# Bioinformatics Crash Course

Ian Misner Ph.D.

Bioinformatics Coordinator

UMD Bioinformatics Core

# The Plan

- Monday
  - Introductions
  - Linux and Python Hands-on Training
- Tuesday
  - NGS Introduction
  - RNAseq with Sailfish (Dr. Steve Mount, CBCB)
  - RNAseq with Tuxedo package
- Wednesday
  - Genome Sequencing Introduction
  - Genome Assembly and QC
  - Metagenomics (Dr. Mihai Pop, CBCB)

- Thursday
  - Genome Annotation
  - PacBio Genome Assembly (Matt Conte)
  - Review Genome Assembly and Annotation
- Friday
  - Cloud computing and Galaxy
    - Variant Detection and RNAseq analysis
- Each day we can have a Q&A session to find out what works or doesn't work as well as try to address any topics we haven't covered.

- http://www.biology.umd.edu/files/biology/bioinformatics/Workshop_July_2014.pdf

# UMD Bioinformatics Core



- Mission: To provide users with the bioinformatic services, support, and education necessary to advance their research program.

- The Core has partnered with the Division of IT to provide the necessary computational resources needed for these demanding analyses.

UNIVERSITY OF
MARYLAND

Bioinformatics Core

# Bioinformatics Core Services

- Raw data processing
- Genome and transcriptome assembly
- RNAseq analysis
- Variant discovery

- Grant writing support
- Experimental design assistance
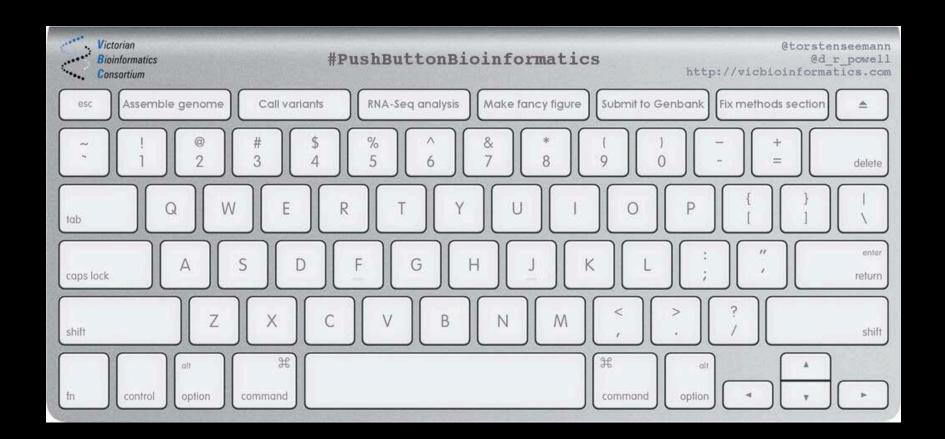- Workflow and pipeline construction
- Custom analyses

# The Goal

- Allow you to conduct your own analysis.

- Get you comfortable using the command line.

- Introduce programing, HPCC's, and DeepThought2.

- Learn some of the best practices with experimental design and data analysis.

- Avoid common pitfalls with data processing.

# What is Bioinformatics?

- Interdisciplinary field combining:
  - Computer science
  - Statistics
  - Mathematics
  - And Biology
- Lots of different areas of expertise:
  - Biological programing
  - Software development
  - Hardware development
  - Experimental design
- Difficult for an individual to be an expert in all areas.
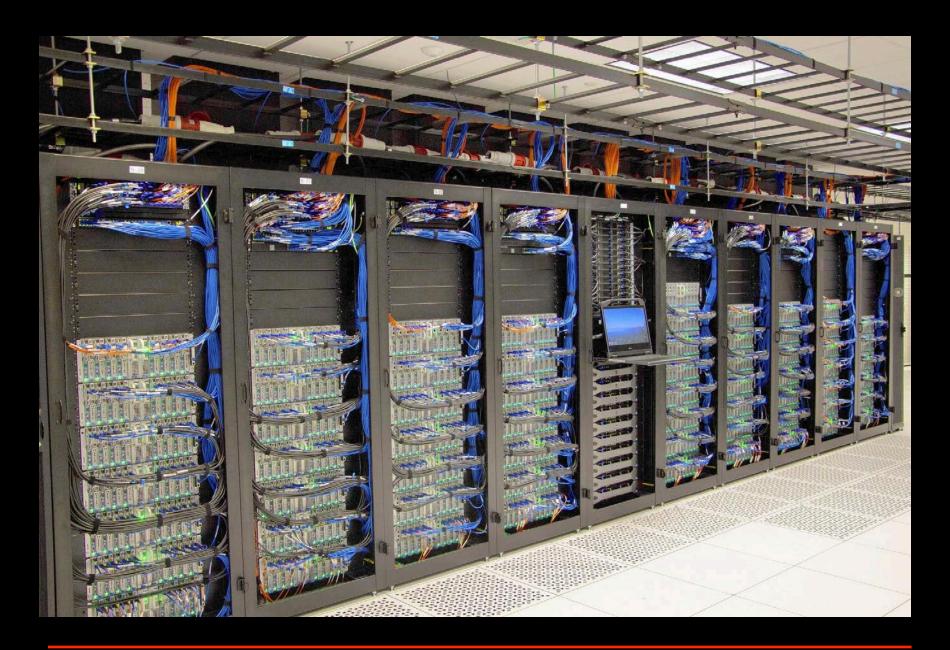- HIGHLY COLLABRATIVE!

# What it isn't…

# What is Bioinformatics?

- Bioinformatics is experimental.
- Tools and packages are under constant development and redesign.
- Best practices are only just starting to be determined.
- Always do your own checking … don't assume a program is producing valid information just because there is some output.
- Garbage in … Garbage out!

# Tools of the Trade

- Mostly open source tools.
  - These have published code that people can review, modify, and correct.
- This does not necessary mean free.
- Computers, lots of computers.

# Tools of the Trade

- Mostly open source tools.
  - These have published code that people can review, modify, and correct.
- This does not necessarily mean free.
- Computers, lots of computers.
- Mountains of Next Generation Sequencing (NGS) Data.

http://www.genomicglossaries.com/images/shenemangenome.gif

# Bioinformatic Platforms

- Linux Command Line
  - Python, Perl, R, bash, etc.
- iPlant, iAnimal etc
  - Grant funded, programmer support, intuitional support
- Galaxy
  - Heavy community support and funding.
- Commercial software Geneious, CLC, etc.

# Basic File Formats

- FASTA
- FASTQ
- SAM
- BAM

# FASTA

```
>My_gene|some description
AGAAAATAGAGAGGCCAGACGATAGATAGAGATCAGCCC
CAGACGCGCGAA
```

- Text based representation of DNA or protein sequence.
- First line starts with a > and is the sequence description.
- The next line is the sequence.
- No standard file extention
  - .fa .fasta .fas
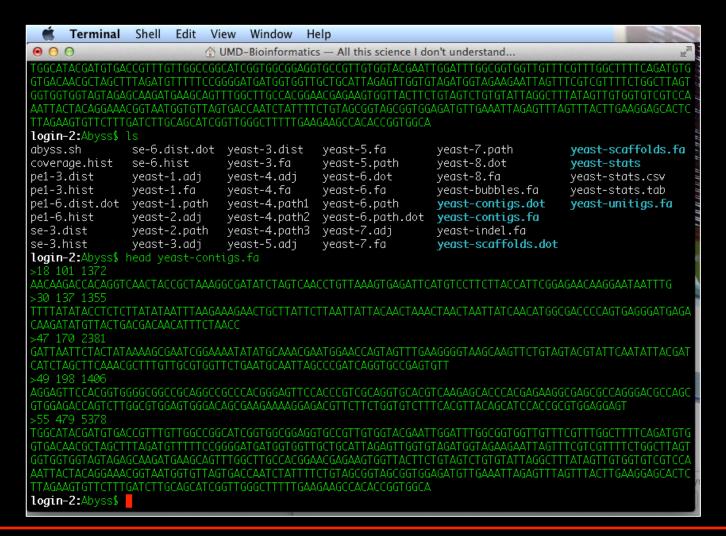
# FASTQ

- Fasta with quality information

```
@HWI-EAS225:3:1:2:854#0/1
GGGGGGAAGTCGGCAAAATAGATCCGTAACTTCGGG
+HWI-EAS225:3:1:2:854#0/1
a`abbbbabaabbababb^`[aaa`_N]b^ab^``a
@HWI-EAS225:3:1:2:1595#0/1
GGGAAGATCTCAAAAACAGAAGTAAAACATCGAACG
+HWI-EAS225:3:1:2:1595#0/1
a`abbbababbbabbbbbabb`aaababab\aa_`
```

# Plan

- That covers just the basics.

- Today we are going to work on computer skills
  - Linux
  - Python

# Linux

# Python

- http://pythonforbiologists.com/