you must have appropriate credentials for the chosen cloud. If a desired cloud is not available and you would like to see it there, please <u>contact us</u>.

Launching servers on the Amazon cloud will incur usage fees from Amazon for their resources. By using this service you acknowledge your sole responsibility for any costs accrued.

Graph this data and manage this system at https://landscape.canonical.com/

System information as of Wed May 2 04:53:21 UTC 2012

 System load:
 0.14
 Processes:
 134

 Usage of /:
 79.5% of 9.84GB
 Users logged in:
 0

 Memory usage:
 135
 IP address for eth0:
 115.146.94.163

 Snep usage:
 0%
 10%
 10%

111 packages can be updated. 56 updates are security updates.

💳 Galaxy (Options v)

Set Data Get Data Send Data DACODE Tools Uft-Deer Test Manipulation Titler and Set Jains. Subtract and Group Convert Foremas Dataset Fostures Fetch Sequences Fetch Asignments Get Genomic Scores Departure on Genomic Into Seatistics

statistics Graph/Display Data

Analyze Data Workflow Shared Data Visualization Help User

Welcome to Galaxy on the Cloud

1 1 1 1 1 1 1

Options v

13: Cuffcompare on # Ø 31 data 3: data 3: and data 6: data 6 tmap file

11: Caffcompare on \$\overline\$ \$\overl 10: Cuffcompare on
 4 0 30
 data 9. data 3. and data 6:
 transcript accuracy

© 2: Cafflinks on data 2: ● Ø 2t assembled transcripts

New release '12.04 LTS' available. Run 'do-release-upgrade' to upgrade to it.

* Documentation: https://help.ubuntu.com/

Get cloud support with Ubuntu Advantage Cloud Guest http://www.ubuntu.com/business/services/cloud ubuntu@server-3772:-\$ qstat -f queuename qtype resv/used/tat. 1 qtype resv/used/tot. load_avg arch states all.q@server-3772.novalocal BIP 0/0/1 0.33 1x24-amd64 all.q@server-3804.novalocal BIP 0/0/1 0.33 1x24-amd64 all.q@server-3805.novalocal BIP 0/0/1 0.60 1x24-amd64 all.q@server-3806.novalocal BIP 0/0/1 0.09 1x24-amd64 all.q@server-3807.novalocal BIP 0/0/1 0.65 1x24-amd64

all.q@server-3808.novalocal BIP 0/0/1 ubuntu@server-3772:~\$ [] 0.56 1x24-amd64

Tutorials and workshops



Researcher activities

Exemplar best practice workflows



Researcher activities

Exemplar workflows

• Variant calling:

- GATK best-practice
- o microbial
- cancer-optimised
- RNA-seq differential expression
- Fusion gene discovery from RNA-seq
- MicroRNA analysis
- De novo genome and transcriptome assembly
- Metagenomics
- ChIP-seq
- Variant annotation
- Pathway analysis
- Methylation



Scalable, reproducible, accessible

- New approaches for dealing with the data
 - Provide an illusion of *unlimited* data storage
 - New model for data accessibility
- Scale and replicate across multiple nodes around the country (and the world)
 - Integrate with public and other academic clouds
- Enable sharing and provide consistency
- Ensure 'upgradeability'
 - While permitting customization

Acknowledgments

Univ. of Queensland

- Xin-Yi Chua
- Michael Pheasant

VLSCI

- Enis Afgan
- Clare Sloggett
- Andrew Lonie

Garvan Institute

- Derrick Lin
- Warren Kaplan

Galaxy Team

