Developing distributed analysis pipelines with shared community resources using CloudBioLinux and CloudMan

Brad Chapman Bioinformatics Core Harvard School of Public Health

22 September 2011

Demonstration

Acknowledgements

CloudBioLinux – Ntino Krampis, Tim Booth, Dawn Field, Pjotr Prins and CloudBioLinux community

CloudMan – Enis Afgan, James Taylor Exome pipeline – HSPH, MGH, Win Hide, Oliver Hofmann

Demonstration

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Follow along

http://www.slideshare.net/chapmanb

Demonstration

Cue the "lots of data" slide

ls -lh fastq/

- 24G 1_110907_AD08A5ACXX_1_fastq.txt
- 21G 1_110907_AD08A5ACXX_2_fastq.txt
- 24G 2_110907_AD08A5ACXX_1_fastq.txt
- 20G 2_110907_AD08A5ACXX_2_fastq.txt

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Rapidly changing tools

Whole genome, deep coverage v1

Warning: the material on this page is considered out of date by the GSA team.

Best Practice Variant Detection with the GATK v2

Warning: the material on this page is considered out of date by the GSA team.

Best Practice Variant Detection with the GATK v3

Demonstration

・ロト・日本・モート モー うへぐ

Science – fundamental challenge

75% one-off experimental 25% reused code

Demonstration

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Unfortunate result

Why scientific programming does not compute (nature.com) 188 points by szany 67 days ago | comments

🛦 ajdecon 67 days ago | link

(Disclaimer: my background is in materials physics, and it may be different in other fields. But I doubt it.) Unfortunately there is very little *direct* incentive for research scientists to write or publish clean, readable code:

http://news.ycombinator.com/item?id=2735537

Demonstration

Hard choices

Computation

Demands flexible, well-architected, scalable code

Science

Requires rapid turn around and

experimentation

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ □臣 = のへ⊙

Demonstration

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

2 solutions (at least)

- I Improve your programming skills
- **2** Utilize community resources

Demonstration

Become a better coder



http://software-carpentry.org/

Demonstration

Community resources

Share painful parts

Base of well-written, scalable code

Start each problem from a higher level of abstraction

Demonstration

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Community components

- CloudBioLinux install software
- CloudMan manage cluster
- Exome analysis pipeline do science

Demonstration

CloudBioLinux

- Amazon image with bioinformatics software and libraries
- Automated build framework
- Community effort to maintain and extend

http://cloudbiolinux.org

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

CloudMan

- SGE cluster plus automation
- Web interface and monitoring
- Persistence and sharing
- Powers the Galaxy Cloud offering

http://wiki.g2.bx.psu.edu/Admin/Cloud

Demonstration

Exome analysis pipeline

Existing algorithms Aligners – Bowtie, BWA Variation – GATK Quality assessment – FastQC, Picard Messaging system – AMQP https://github.com/chapmanb/bcbb/

tree/master/nextgen

Fastq lane processing



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Sample processing



Variant calling



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Parallelization



Solution

Implementation



◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

Demonstration

Amazon

Virtual machines

- Share
- Reproduce
- Coordinate
- Accessibility



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Demonstration

What are we going to do?

- Use AWS console to boot CloudBioLinux
- Setup CloudMan in AWS console
- Boot CloudMan instance with demo data

Demonstration

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

What are we going to do? continued

- Manage cluster with CloudMan interface
- Setup messaging queue
- Run pipeline, examine results
- Share cluster

CloudBioLinux

Select and launch CloudBioLinux AMI from AWS console Connect

FreeNX graphical clientssh

Full tutorial PDF: http://j.mp/nnh5TE

Prep work

- Signup for AWS account: http://aws.amazon.com/
- Create login key pair in AWS ConsoleInstall NX client:

http://www.nomachine.com/select-package-client.php



https://console.aws.amazon.com/ec2/



Select CloudBioLinux image from Community AMIs

Request Inst	ances Wizard G	uncel 🗙 🚽
×	0	
CHOOSE AN AMI	INSTANCE DETAILS CREATE KEY PAIR CONFIGURE FIREWALL REVIEW	
Number of In	stances: 1	- 1
Availability Zo	ne: No Preference	fr
Advanced I	istance Options	
Here you can cl Monitoring or e	noose a specific kernel or RAM disk to use with your instances. You can also choose to enable CloudWatch Detailed nter data that will be available from your instances once they launch.	
Kernel ID:	Use Default V RAM Disk ID: Use Default V	1
Monitoring:	Enable CloudWatch detailed monitoring for this instance (additional charges will apply)	
User Data:	freenxpass: demo	
@as text		
⊎as file	base64 encoded	
Termination Protection:	Prevention against accidental termination.	
Shutdown	Stop Choose the behavior when the instance is shutdown from within the instance.	
Denavior.		
		- 1
		-
< Back	Continue	

enter NX password in user-data (freenxpass: secret)

Request filstances w	Zaru			Curren IA
<u> </u>	¥	~	O	
CHOOSE AN AMI INSTANCE	DETAILS CREATE KEY PAIR	CONFIGURE FIREWALL	REVIEW	
Please review the informat	ion below, then click Launc	n.		
AMI:	🗘 Ubuntu AMI ID ami-a	ad8e4ec4 (x86_64)	Edit AMI	
Number of Instances:	1			
Availability Zone:	No Preference			
Instance Type:	Micro (t1.micro)			
Instance Class:	On Demand		Edit Instance Details	
Monitoring:	Disabled	Termination Protection: Disabled	l.	
Tenancy:	Default			
Kernel ID:	Use Default Shutdow	n Behavior: Stop		
RAM Disk ID:	Use Default			
User Data:	freenxpass: demo		Edit Advanced Details	
Key Pair Name:	kunkel-keypair		Edit Key Pair	
Security Group(s):	sg-c1035aa8		Edit Firewall	
Back		Launch N		

Launch CloudBioLinux server

Navigation	My Instan	ces								
Region:	The Launch Instance Actions V							💭 Show/Hide 🛛 Refresh 🥝 Help		
US East (Virginia) 👻	Viewing: Al	I Instanc	es	•	All Instance	e Types 🔹	Search		≪ ≪ 1 to 1 of 1	Instances > >
EC2 Dashboard	Name	a 🦈 In	stance	AM	I ID	Root Device	Туре	Status	Security Groups	Key Pair Name
INSTANCES Instances Shot Requests	🗹 empt	v	1-24903744	ami	-ad8e4ec4	ebs	t1.micro	running	CloudBioLinux	kunkel-keypair
Reserved Instances										
 IMAGES AMIs Bundle Tasks 										
 ELASTIC BLOCK STORE Volumes Snapshots 										
 NETWORK & SECURITY Security Groups Elastic IPs 										
Placement Groups	Block	Devices		sda1						-
Load Balancers Key Pairs	Public	DNS:		ec2-1	184-73-64-4	7.compute-1.am	azonaws.com			
	Private	DNS:		ip-10	-244-166-1	73.ec2.internal				
	Private	IP Add	dress:	10.2	44.166.173					-
	Launc	n Time:		2011	-09-13 09:2	3 EDT				
	State 1	ransiti	on Reason:							
	Termin	ation P	Protection:	Disah	led					-

Get external hostname from Instances page

Session		Desktop	
NEMACHINE	Insert name of the session. Your configuration settings will be saved with this name. Session [cbi-demo Insert server's name and port where you want to connect. Host [cc2:184-73-64-47.compute-] Port [22] Select type of your internet connection. MODEM ISDN ADSL WAN LAN	NCMACHINE US US US SE US SE US SE SE SE SE SE SE SE SE SE SE SE SE SE	Ing NX Client you can run RDP, VNC and X sktops, depending on what the service provider has de available. nix GNOME Settings lect size of your remote desktop. D242768 H : 500 H : 500 H : 500 H D4 wr 200 H : 500 H : 500 H H : 500 H
	< Back Next > Cancel		< Back Next > Cancel
	Login ubuntu Password **** Session cbl-dem Configure	D Close	- - - -

Connect using NX client, with ubuntu user and secret password



```
kunkel:~ $ ssh -i ~/.ec2/id-kunkel.keupair ubuntu@ec2-184-73-64-47.compute-1.amazonaws.com
Warning: Permanentlu added 'ec2–184–73–64–47.compute–1.amazonaws.com' (RSA) to the list of kn
Welcome to Ubuntu 11.04 (GNU/Linux 2.6.38–8–virtual x86_64)
 * Documentation: https://help.ubuntu.com/
Sustem information disabled due to load higher than 1.0
At the moment, only the core of the system is installed. To tune the
system to your needs, you can choose to install one or more
predefined collections of software by running the following
command:
   sudo tasksel --section server
47 packages can be updated.
26 updates are security updates.
Last login: Tue Sep 13 13:40:34 2011 from 209-6-39-30.c3-0.smr-ubr1.sbo-smr.ma.cable.rcn.com
ubuntu@ip-10-244-166-173:~$ bowtie
No index, queru, or output file specified!
Usage:
  bowtie [options]* <ebwt> (-1 <m1> -2 <m2> | --12 <r> | <s>} [<hit>]
```

《曰》 《國》 《臣》 《臣》 臣

Connect with ssh, using private ssh key-pair

AWS Management Console > Amazon Elastic Compute Cloud (EC2) 4 Brad Chapman V Help V 📑 Show/Hide 🥔 Refresh 🥹 Help Region: Launch Instance Instance Actions -US East (Virginia) -Viewing: All Instances All Instance Types • Search |< < 1 to 1 of 1 Instances > >| EC2 Dashboard Name 10 Instance AMI ID **Root Device** Type Status Security Groups Key Pair Name - INSTANCES empty 1-24903744 ami-ad8e4ec4 t1.micro CloudBioLinux kunkel-kevpair ebs running Instances Instance Management Spot Requests Connect Reserved Instances Get System Log IMAGES Create Image (EBS AMI) AMIS Add/Edit Tags Change Security Groups Bundle Tasks FLASTIC BLOCK STORE Launch More Like This Volumes Disassociate IP Address Snapshots Change Termination Protection NETWORK & SECURITY View/Change User Data Security Groups Change Shutdown Behavior Elastic IPs Placement Groups sda1 Block Devices: Instance Lifecycle Load Balancers Public DNS: ec2-184-73-64-47.compu Terminate Key Pairs Reboot Private DNS: ip-10-244-166-173.ec2.ir Stop Private TP Address: 10 244 166 173 Launch Time: 2011-09-13 09:23 EDT **CloudWatch Monitoring** State Transition Reason: Enable Detailed Monitoring Termination Protection: Disabled

Terminate the server when finished

Demonstration

Setup CloudMan in AWS console

Create a custom security group

Full tutorial:

http://wiki.g2.bx.psu.edu/Admin/Cloud

Navigation	Security Groups	
Region:	🏷 Create Security Group 🕌 Delete	🎲 Show/Hide 💐 Refresh 🥝 Help
US East (Virginia) 🔻	Viewing: EC2 Security Groups V Search	🛛 🐇 🐇 1 to 2 of 2 Items 🔉 🔊
EC2 Dashboard	Name VPC ID Description	
INSTANCES Instances	Create Security Group	
Reserved Instances IMAGES AMIS Bundle Tasks ELASTIC BLOCK STORE Volumes Snapshots NETWORK & SECURITY Security Groups	Description: CloudBioLinux VPC: No VPC • Cancel Yes, Create	
Elastic IPs Placement Groups Load Balancers Key Pairs	Security group selected Security Group: default Details Inbound	
	Group Name: default	
	Group Description: default group	

Create security group rules following wiki instructions
Amazon Elastic Compute Cloud (EC2) Brad Chapman 🔻 | Help 🔻 Security Groups Region: 🏠 Create Security Group 🚺 Show/Hide 🏾 😂 Refresh 🛛 🥹 Help US East (Virginia) + Viewing: EC2 Security Groups * 1< 1 to 2 of 2 Items > EC2 Dashboard VPC ID Name Description - INSTANCES ~ CloudBioLinux CloudBioLinux Instances default default group Spot Requests Reserved Instances - IMAGES AMIC **Bundle Tasks 1** Security Group selected ELASTIC BLOCK STORE Security Group: CloudBioLinux Volumes Snapshots Details Inbound NETWORK & SECURITY Create a Custom TCP rule . Security Groups new rule: Elastic IPs sg-c1035aa8 Placement Groups Port range: 0 - 65535 Delete (CloudBioLinux) (e.g., 80 or 49152-65535) Load Balancers 20 - 210.0.0.0/0 Delete Key Pairs Source: 0.0.0/0 22 (SSH) 0.0.0.0/0 Delete (e.g., 192.168.2.0/24, sg-47ad482e, or 1234567890/default) 80 (HTTP) 0.0.0.0/0 Delete Add Rule 30000 - 30100 0.0.0.0/0 Delete 42284 0.0.0.0/0 Delete Apply Rule Changes

Final security group specifications

Boot CloudMan instance with demo data

- Start server
- Pass in CloudMan user data
- Load shared CloudMan image

HOOSE AN AMI	NSTANCE DETAILS CREATE KEY PAIR CON	NFIGURE FIREWALL REVIEW	
Choose an Amazon Quick Start M	Machine Image (AMI) from one of the tabbe	d lists below by clicking its Select button.	
Viewing: All Image	es CloudBioLinux	≪ ≪ 1	to 2 of 2 Items 🔉 🔌
AMIID	Root Device Manifest	Platform	
ami-90d32af9	ebs 678711657553/CloudBioLin	nux Ubuntu 10.04 LTS 64t 🥠 Ubuntu	Select 🔽
ami-ad8e4ec4	ebs 678711657553/CloudBioLin	nux Ubuntu 11.04 64bit 20 🥠 Ubuntu	Select

Follow same procedure as CloudBioLinux

Create CloudMan user-data file

```
cluster_name: cbldemo
password: cbl
access_key: your_access_key
secret_key: your_long_AWS_secret_key
```

◆□▶ ◆母▶ ◆臣▶ ◆臣▶ 臣 の�?

Request Inst	ances Wizard	Cancel 🗙 늘
¥	0	
CHOOSE AN AMI	INSTANCE DETAILS CREATE KEY PAIR CONFIGURE FIREWALL REVIEW	- 1
Number of In	stances: 1	1
Availability Zo	ne: us-east-1c	fre
Advanced I	istance Options	
Here you can c Monitoring or e	noose a specific kernel or RAM disk to use with your instances. You can also choose to enable CloudWatch Detailed ter data that will be available from your instances once they launch.	
Kernel ID:	Use Default RAM Disk ID: Use Default	15
Monitoring:	Enable CloudWatch detailed monitoring for this instance (additional charges will apply)	u
User Data:	Choose File cloudman-cbldemo.txt	
⊜ as text		- H
eas file	Base64 encoded	
Termination Protection:	Prevention against accidental termination.	- 1
Shutdown Behavior:	Stop Choose the behavior when the instance is shutdown from within the instance.	- 1
		- 1
		- 1
Beel	Continue	

▲□▶ ▲圖▶ ▲필▶ ▲필▶ - 필 -

Provide user-data from file

equest inst	tances Wizard				Cance
¥	¥	×0.	0		
HOOSE AN AMI	INSTANCE DETAILS	CREATE KEY PAIR	CONFIGURE FIREWALL	REVIEW	
Security group	s determine whether	a network port is or	en or blocked on your ins	tances. You may use an existi	ng security group, or
we can help yo	u create a new securi	ty group to allow ac	cess to your instances usi	ng the suggested ports below	Add additional ports
now or update	your security group t	inythine dailing the 5	ecurity oroups page.		
Choose o	ne or more of yo	our existing Sec	curity Groups		
sg-c1035aa8 -	CloudBioLinux				
sg-4752b62e -	default				
(Selected aro	ups: sg-c1035aa8)				
(beneeded gro	app. by crossado,				
© Create a	new Security Gro	oup			

Choose created security group

$\leftarrow \rightarrow \otimes$	() ec2-67-202-14-208.compute-1.	amazonaws.com/cloud
		The server ec2-67-202-14-208.compute-1.amazonaws.com:80 requires a username and password. The server says: CM Administration.
		User Name:
		Password: •••
		Cancel

Login to instance with password from user-data

CloudMan share-an-instance

Persist data in a CloudMan clusterEasily sharable

For this demo

cm-b53c6f1223f966914df347687f6fc818/shared/2011-10-07-14-00



Admin | Report bugs | Wiki | Screencast

값 🔛 🔧

Welcome to Galavi	Initial Cluster Configuration	rovided
within. If this is yo	Initial cluster configuration	ie data
store is configured		las
worker' nodes on	Welcome to Galaxy Cloudman. This application will allow you to manage this cluster and the services provided within. To get started, choose the type of cluster you'd like to work with and specify the size of your persistent data	
	storage, if any.	
Status		
	Start a full Galaxy Cluster. Specify initial storage size (in Gigabytes)	
Cluster name		
Disk status:	GB	ng is off
Disk Status.		on?
worker statu	Share-an-instance	
Service statu		
	cm-0011923649e9271f17c4f83ba6846db0/shared/2011-08 Shared instance	
	bucket path	
Cluster stat	8	•
	Data volume and SGE only. Specify initial storage size (in Gigabytes)	
	GB	
	CCE Only Me perdatent stars an available	
	© SGE Only. No persistent storage created.	
	Hide extra ontions	
	Start Churter	
	Start Cluster	

Import shared instance with demo data

Demonstration

Manage cluster with CloudMan

- Web-based console
- Monitor running processes
- Add nodes to cluster as needed

Galaxy Cloudman Console

Warning: You are running out of disk space. Use the disk icon below to increase your volume size.

Welcome to Galaxy Cloudman. This application will allow you to manage this cloud instance and the services provided within. If this your first time running this cluster, you will need to select an initial data volume size. Once the data store is configured, default services will start and you will be able to add and remove additional services as well as 'worker' nodes on which jobs are run.

	Terminate	cluster	Add nodes 🔻	Remove nodes	Access Galaxy
Sta	atus				
c	luster name:	cbidemo <	ĩ.		
)isk status:	17G / 20G (83%) 🚱		Autoscaling is off.
1	Vorker status:	Idle: 0 Ava	ilable: 0 Requested: 0		Turn on?
5	ervice status:	Applications	🗧 Data 😑		

	0
Cluster status log 17:36:40 - Master starting 17:36:42 - Completed initial cluster configuration. 17:37:02 - Prerequisites OK; starting service 'SGE' 17:37:08 - Configuring SGE. 17:37:16 - Successfully setup SGE; configuring SGE 17:37:16 - Successfully setup SGE; configuring SGE 17:37:16 - Successfully setup SGE; configuring SGE	
17:37:17 - Saved file 'cm_ boot.py' to bucket 'cm-56133hab111:04c3841ed40368b54e6' 17:37:17 - Saved file 'cm_targ' to bucket 'cm-56133hab111:04c3841ed40368b54e6' 17:37:17 - Saved file 'cmLarg' to bucket 'cm-56133hab111:04c3841ed40368b54e6' 17:45:09 - Retrieved file 'shared/2011-06:19-21-00;hared_Instance_file_list.txt'. 17:45:09 - Retrieved file 'sristend' data-aymi' from bucket 'cm- 0011923649e92711724f83ba6846db0' to 'shared_Instance_file_list.txt'.	

CloudMan console to interact with cluster

æ

💳 Galaxy Cloudman

Galaxy Cloudman Console

Warning: You are running out of disk space. Use the disk icon below to increase your volume size.

Welcome to Galaxy Cloudman. This application will allow you to manage this cloud instance and the services provided within. If this is your first time running this cluster, you will need to select an initial data volume size. Once the data store is configured, default services will start and you will be able to add and remove additional services as well as 'worker' nodes on which jobs are run.

Terminate c	luster	Add nodes 🔻	Remove nodes	Access Ga	laxy
tatus		Add	nodes		
		Number of n	odes to start:		
Cluster name:	cbidemo 🗧	1			
Disk status:	17G / 20G	C	Ж	Auto	scaling is off.
Worker status:	Idle: 0 Av	Type of	node(s):		furn <u>on</u> ?
Service status:	Application	(master node	type: m1.large)		
		Same as Master	•		
		Start Addit	tional Nodes		
					0
17:36:40 - Master	starting				
Cluster status 17:36:40 - Master 17:36:42 - Comple	log starting ted initial clust	Start Addit	tional Nodes		

Add node to cluster

Setup messaging communication

- Command line access to server
- Adjust RabbitMQ configuration
- Setup messaging queue

Command line access to server

ssh -i ~/.ec2/id-kunkel.keypair
 ubuntu@ec2-67-202-14-208.compute-1.amazonaws.com

Follow approach used to connect to CloudBioLinux cluster; can also connect via NX

Edit /export/data/galaxy/universe_wsgi.ini configuration file to add internal host name.

```
[galaxy_amqp]
host = ip-10-125-10-182.ec2.internal
port = 5672
userid = biouser
password = tester
```

Setup messaging queue

```
$ sudo rabbitmqctl add_user biouser tester
creating user 'biouser' ...
...done.
$ sudo rabbitmqctl add_vhost bionextgen
creating vhost 'bionextgen' ...
...done.
$ sudo rabbitmqctl set_permissions -p bionextgen
biouser ".*" ".*"
setting permissions for user 'biouser' in vhost 'bionextgen' ..
...done.
```

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Run pipeline, examine results

- Ready to run distributed pipeline
- Demo data two paired end fastq lanes
- Variant calling workflow

Input sequence data

\$ ls -1 /export/data/exome_example/fastq/ 7_100326_FC6107FAAXX_1-chr22.fastq 7_100326_FC6107FAAXX_2-chr22.fastq 8_100326_FC6107FAAXX_1-chr22.fastq 8_100326_FC6107FAAXX_2-chr22.fastq

Run level: YAML Configuration

\$ cat /export/data/exome_example/config/run_info.yaml fc_date: '100326' fc_name: FC6107FAAXX details: - files: [7_100326_FC6107FAAXX_1-chr22.fastq, 7_100326_FC6107FAAXX_2-chr22.fastq] lane: 7 description: Test replicate 1 analysis: SNP calling genome_build: hg19 algorithm: quality_format: Standard hybrid_bait: hybrid_selection/baits.bed hybrid_target: hybrid_selection/targets.bed

System level: YAML Configuration

```
$ cat /export/data/galaxy/post_process.yaml
program:
  bowtie: bowtie
  bwa: bwa
  ucsc_bigwig: wigToBigWig
  picard: /usr/share/java/picard
  gatk: /usr/share/java/gatk
  snpEff: /usr/share/java/snpeff
  fastqc: fastqc
distributed:
  cluster_platform: sge
  platform_args: '-q all.q'
  cores_per_host: 1
  rabbitmq_vhost: bionextgen
```

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Run exome pipeline

- \$ cd /export/data/work
- \$ distributed_nextgen_pipeline.py
 - /export/data/galaxy/post_process.yaml
 - /export/data/exome_example/fastq
 - /export/data/exome_example/config/run_info.yaml

What just happened?



Demonstration

Monitoring: SGE queues

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Monitoring: Analysis directory

\$ cd /export	t/data	a/work
\$ ls -lh		
drwxr-xr-x	4.0	alignments
-rw-rr	2.0K	automated_initial_analysis.py.o11
drwxr-xr-x	33	log
-rw-rr	15K	nextgen_analysis_server.py.o10
-rw-rr	15K	nextgen_analysis_server.py.o9
drwxr-xr-x	102	tmp

・ロト・日本・モート モー うへぐ

Monitoring: Log files

\$ less nextgen_analysis_server.py.o10 INFO: nextgen_pipeline: Processing sample: Test replicate 2; lane 8; reference genome hg19; researcher ; analysis method SNP calling INFO: nextgen_pipeline: Aligning lane 8_100326_FC6107FAAXX with bwa aligner INFO: nextgen_pipeline: Combining and preparing wig file [u'', u'Test replicate 2'] INFO: nextgen_pipeline: Recalibrating [u'', u'Test replicate 2'] with GATK

Retrieve results: Copy files

\$ upload_to_galaxy.py /export/data/galaxy/post_process.yaml /export/data/exome_example/fastq /export/data/work /export/data/exome_example/config/run_info.yaml

Final files copied into new directory; allows cleanup of analysis directory

Retrieve results: Output directory

\$ ls -lh /export/data/galaxy/storage/100326_FC6107FAAXX/7
-rw-r--r- 38M 7_100326_FC6107FAAXX.bam
-rw-r--r- 22M 7_100326_FC6107FAAXX-coverage.bigwig
-rw-r--r- 72M 7_100326_FC6107FAAXX-gatkrecal.bam
-rw-r--r- 109K 7_100326_FC6107FAAXX-snp-effects.tsv
-rw-r--r- 827K 7_100326_FC6107FAAXX-snp-filter.vcf
-rw-r--r- 1.6M 7_100326_FC6107FAAXX-summary.pdf

Share results

- Share-an-instance
- Uses CloudMan web interface
- Reproducible research
 - CloudBioLinux AMI software
 - CloudMan data and configuration

💳 Galaxy Cloudman

Galaxy Cloudman Console

Warning: You are running out of disk space. Use the disk ice

Welcome to Galaxy Cloudman. This application will allow you to manage within. If this is your first time running this cluster, you will need to sel store is configured, default services will start and you will be able to add 'worker' nodes on which jobs are run.



《曰》 《聞》 《臣》 《臣》 臣

Status

 Cluster name:
 cbldemo ≤ Share this cluster instance

 Disk status:
 17G / 20G (83%) (2)

 Worker status:
 Idle: 0 Available: 0 Requested: 0

 Service status:
 Applications ● Data ●

CloudMan console enables push button sharing

🚾 Galaxy Cloudman

Currently shared instances

Share-an-instance

This form allows you to share this cluster instance, at its current state, with others. You can make the instance public or share it with specific users by providing their account information below. You may also share the instance with yourself by specifying your own credentials, which will have the effect of saving the instance at its current state.

While setting up an instance to be shared, all currently running cluster services will be stopped. Then, a snapshot of your data volume and a folder in your cluster's bucket will be created (under 'shared/[current date and time]); this folder will contain your cluster's current configuration. The created snapshot and the folder will be given READ permissions to the users you choose (or make it public). This will enable those users to instantiate their own instances of the given cluster instance. This implies that you will only be paying for the created snapshot while users deriving a cluster from yours will incur costs for running the actual cluster. After the sharing process is complete, services on your cluster will automatically resume.

Public
 Shared

Share-an-instance

Can make public or available to specific collaborators

Admin | Report bugs | Wiki |

EC2 Cluster Configuration Are you sure you want to power the cluster off? This action will shut down all services on the cluster and terminate any worker nodes (instances) associated with this cluster. By default, the master instance will be left alive and should be terminated manually (using the AWS console). Automatically terminate the master instance? If checked, this master instance will automatically terminate after all services have been shut down. Also delete this cluster? If checked, this cluster will be deleted. This action is irreversible! All your data will be deleted. Yes, power off Worker status: Idle: 0 Available: 0 Requested: 0 Service status: Applications Data

When finished, turn everything off through CloudMan

🗖 Galaxy Cloudman



CloudBioLinux

- Shared machine image of biological software
- Boot from AWS console
- Connect with NX graphical client and ssh



CloudMan

- Cluster setup and management
- Boot from share-an-instance
- Manage cluster through web interface
- Share final results



Exome pipeline

- Parallel framework for running analyses
- Run using automated scripts
- Extract alignments, variant calls and summary information

Future: interfaces make it easier

Analysis type Datasets Parameters Summary	Alignment Exome variant calling	Align reads with BWA, recalibration, realign and variant call with GATK, Returns sorted BAM file, variant cals in VCF format, tab separated file of predicted effects and PDF with quality statistics.
		Next

https://bitbucket.org/hbc/galaxy-central-hbc

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Future: Simplified file selection

alysis type	Available	To process
tasets		Read
meters	1_110824_V0P8FM-gatkrecal.bam	
mary		7_100326_FC6107FAAXX_1-chr22.fastq
		Read pair (optional)
		7_100326_FC6107FAAXX_2-chr22.fastq
Demonstration

Future: Top level parameters

Analysis type	Include read	is that map to multiple loca	tions?	
Datasets	yes	-		
Param Parameter	S Target file (from history; bed format)		
Summary	targets.bed	-		
	Bait file (fro	om history; bed format)		
	baits.bed			
	Ownerstern			
	Urganism	2000 (CDCh27/he10) (h		
	Human reb.	2009 (GKCh57/hgra) (h		
	Barcodes		Samples	
	Illumina	454-Rapid		
	1 : ATCACG		Add new sub-sample	
	2 : CGATGT			
	3 : TTAGGC			
	4 : TGACCA			
	5 : ACAGTG			
	6 : GCCAAT			
	7 : CAGATC			
	8 : ACTTGA			
	9 : GATCAG			
	10 : TAGCTT			
	11 : GGCTAC			
	13 · CTTCTA			Ŀ
			Provinue	Next
			Frevious	Next

◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

Demonstration

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Future: Galaxy data libraries

Data Library "bchapman@hsph.harvard.edu"

Name	Message	Uploaded By	Date	File Size
□ ▼ ≧ 110824 V0P8FM ▼				
iii 🔻 📴 <u>1</u> 👻	1			
<u>1 110824 VOP8FM.bam</u> *		bchapman@hsph.harvard.edu	2011-08-25	37.8 Mb
1 110824 VOP8FM-coverage.bigwig *		bchapman@hsph.harvard.edu	2011-08-25	21.9 Mb
1 110824 VOP8FM-gatkrecal.bam *		bchapman@hsph.harvard.edu	2011-08-25	71.6 Mb
<u>1 110824 VOP8FM-snp-effects.tsv</u> =		bchapman@hsph.harvard.edu	2011-08-25	108.8 Kb
1 110824 VOP8FM-snp-filter.vcf *		bchapman@hsph.harvard.edu	2011-08-25	826.2 Kb
1 110824 VOP8FM-summary.pdf *		bchapman@hsph.harvard.edu	2011-08-25	1.6 Mb
For selected datasets: Import to current history v Go				

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Future: Galaxy analysis

💳 Galaxy	Analyze Data Workflow Shared Data	Visualization Admin Help	User	
Tools Options			<u>^</u>	History Options -
SCDE TOOLS				00
Gene List Comparison			-	Unnamed history 218.8 Mb
GALAXY TOOLS	: 1 (1_110824_V0P8FM-sort)			17:1 110824 VOP8FM- ● Ø X
Get Data	Reference organism	hg19		1.6 Mb
Send Data	Total	688,890 76bp paired		format: pdf_database: hg19
ENCODE Tools	Aligned Dairs aligned	679,724 (98,7%)		Info: uploaded pdf file
Lift-Over	Alignment combinations	339,194		
Text Manipulation	Pair duplicates	31,229 (9.2%)		
Filter and Sort	Insert size	151.7 +/- 30.5		Image in pdf format
Join Subtract and Group	Near bait bases	5,336,058 (11.8%)		
Some Some Some Some Source Sou	Off bait bases	18,249,851 (40.3%)		16.1 110034 VODOEN @ 0 %
Convert Formats	Mean bait coverage	231.8		16:1 110824 V0P8FM- @ / &
Extract Features	On target bases	244v (38.9%)		Varkiesalipalli
Fetch Sequences	10x coverage targets	93.1%		15-1 110924 VOD9EM- @ / M
Fetch Alignments	Zero coverage targets	4.3%		coverage bigwig
Get Genomic Scores 4	Fold enrichment	16056x		* 21.9 Mb
Operate on Genomic Intervals	In dbSNP	100.0%		format: bigwig, database: hg19
Statistics	Transition/Transversion (all)	-1.00		Info: uploaded bigwig file
Graph/Display Data	Transition/Transversion (dbSNP)	-1.00		
Multiple regression	Transition/Transversion (novel)	-1.00		display at UCSC main
Multivariate Analysis	Table 1: Summary of	f lane results		Discourse Discourse and the second
Multivariate Analysis				Binary ocsc Bigwig Tile
FASTA manipulation				

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Future: External UCSC visualization

move and a		zoom in 15x	3v 10v base	700m out 1 5x 3	3v 10v
				Zoom out 1.5x	
				the for a well of the	april 1
position/search ch	22:29.126.614-29.187.983	gene	iump cle	ear size 61.370 bp.	configur
1		and a second second			
chr22 (q12,1)	22n13 22n12 22n11.2	g11.21	a12.1 12.2 22a12.3	013.1 013.2 013.31	- 12k
Scale		28 KD			TIT
chr22:	29135000 29140000 2914500	0 29150000 29155000 2	9160000 29165000 2	9170000 29175000 29180	0000 29185
		1 110804 U008	EM coulonade biguid		
396 _		A_AAV067_VVI V	r H= COver age / D 190 19		
396 -		1_110024_0000	rn-coverage.bigwig		
396 -	ĺ a.	1_110024_0000	rn-coverage.bigwig		
396 -	Ì h	1_110024_0000	- cover age to 190 19		
396 -		1_11004_0010	rn-coverage.bigeig		
396 _ 110824 V0P8FM-c	lh.	1_110024_0000	ni-coverage.biguig		
396 _ 110824_V0P8FM-c		1_110041_0010	ni-coverage.bigwig		
396 _ 110824_V0P8FM-c		1_110041_0010	ne coverage to igo ig		
396 _ 110824_V0P8FM-c		1	n-coverage.bigwig		
396 _ 110824_V0P8FM-c			n-coverage, b 196 19		
396 _ 110824_V0P8FM-c		1_110044_0000	n-coverage.bigaig		
396 _ 110824_V0P8FM-co 1					
396 _ 110824_V0P8FM-c	UCSC Genes Bar	sed on RefSeq, Unifrot	, GenBank, CCDS and	Comparative Genomics	
396 _ .110824_V078FM-c. 1	UCSC Genes Bar	sed on RefSeq, Unifrot	, GenBank, CCDS and	Comparative Genomics	

Read more

- Step-by-step instructions http://j.mp/rp69nx
- Approaches to parallelism http://j.mp/nPQHcm
- Future work

http://bcbio.wordpress.com