

# Molecular Evolution, Genomic Analysis and FreeBSD

Joseph Mingrone

Department of Mathematics and Statistics  
Dalhousie University

# About the Bielawski Group



- ▶ Small group (5 - 10 members)
- ▶ Molecular Evolution and Genomic Analysis
  - ▶ We analyze DNA sequences to make inferences
  - ▶ The Research is purely computational
- ▶ Multidisciplinary

# Outline

- ▶ Hardware
- ▶ Research tracks
  - ▶ Modelling evolution at the molecular level
  - ▶ MicroBiome / Metagenomics
- ▶ Software, Design decisions, Observations

# Hardware

## Computing Cluster: Awarnach



- ▶ Purchased from Sun in 2006
- ▶ Sun Fire V40z master node (two dual core Opteron 870, 16 GB ECC RAM)
- ▶ Twenty compute nodes (X4100), two dual core Opteron 270, 4 GB ECC RAM, two 73 GB disks and four gigabit ethernet ports
- ▶ 48-port gigabit SMC switch



# Hardware

## New Compute Node

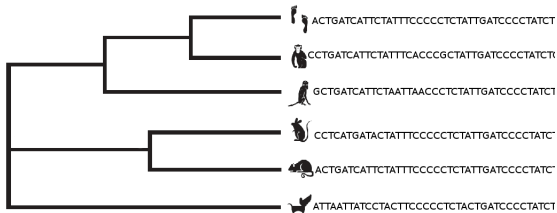
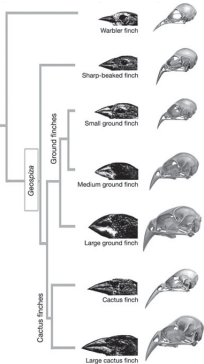


- ▶ New Compute Node
- ▶ Four 12-core 6348 CPUs
- ▶ 256 GB RAM

# Track 1: Background

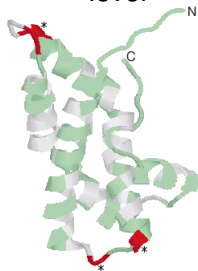


- ▶ Studying evolution is studying the past
  - ▶ Clues in present day to infer past processes
  - ▶ Morphology from fossil record
  - ▶ Studies with some organisms
  - ▶ Genetic material contains many markers of events in evolutionary history



# Track 1: Molecular Evolution Modelling

- ▶ Classify selection pressure at the protein and amino acid level



- ▶ Purifying Selection

- ▶ Neutral Evolution

- ▶ Positive Selection

- ▶ Usually most interested in detecting sites under positive selection

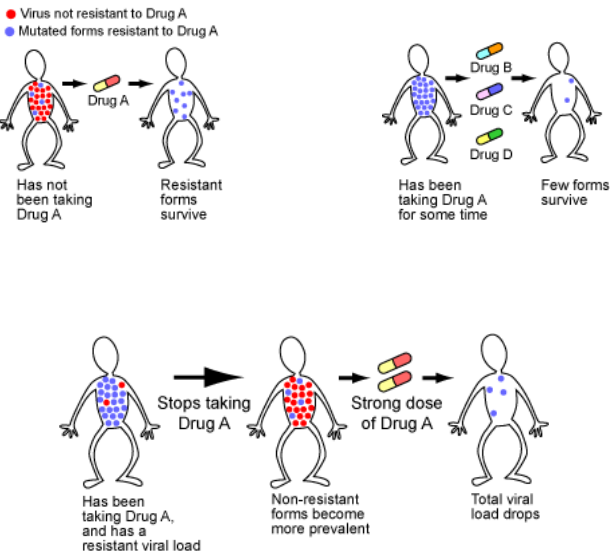


# Track 1: Arms Race

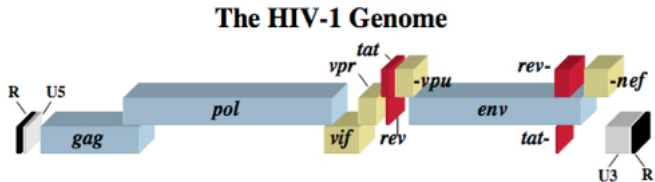
The law and order arms race...



# Track 1: Arms Race



# Track 1: Analysis of HIV Genes

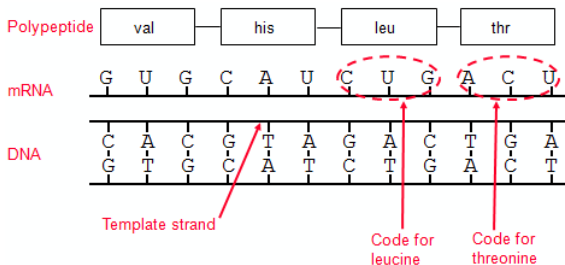


- ▶ Viral Envelop (Env) 3/91
- ▶ DNA polymerase (Pol) 11/947
- ▶ Viral infectivity factor (Vif) 10/192

# Track 1: Background

## The Central Dogma of Biology

- ▶ DNA is made up of four nucleotides with four different bases: Adenine (A), Cytosine (C), Guanine (G), or Thymine (T)
- ▶ Codons specify amino acids, the building blocks of protein



# Track 1: Redundant Code

	T	C	A	G
T	T T T phe	T C T ser	T A T tyr	T G T cys
	T T C phe	T C C ser	T A C tyr	T G C cys
	T T A leu	T C A ser	T A A stop	T G A stop
	T T G leu	T C G ser	T A G stop	T G G trp
C	C T T leu	C C T pro	C A T his	C G T arg
	C T C leu	C C C pro	C A C his	C G C arg
	C T A leu	C C A pro	C A A gln	C G A arg
	C T G leu	C C G pro	C A G gln	C G G arg
A	A T T ile	A C T thr	A A T asn	A G T ser
	A T C ile	A C C thr	A A C asn	A G C ser
	A T A ile	A C A thr	A A A lys	A G A arg
	A T G met	A C G thr	A A G lys	A G G arg
G	G T T val	G C T ala	G A T asp	G G T gly
	G T C val	G C C ala	G A C asp	G G C gly
	G T A val	G C A ala	G A A glu	G G A gly
	G T G val	G C G ala	G A G glu	G G G gly

- ▶  $4^3 = 64$  possible codons (61 sense), 20 amino acids
- ▶ Code is redundant
- ▶ Nucleotide substitution may or may not mean change in amino acid

# Track 1: A Measure of Selection Pressure

$d_S$  = rate of synonymous substitutions

$d_N$  = rate of nonsynonymous substitutions

$d_N/d_S = 1$       Neutral Evolution

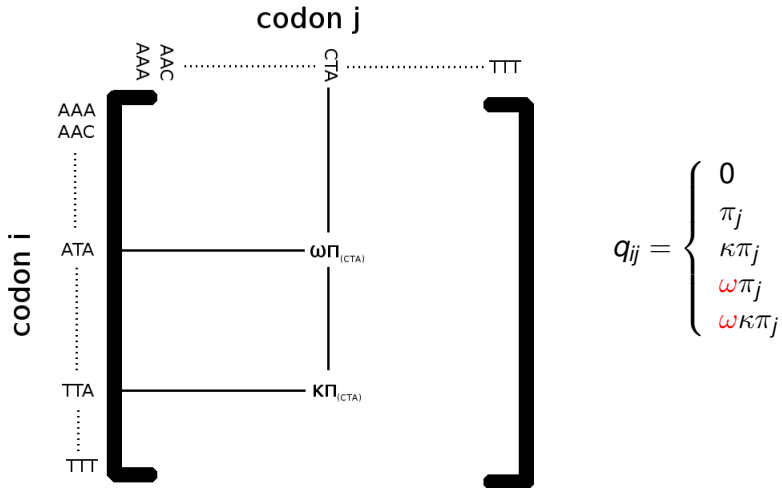
$d_N/d_S < 1$       Purifying Selection

$d_N/d_S > 1$       Positive Selection

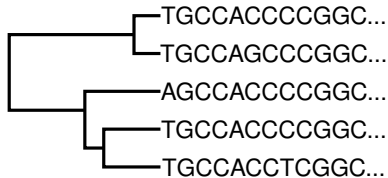
$$d_N/d_S = \omega$$

A measure of the strength and direction of selection pressure

# Track 1: Markov Process



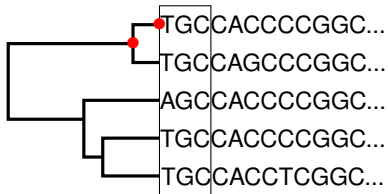
# Track 1: Estimating Model Parameters is Computationally Intense



$$q_{ij} = \begin{cases} 0 \\ \pi_j \\ \kappa\pi_j \\ \omega\pi_j \\ \omega\kappa\pi_j \end{cases}$$

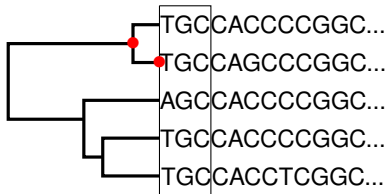


# Track 1: Estimating Model Parameters is Computationally Intense



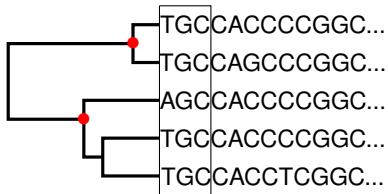
$$q_{ij} = \begin{cases} 0 \\ \pi_j \\ \kappa\pi_j \\ \omega\pi_j \\ \omega\kappa\pi_j \end{cases}$$

# Track 1: Estimating Model Parameters is Computationally Intense



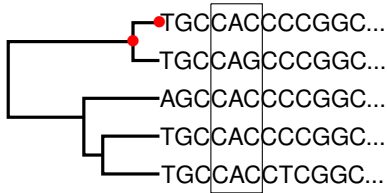
$$q_{ij} = \begin{cases} 0 \\ \pi_j \\ \kappa\pi_j \\ \omega\pi_j \\ \omega\kappa\pi_j \end{cases}$$

# Track 1: Estimating Model Parameters is Computationally Intense



$$q_{ij} = \begin{cases} 0 \\ \pi_j \\ \kappa\pi_j \\ \omega\pi_j \\ \omega\kappa\pi_j \end{cases}$$

# Track 1: Estimating Model Parameters is Computationally Intense



$$q_{ij} = \begin{cases} 0 \\ \pi_j \\ \kappa\pi_j \\ \omega\pi_j \\ \omega\kappa\pi_j \end{cases}$$

## Track 2: MicroBiome / Metagenomics

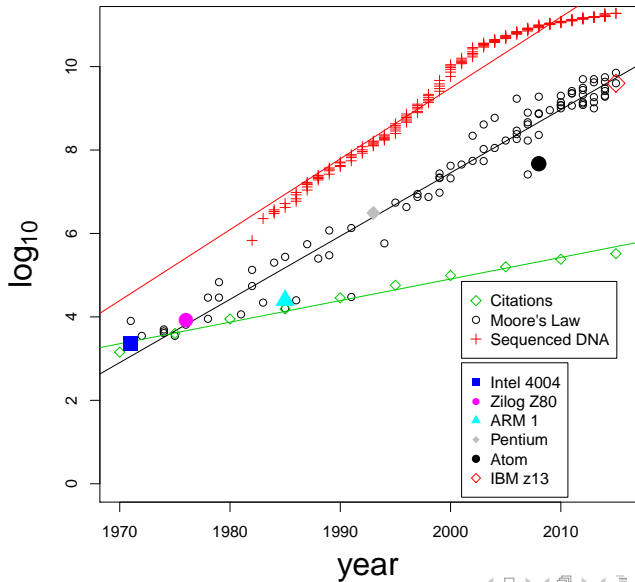
- ▶ Culturing diverse microbial communities in the lab can be difficult
- ▶ Sequence DNA as it exists in the natural environment (metagenomics)
- ▶ Understand associations between changes in microbiomes and changes in complex systems
  - ▶ BiomeNET BioMiCo: group microbial communities according to properties

El-Swaiss, Heba, et al. "Seasonal assemblages and short-lived blooms in coastal north-west Atlantic Ocean bacterioplankton." *Environmental microbiology* (2014).

Shafiei, Mahdi, et al. "BioMiCo: a supervised Bayesian model for inference of microbial community structure." *Microbiome* 3.1 (2015): 1-15.

Shafiei, Mahdi, et al. "BiomeNet: A Bayesian model for inference of metabolic divergence among microbial communities." (2014): e1003918.

# Research and Computing Technology are Coupled



# Software

- ▶ Originally Solaris, since 2009 FreeBSD release (7.x - FreeBSD 10.1), STABLE on storage / new compute node

# Software

- ▶ Originally Solaris, since 2009 FreeBSD release (7.x - FreeBSD 10.1), STABLE on storage / new compute node
- ▶ Poudriere, VirtualBox (MatLab)



# Software

- ▶ Originally Solaris, since 2009 FreeBSD release (7.x - FreeBSD 10.1), STABLE on storage / new compute node
- ▶ Poudriere, VirtualBox (MatLab)
- ▶ ZFS (master node and storage server), NFS

# Software

- ▶ Originally Solaris, since 2009 FreeBSD release (7.x - FreeBSD 10.1), STABLE on storage / new compute node
- ▶ Poudriere, VirtualBox (MatLab)
- ▶ ZFS (master node and storage server), NFS
- ▶ math/R, lang/perl5.20

# Software

- ▶ Originally Solaris, since 2009 FreeBSD release (7.x - FreeBSD 10.1), STABLE on storage / new compute node
- ▶ Poudriere, VirtualBox (MatLab)
- ▶ ZFS (master node and storage server), NFS
- ▶ math/R, lang/perl5.20
- ▶ print/texlive-full

# Software

- ▶ Originally Solaris, since 2009 FreeBSD release (7.x - FreeBSD 10.1), STABLE on storage / new compute node
- ▶ Poudriere, VirtualBox (MatLab)
- ▶ ZFS (master node and storage server), NFS
- ▶ math/R, lang/perl5.20
- ▶ print/texlive-full
- ▶ biology/paml / Proteus

# Software

- ▶ Originally Solaris, since 2009 FreeBSD release (7.x - FreeBSD 10.1), STABLE on storage / new compute node
- ▶ Poudriere, VirtualBox (MatLab)
- ▶ ZFS (master node and storage server), NFS
- ▶ math/R, lang/perl5.20
- ▶ print/texlive-full
- ▶ biology/paml / Proteus
- ▶ biology/ncbi-blast+, Diamond, Humann, Qiime, Metaphlan, Biomenet, BioMico, BMTagger, PEAR

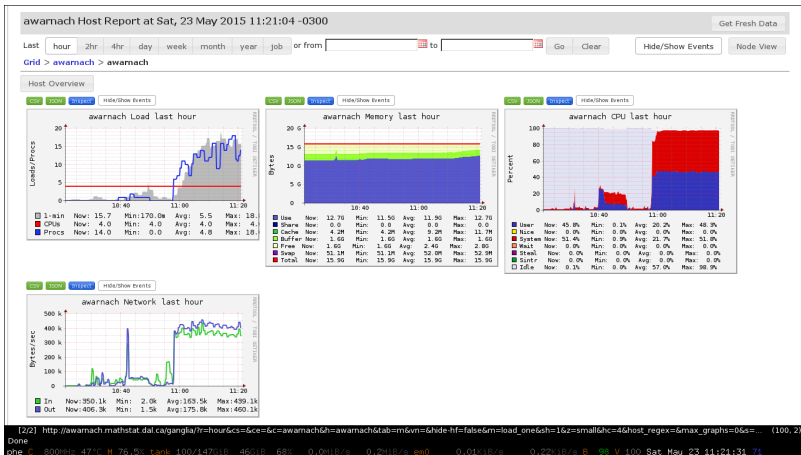
# Software

security/tmux-cssh

```
gw1 gw2 gw3 gw4
# # # #
gw5 gw6 gw7 gw8
# # # #
gw9 gw10 gw11 gw12
# # # #
gw13 gw14 gw15 gw16
# # # #
gw18 gw19 gw20
# |
[tmux-cssh1:ssh*
[0] 1:gly 2:phe- 3:awarnacht 4:trp
phe © 2200Hz 49°C W 45.8h tank 108/1470.0 38.0B 74% 0.001B/s 0.891B/s ewd 1.471B/s 1.201B/s B 90 V 100 Tue Jun 09 22:32:11
```

# Software

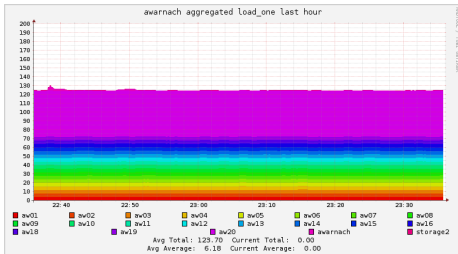
sysutils/ganglia-monitor-core sysutils/ganglia-webfrontend



# Software

sysutils/ganglia-monitor-core sysutils/ganglia-webfrontend

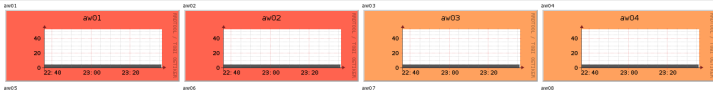
Stacked Graph - load\_one



awarnach **load\_one** last hour sorted by name

Metric:  Show Hosts Scaled:  Same  Size:  Columns:  (0 = metric + reports)

Show only nodes matching  Filter: Max graphs to show  Sorted:



[15/15] [http://awarnach.mathstat.dal.ca/ganglia/?r=hour&cs=6ce=6m=load\\_one&s=by+name&c=awarnach&tab=m&vn=6hide-N=false](http://awarnach.mathstat.dal.ca/ganglia/?r=hour&cs=6ce=6m=load_one&s=by+name&c=awarnach&tab=m&vn=6hide-N=false) (100, 48)  
Done  
php 0.2500/s 49.1% M 45.8% tank 108/147C 16 38.1B 74% 0.001B/s 0.001B/s ew0 0.12/1B/s 0.19/1B/s B 96 V 100 Tue Jun 09 22:36:26



# Software

## Basic Unix Tools

```
#!/bin/sh
```

```
fping=/usr/local/sbin/fping
```

```
nodes="aw1 aw2 aw3 aw4 aw5 aw6 aw7 aw8 aw9 aw10 aw11 aw12 aw13 aw14 aw15 aw16 aw18 aw19 aw20"
```

```
rdo_user=root
```

```
for host in `${fping} -a ${nodes}`; do
```

```
    printf "*** ${host} ***\n"
```

```
    ssh ${rdo_user}@${host} $*
```

```
done
```

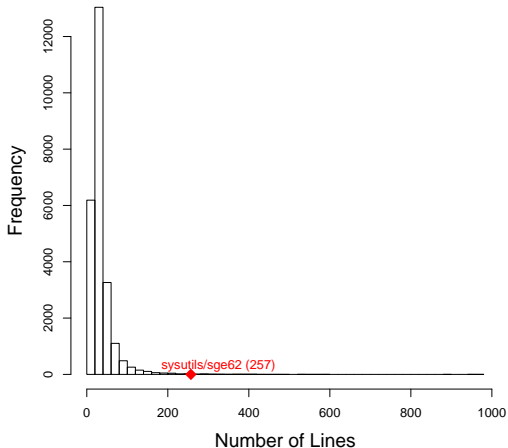
# Software

Batch Submission, Resource Management with Sun Grid Engine

# Software

Batch Submission, Resource Management with Sun Grid Engine

Distribution of Port Makefile Size



# Software

sysutils/slurm-hpc (slurm-wlm)

- ▶ Despite the name, Simple Linux Utility for Resource Management (Slurm): “Portability: Written in C with a GNU autoconf configuration engine. While initially written for Linux, Slurm has been ported to a diverse assortment of systems.”
- ▶ “Sequoia, an IBM BlueGene/Q system at Lawrence Livermore National Laboratory with 1.6 petabytes of memory, 96 racks, 98,304 compute nodes, and 1.6 million cores, with a peak performance of over 17.17 Petaflops.”

# Software

Porting Software Written by Biologists





# Thank You

# Questions?

Image credits:

Arms race: [www.inkcintc.com.au](http://www.inkcintc.com.au)

HIV images: [http://evolution.berkeley.edu/evolibrary/article/medicine\\_04](http://evolution.berkeley.edu/evolibrary/article/medicine_04)