# Cloud BioLinux: Pre-configured and On-demand Bioinformatics Computing for the Genomics Community

Ntinos Krampis
Asst. Professor
J. Craig Venter Institute
kkrampis@jcvi.org

http://www.jcvi.org/cms/about/bios/kkrampis/
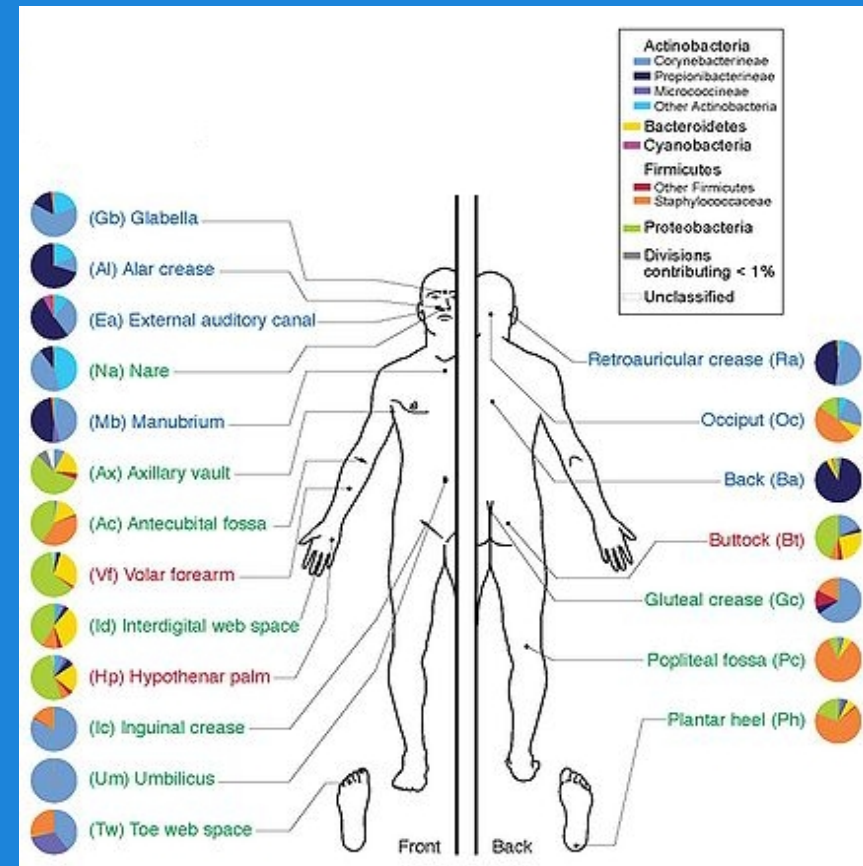
J. Craig Venter™

INSTITUTE

# J. Craig Venter Institute (JCVI)

## Large-scale genome sequencing and bioinformatics computing

- Human Microbiome Project (HMP): sequencing and assembly of 1000 reference microbe genomes from the human body

- Global Ocean Sampling (GOS) survey: metagenomic sequencing of microbes sampled from oceans around the world



2003 – 2008 Routes    2009 – 2010 Route



J. Craig Venter™
I N S T I T U T E

# JCVI: sequencing and computing infrastructure

## sequencing laboratory: 454, Solexa, HiSeq, and IonTorrent on the way

| Vendor: | Roche | | | Illumina | | | ABI | | |
|---|---|---|---|---|---|---|---|---|---|
| Technology: | 454 | | | Solexa GA | | | SOLiD | | |
| Platform: | GS20 | FLX | Ti | I | II | IIx | 1 | 2 | 3 |
| Images: (TB) | | 0.01 | 0.03 | 1 | 2.2 | 5.6 | 3.6 | 5 | 3.8 |
| PA Disk: (GB) | | 3 | 15 | 350 | 500 | 550 | 600 | 1500 | 2400 |
| PA CPU: (hr) | | 140 | 220 | 160 | 120 | NA | NA | NA | NA |
| SRA: (GB) | | 1 | 4 | 60 | 100 | 3.5 | 200 | 280 | 1200 |

J. Craig Venter™
I N S T I T U T E

# JCVI: sequencing and computing infrastructure

- large-scale sequencing needs large-scale informatics

- workhorse : ~1000 node Sun Grid Engine (SGE) cluster

- research in data processing and software development model with Hadoop / MapRecuce and a small private cloud

- bioinformatics department (57 bioinformaticians + software developers)



J. Craig Venter™
INSTITUTE

# A new paradigm:
## Low-cost, bench-top sequencers

- small-scale sequencers available: GS Junior by 454, MiSeq by Illumina

- complete sequencing of bacterial, viral, small fungal genomes

- RNAseq (gene expression), ChiPseq (protein interactions), gene variant discovery

- sequencing as a standard technique in basic genetics research - like PCR ?



http://www.gsjunior.com/                    http://www.illumina.com/systems/miseq.ilmn

J. Craig Venter™
INSTITUTE

# Sequencers shipped with minimal computational capacity

- <u>Problem 1</u>: sequence analysis requires  plenty of computational capacity

   For example: genome assembly, BLAST and genome annotation

- <u>Problem 2</u>: bioinformatics tools need expertise with unix/linux operating systems, software libraries, compiling source code etc.

   Difficult to install and use for biologists

**???**

J. Craig **Venter** ™
INSTITUTE

# Each lab with a sequencer building an informatics infrastructure ?

- difficult for individual PIs to get additional funds to build clusters

- funds for personnel to maintain the clusters and software

- duplication of effort across labs

- sub-optimal utilization of the hardware

- few sequencing runs per year

# Solution ?  Large sequencing centers offering bioinformatics services

- Bioinformatic Resource Centers (BRCs) by NIAID

- bioinformatic services usually coupled with sequencing of a genome

- provide data access and on-line tools

- cannot provide bioinformatic support for every PI in a lab acquiring a sequencing instrument

- need end-to-end solutions, users submit sequence data and get final annotation

# Solving Problem 1: sequence analysis requires computational capacity

- computational capacity on-demand without investment on hardware

- Amazon Elastic Compute Cloud (EC2), pay-by-the-hour computing

- cloud servers cost $0.085 - $2 per hour

- max capacity per server 64GB RAM / 8 CPU (but a PI can run thousands of servers)

- access to computing resources without institutional, economic or national boundaries

750 hours free for new users:
http://aws.amazon.com/free/

# Cloud Computing and Virtualization

- operating system, bioinformatics software and data, are pre-installed on a Virtual Machine (VM)

- a VM is a full-featured unix/linux server, in the form of a single, executable binary file

- the cloud provides the physical computational resources and virtualization layer to run the VM



Credit: VMware Inc.

# Solving Problem 2: bioinformatics tools need software engineering expertise

- a VM with pre-installed bioinformatics software publicly accessible on the cloud

- no need to compile source code, set-up configuration files, or other software dependencies

- PIs rent computational capacity to run the VM

- bioinformatics software can be accessed from anywhere in the world via a local computer with Internet access

- no need for sequencing informatics infrastructure at each laboratory

Amazon EC2 cloud

VM

VM

VM

Internet

local desktop computer at laboratory

J. Craig Venter ™
I N S T I T U T E

# Solving Problems 1 & 2: the Cloud BioLinux project

- Cloud BioLinux: a publicly accessible Virtual Machine (VM) on the Amazon EC2 cloud

- 100+ pre-installed bioinformatics tools on the VM with a graphical interface for non-technical users

- sequence analysis, genome assembly, annotation, phylogeny, molecular modeling, gene expression

- a researcher can initiate a practically unlimited number of Cloud BioLinux VMs for large-scale data analysis

**Krampis K.**, Booth T., Chapman B., Tiwari B., Bicak M., Field D. and Nelson K.E. *(2012) BMC Bioinformatics (in review)*, "Cloud BioLinux: pre-configured and on-demand bioinformatics computing for the genomics community"

**J. Craig Venter**™
INSTITUTE

# Cloud BioLinux for Bioinformatics

• how the Cloud BioLinux project came to be, what it can offers to small labs for genome sequence analysis

• where and how do I run Cloud BioLinux , especially if I am not a computer expert

• besides end-users, how bioinformatics developers are provided a framework for modifying and sharing VM configurations and data

http://www.cloudbiolinux.org

http://tinyurl.com/BioLinux-NEBC

J. Craig Venter™
I N S T I T U T E

# The making of Cloud Biolinux

tinyurl.com/BioLinux-NEBC

**+**

**=**

http://www.cloudbiolinux.org

- JCVI  bioinformatics cloud computing research
- NEBC BioLinux software repository
- community effort at BOSC 2009 – 11

- initially: a VM on Amazon EC2 with the tools copied and installed from the NEBC repository

- now: developer's framework for creating a customized cloud VM for bioinformatics

- main contributors:

HARVARD SCHOOL OF PUBLIC HEALTH

Galaxy

debian

J. Craig Venter ™
INSTITUTE

NEBC

# Accessing Cloud BioLinux



Account on the Amazon EC2 cloud   http://aws.amazon.com/ec2

J. Craig Venter
INSTITUTE

# Launch Cloud BioLinux through the EC2 cloud console



http://tinyurl.com/cloud-biolinux-tutorial

# Cloud BioLinux and VM launch wizard



Community AMIs, search for Cloud BioLinux VM identifier
(most recent update: cloudbiolinux.org)

select computational capacity for the VM

J. Craig Venter™
I N S T I T U T E

- remote desktop connection client

- free and open-source : http://nomachine.com

# Cloud BioLinux with remote desktop connection

Applications   Places   Syste

Sat May 7, 10:09 PM

- Accessories
- Bioinformatics
- Graphics
- Internet
- Programming
- Science
- Sound & Video
- System Tools
- Ubuntu Software Center

- CLC Sequence Viewer
- clustalw
- clustalx
- Cn3D
- dendroscope
- entrez
- Exchanger
- fastDNAml
- Forester ATV
- gap4
- jalview
- Jemboss
- maxdLoad2
- mesquite
- Mr Bayes Multi
- oligoarray
- omegamap
- pfaat
- pregap4
- sequin
- squint
- taverna
- TaxInspector
- tetra
- treeview
- trev

**Jemboss**

File   Preferences   Tools   Favourites   Help

ALIGNMENT
DISPLAY
EDIT
ENZYME KINETICS
FEATURE TABLES
INFORMATION
NUCLEIC
PHYLOGENY
PROTEIN
UTILS

GoTo:

abiview
aligncopy
aligncopypair
allversusall
antigenic
backtranambig
backtranseq
banana
biosed
btwisted
cai
cathparse
chaos
charge
checktrans
chips
cirdna
codcmp
codcopy
coderet
compseq
consambig

Keyword Search GO

● AND  ○ OR

Jemboss

🔧  (No Current Jobs)  ➤

I N S T I T U T E      RESEARCH COUNCIL

[ubuntu@domU-12-31...

# Cloud BioLinux:

## sharing data & results with VM snapshots

- access rights to the "snapshot" VM: public or for specific user

- other researchers access the VM with all the software, data, analysis results directly on the cloud

- storage cost: 0.10$ / GB / month

**Set AMI Permissions**

This image is currently Private

○ Public ⊙ Private

**Add Launch Permission:**
AWS Account Number 1: [_____]  add additional user

**Remove Launch Permission:**
No user permissions

Save

---

aws.amazon.com    AWS | Products | Developers | Community | Support | Account

| Amazon S3 | Amazon EC2 | Amazon VPC | Amazon Elastic MapReduce | Amazon CloudFront | Amazon RDS |

**Navigation**

Region: 🇺🇸 US East ▼

› EC2 Dashboard

INSTANCES
› Instances
› Spot Requests

IMAGES
› AMIs
› Bundle Tasks

**My Instances**

Launch Instance | Instance Actions ▼ | Reserved Instances ▼ | Show/Hide | Refresh | Help

Viewing: All Instances

**Instance Management**
- Connect
- Get System Log
- Get Windows Admin Password
- Create Image (EBS AMI)
- Add/Edit Tags
- Bundle Instance (S3 AMI)
- Launch More Like This
- Disassociate IP Address

| | | N₁ Instance | | Root Dev | Type | Status | | Lifecycle | Public D | Security | Key Pair | Moni |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | | i-4920: | ebs | | m1.large | 🟢 | running | normal | ec2-67-2( | default | jcvi_key: | disab |
| ☐ | er | i-f734( | ebs | | m1.large | 🟢 | running | normal | ec2-184-' | default | jcvi_key: | disab |
| ☐ | er | i-795b( | ebs | | m1.large | 🔴 | terminated | normal | | default | jcvi_key: | disab |
| ☐ | er | i-f933( | ebs | | m1.large | 🔴 | terminated | normal | | default | jcvi_key: | disab |

1 to 4 of 4 Instances

# Research at JCVI with Cloud BioLinux

- bioinformatics data analysis pipelines are complex

- approach: pre-install pipelines and all their dependencies within a VM

- make VM available on Amazon EC2

- we use private clouds, Eucalyptus & OpenStack

- open-source cloud platforms, fully compatible with Amazon EC2 (identical API)

- easy to set up on a local computer cluster, comes with Ubuntu Linux server edition

- also can run on your laptop with VirtualBox

instructions and VM at  http://www.cloudbiolinux.org

J. Craig Venter™
I N S T I T U T E

JCVI's Viral Genome
Sequencing Pipelines

Phase I-a
Sequencing & Assembly

Credit: Tim Stockwell, JCVI Viral Informatics

JCVI's Viral Genome
Sequencing Pipelines

Phase I-b
Sequencing & Assembly

Credit: Tim Stockwell, JCVI Viral Informatics

# JCVI's Viral Genome Sequencing Pipelines

## Phase II
## Annotation

- Assembled genomes as input to
  Viral Genome ORF Reader (VIGOR)

  Wang et al. BMC Bioinformatics 2010, 11:451

- detect coding regions, frame shifts,
  overlapping and embedded genes

- successfully used for annotating the
  influenza virus, rotavirus, rhinovirus,
  coronavirus and subtypes

## Phase III
## Annotation Visualization & Editing



J. Craig Venter
I N S T I T U T E

# Research at JCVI with Cloud BioLinux

- Funded by NIAID until 2013, focus on Viral,
  end-to-end, sequencing-to-annotation pipelines

- approach: pre-install pipelines and all their software
  dependencies in a VM

- export VM on Amazon EC2: pipelines ready to
  execute, no need to purchase hardware

- users simply need a web browser

- benefits small laboratories that lack resources or
  expertise

- if you own a cluster: download and run VM on your
  private Eucalyptus or Openstack cloud

**JCVI - GSC**

## National Institute of Allergy and Infectious Diseases
Leading research to understand, treat, and prevent infectious, immunologic, and allergic diseases.

**J. Craig Venter™**
INSTITUTE

# Scalable Data Analysis with Cloud BioLinux

- Sun / Oracle Grid Engine (GE) cluster: computational task scheduling

- Cloud BioLinux VM, dual role: Master or Worker

- Master VM coordinates distribution of computational tasks, Workers runs the computes

- The Master VM contains all code needed to start Workers and assemble a virtual cluster on the cloud

- Currently works on Amazon EC2

Master VM

Worker VM

Worker VM

Worker VM

Worker VM

# Scalable Data Analysis with Cloud BioLinux



- Galaxy Cloudman: users can control size of cluster, storage through a web-browser accessible interface

- Currently in the process of porting to Eucalyptus

- Users can download a VM which can bootstrap GE clusters on their private cloud

- Elastic capacity, size of virtual cluster

Afgan et al. BMC Bioinformatics 2010 11(Suppl 12):S4



J. Craig Venter™
I N S T I T U T E

# Cloud BioLinux for Software Developers

- Issue 1: for researchers with sensitive data a public cloud might not be an option
    moving VMs across clouds is not trivial, need low level operations

- Issue 2: bioinformatic specializations (ex. sequencing, phylogeny, protein structure)
    over-sized VM with too much software for all specializations

- Cloud BioLinux VM deployment framework

# Framework for Cloud Software Developers

- open-source framework to customize cloud Virtual Machines

- python Fabric automated deployment tool    ( DevOps )

- software installed in the VM listed in simple text configuration files

- Fabric scripts automatically pull and install software from repositories

- available from: https://github.com/chapmanb/cloudbiolinux

```
100644  39 lines (38 sloc)  0.668 kb

1   ---
2   # Top level configuration file that specifies w...
3   # should be installed. New sections that are ad...
4   # files should go here. Comment out any groups...
5   # installed.
6   packages:
7     - desktop
8     - programming
9     - distributed
10    - amazon
11    - python
12    - r
13    - ruby
14    - perl
15    - java
16    - erlang
17    - haskell
18    - databases
19    - math
20    - viz
21    - web
22    - bio_general
23    - bio_search
24    - bio_alignment
25    - bio_nextgen
26    - bio_sequencing
27    - bio_annotation
28    - bio_microarray
29    - bio_visualization
30    - bio_utils
31    - phylogeny
```

**software domains in Cloud BioLinux:**

Genome sequencing, *de novo* assembly, annotation, phylogeny, molecular structures, gene expression analysis

high-level configuration describing software groups for each group individual bioinformatics tools

```
516      - apache2
517   bio_general:
518      - emboss
519      - emboss-data
520      - emboss-lib
521      - primer3
522      - readseq
523      - bio-linux-taverna
524      - bio-linux-xcut
525   bio_search:
526      - blast2
527      - hmmer
528      - ncbi-tools-bin
529      - bio-linux-blast+
```
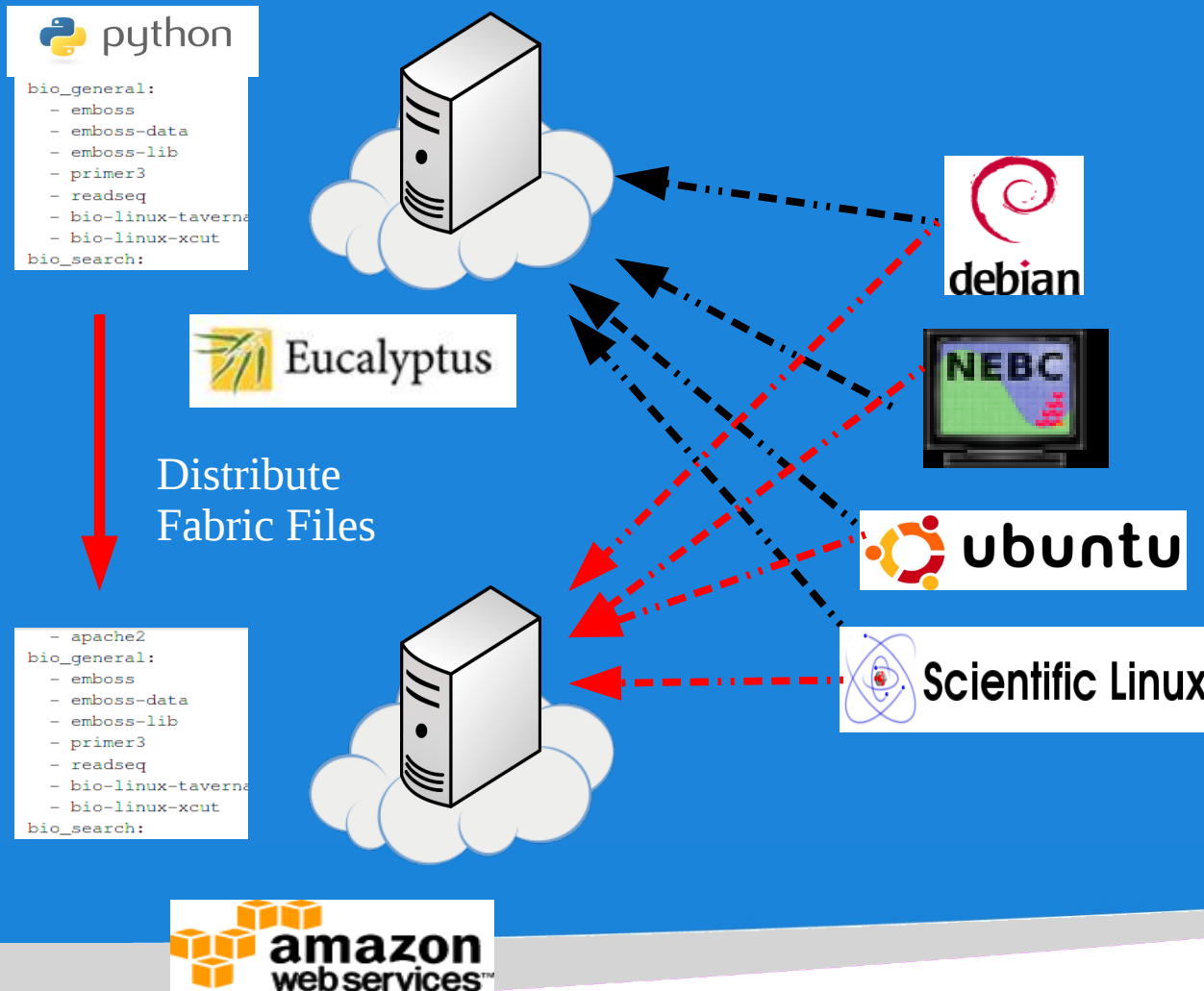
J. Craig Venter™
INSTITUTE

# Framework for Cloud Software Developers

**Customize Fabric Files**

**Custom VM**



```
bio_general:
  - emboss
  - emboss-data
  - emboss-lib
  - primer3
  - readseq
  - bio-linux-taverna
  - bio-linux-xcut
bio_search:
```

**Distribute Fabric Files**

```
  - apache2
bio_general:
  - emboss
  - emboss-data
  - emboss-lib
  - primer3
  - readseq
  - bio-linux-taverna
  - bio-linux-xcut
bio_search:
```

- start a fresh VM on Amazon or private cloud

- edit Fabric files to mix and match software from repositories – customized VM

- use source code repository to share configuration files

- share configuration of VM as source code

J. Craig Venter™
I N S T I T U T E

# Acknowledgments & Credits

*Brad Chapman*   -   development of the Fabric scripts, website

*Tim Booth, Mesude Bicak, Dawn Field*   –   BioLinux 6.0 development

*Enis Afgan*   –   Cloudman and Cloud BioLinux integration

*Members of the Cloud Biolinux community* - http://groups.google.com/group/cloudbiolinux

*Alex Richter* – porting Cloudman to Eucalyptus open-source cloud

JCVI  IT dept. - technology support


*Maria Giovanni, Punam Mathur* - NIAID / GSC funding

*Karen Nelson* – JCVI support for Cloud BioLinux

NIAID / OCICB – Bioinformatics Festival

## *Thank you !*

*kkrampis@jcvi.org*
*http://www.cloudbiolinux.org*
*http://www.slideshare.com/agbiotec*

**J. Craig Venter**™
I N S T I T U T E