

Institute for Systems Genomics

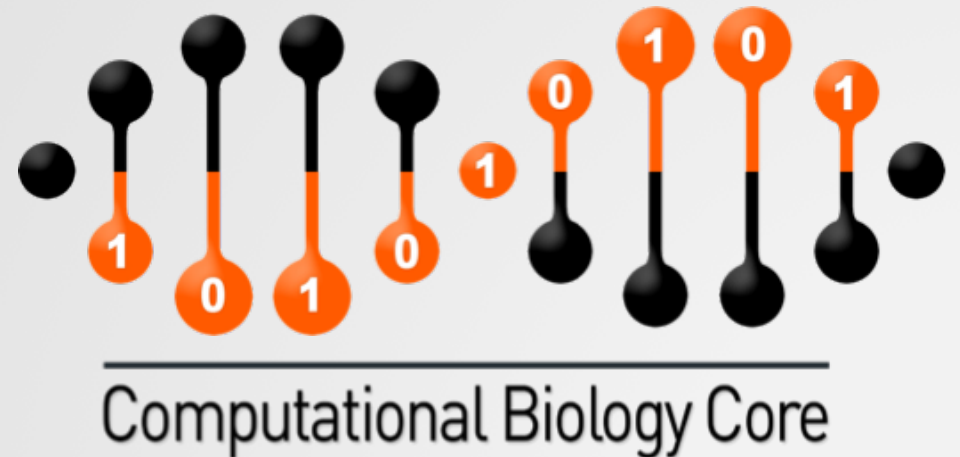
Computational Biology Core

Open House 2018

Jill Wegrzyn

Vijender Singh

Neranjana Perera



Overview of Informatics Services

- Access to HPC
 - Software/Hardware Support
- Bioinformatics
 - Data Therapy
 - Consultations
 - Workshops
 - Full Project Support



CBC Team

- **Vijender Singh** – Lead Bioinformatics Scientist
- **Neranjana Perera** – Postdoctoral Scholar
- **Mike Wilson** – HPC and Software Support

- **Ion Moraru** – Director UCH HPC Facility Director
- **Stephen King** – HPC Administrator



HPC Access



Type the commands



Actual work being done



- Large data is often too much for your small laptop or desktop to handle
- Move to large computers or clusters of computers
 - ie “remote” computers
- Vocabulary
 - “High performance computing” (HPC)



HPC Access

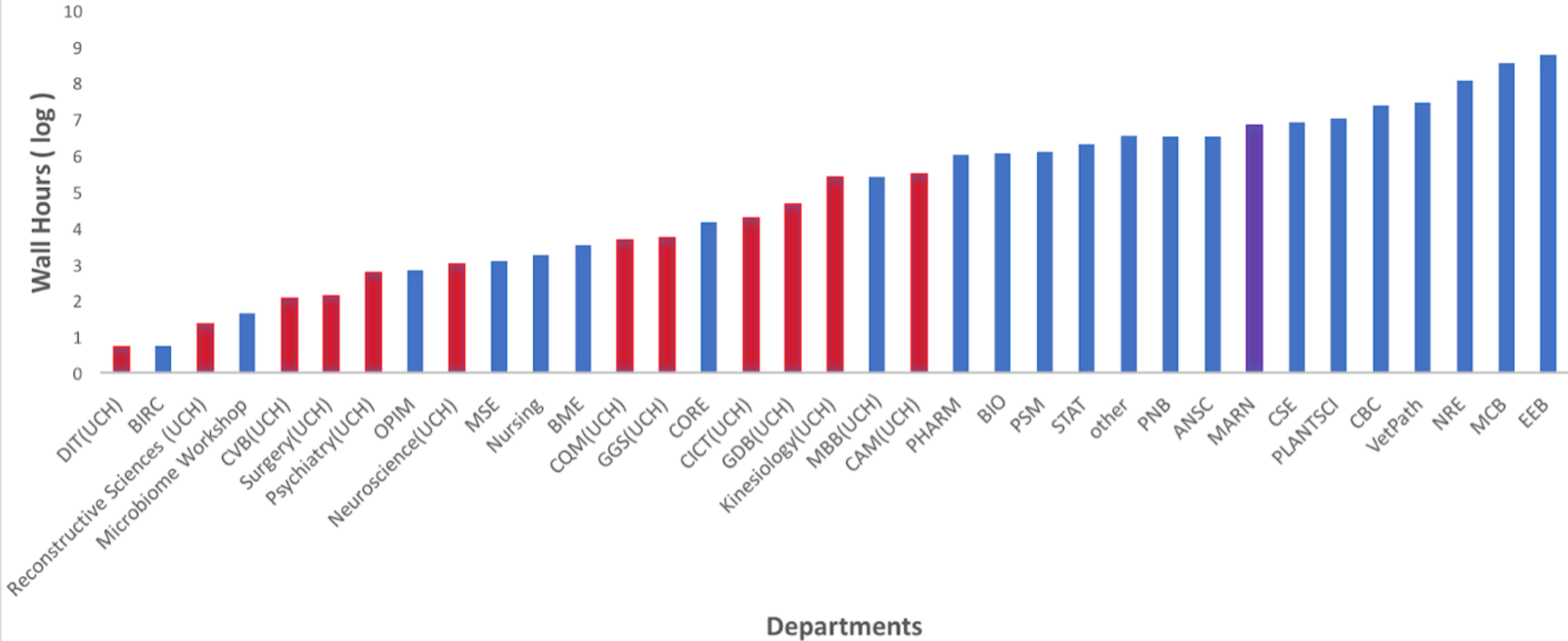
- CBC provides FREE, unlimited access to HPC computing
 - Unlimited account access for students, staff, postdocs, and faculty
 - Unlimited (and redundant) storage
- Provided through the Xanadu cluster
 - Administered by the UCH HPC Facility
 - Transition from BBC Cluster in 2018



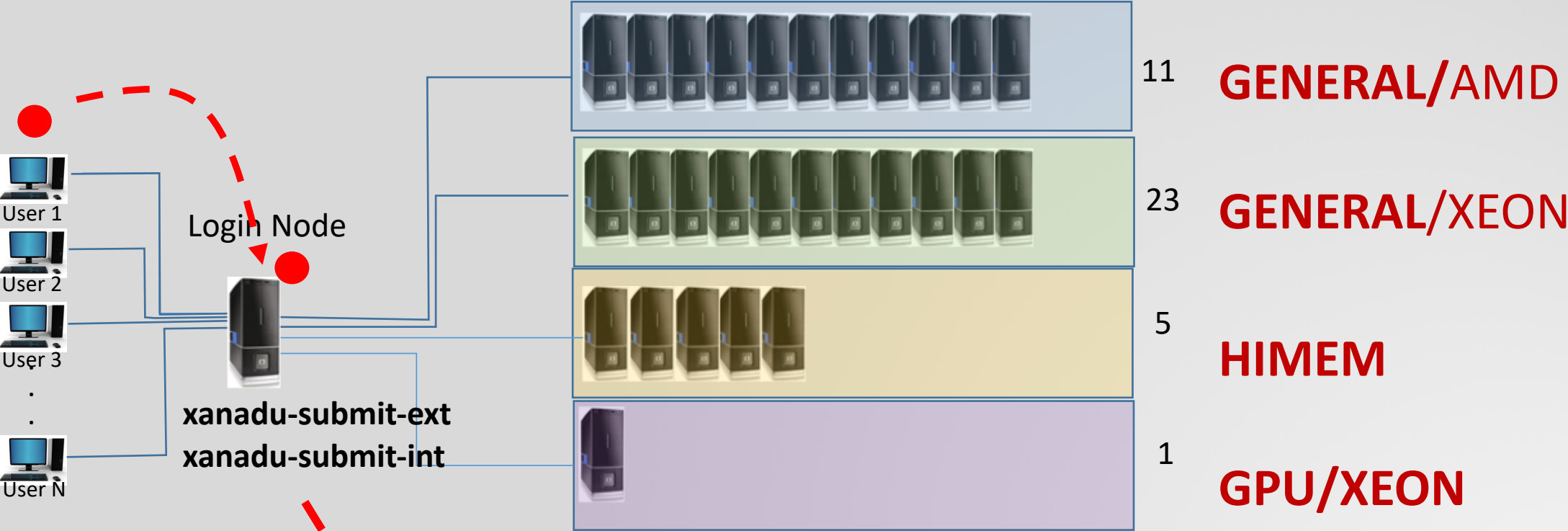
UCHC High Performance Computing Utilization (%)



WALL CLOCK TIME BY DEPARTMENT (APRIL 2017 - APRIL 2018)



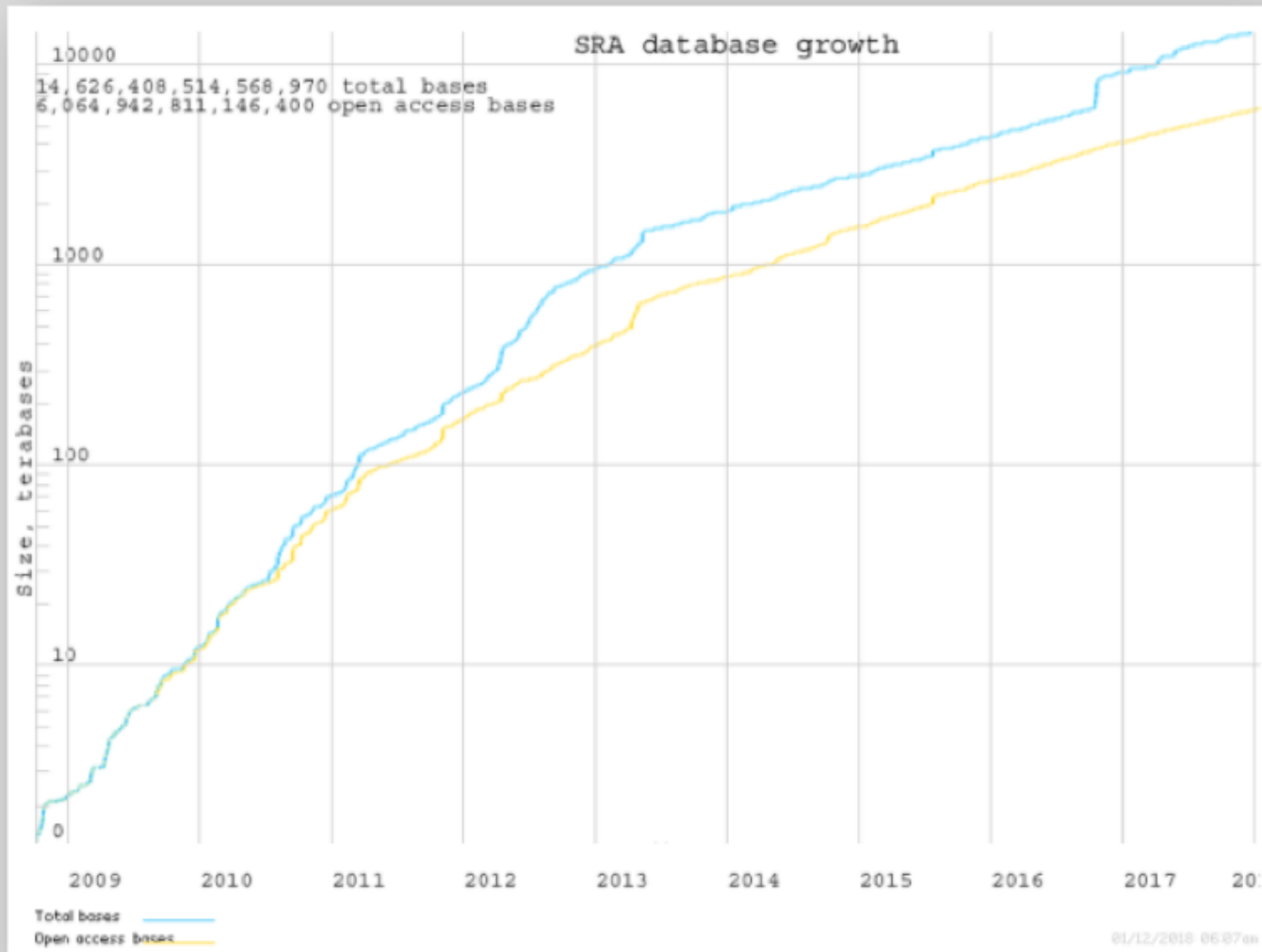
HPC Access – Xanadu Cluster – 40 nodes > 1500 cores



Xanadu Session at 12:30pm today!



HPC Access – Genomics Data



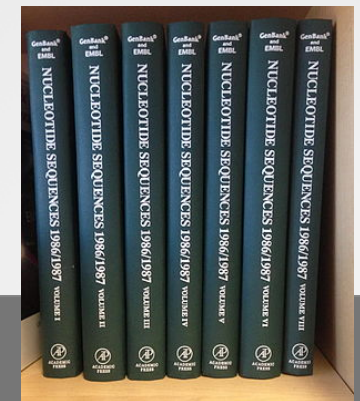
Big Data: Astronomical or Genomical?

Zachary D. Stephens, Skylar Y. Lee, Faraz Faghri, Roy H. Campbell, Chengxiang Zhai, Miles J. Efron, Ravishankar Iyer, Michael C. Schatz, Saurabh Sinha, Gene E. Robinson

Published: July 7, 2015 • <https://doi.org/10.1371/journal.pbio.1002195>

By 2025! “Our estimates show that genomics is a “four-headed beast” — it is either on par with or the most demanding of the domains analyzed here in terms of data acquisition, storage, distribution, and analysis”

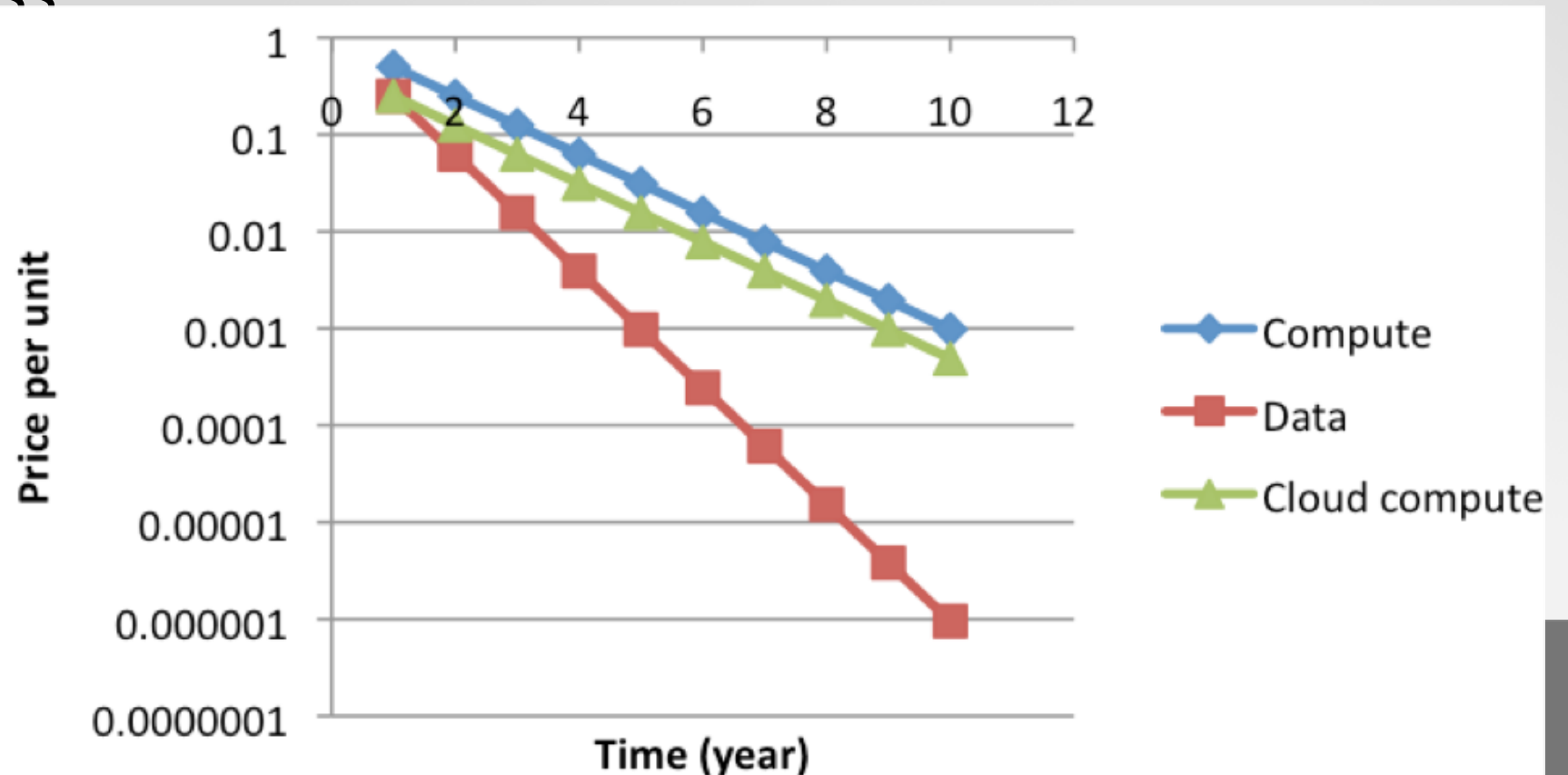
- YouTube
- Astronomy
- Twitter



Institute for Systems Genomics:
Computational Biology Core
bioinformatics.uconn.edu

HPC Access - Storage

- /linuxshare is available on Xanadu to provide archival storage
 - Not fast access
 - Redundant and
 - **Compressed**



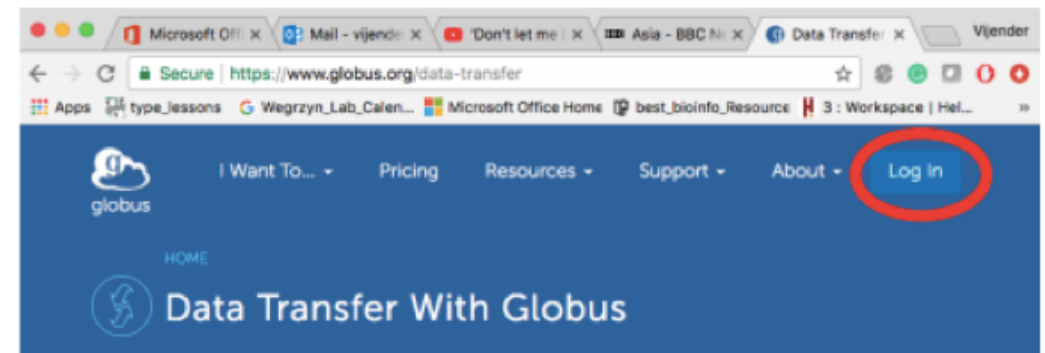
HPC Access – File Transfer

- File transfer can be challenging!
 - Globus endpoint at UCH
 - Globus endpoint on HPC Storrs
 - Drag and Drop
 - Local computer
 - FTP sites
 - Other servers
 - Not compatible?
 - Details on the website

GLOBAL FILE TRANSFER TUTORIAL

STEPS:

1. Go to <https://www.globus.org/data-transfer> and click on login option (right top).



Globus provides a secure, unified interface to your research data. Use Globus to 'fire and forget' high-performance data transfers between systems within and across organizations.

Get Started

Move files now, start a free trial, add new endpoints, and more.



Software Support

- CBC provides support for various software packages
 - Majority of bioinformatics software is available via command line access
 - Users can request software not currently installed
 - **Module** system to manage software and user knowledge of local environment

Home People Hardware **Software** Databases Resources ▾ FAQ Contact Us ▾ Open House

Software

Software available on the clusters are freely available for use. Users must be comfortable in a [Unix environment](#) and understand how to properly submit jobs. Request additional information, accounts, and access [here](#).

1. [Annotation](#)
2. [Chip Seq](#)
3. [Genome Assembly](#)
4. [Metagenomics](#)
5. [Molecular Structure](#)
6. [Multiple Sequence Aligners](#)
7. [Phylogenetics](#)
8. [Population Genetics](#)
9. [Proteomics](#)
10. [RAD-Seq](#)
11. [Repeat Analysis](#)
12. [RNA-Seq](#)
13. [Sequence Clustering](#)
14. [Sequence File Manip./Quality Control](#)
15. [Sequence Submission](#)
16. [Short Read Aligners](#)
17. [Single Cell Genomics](#)
18. [Statistics](#)
19. [Transcriptome Assembly](#)
20. [Variation Detection](#)
21. [Visualization Tools](#)



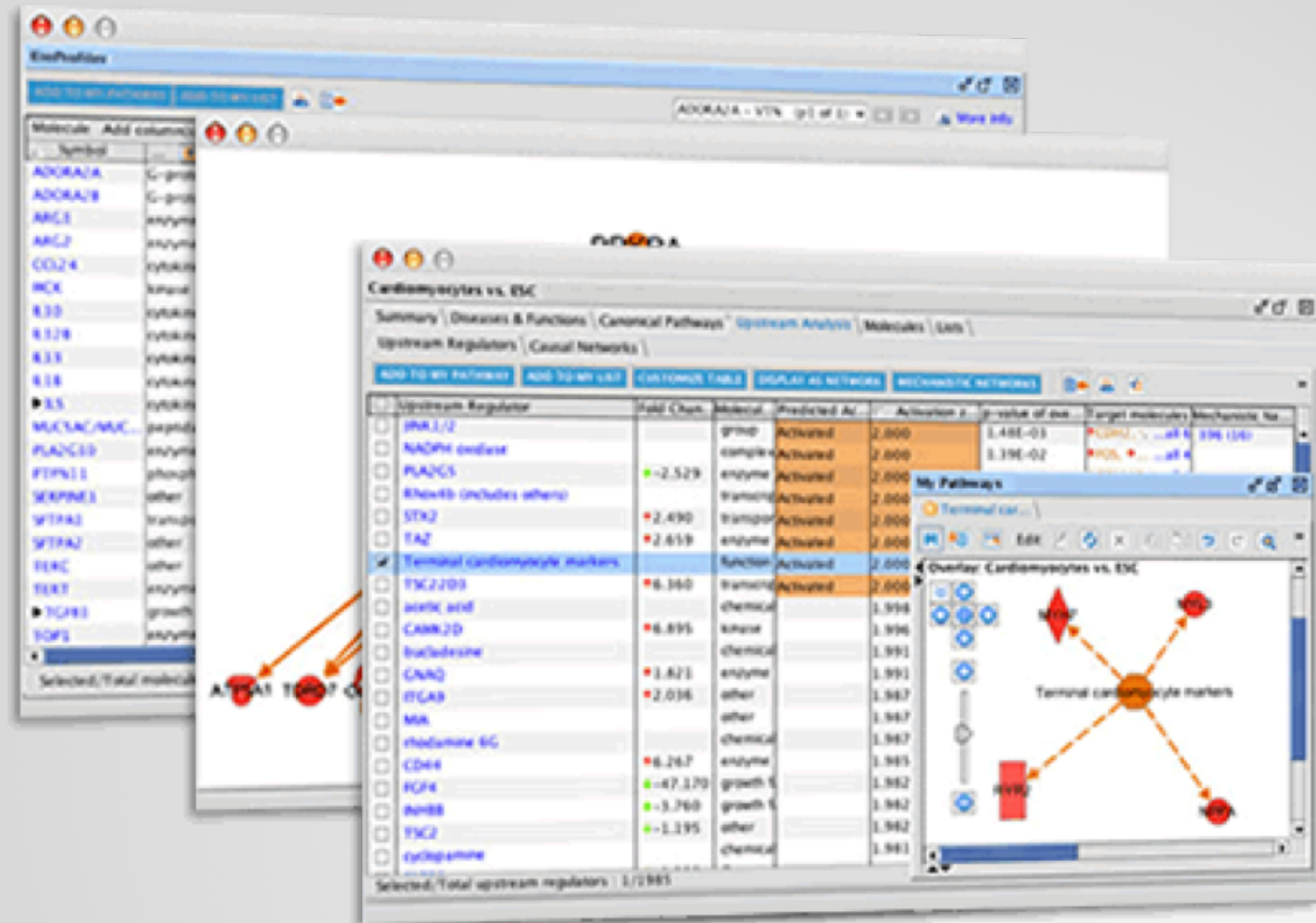
Software Support (Commercial Packages)

- CBC provides support for two commercial packages
 - **Geneious** package (server license)
 - Local instance or dedicated node with increased resources
 - Geneious Training Today!
 - **IPA**
 - Qiagen's Pathway Analysis Package
 - RNA-Seq to Pathway
 - Custom databases
 - Limited license (1 concurrent user + 5 licenses)



Software Support (Commercial Packages)

- Mechanistic Pathways
- Upstream Regulators
- Integrated Expression Results



Database Support

- Pre-Indexed Databases!
- NCBI
 - RefSeq (protein, rna, genomic)
 - NR
 - NT
- Ensembl
 - SwissProt
- Diamond (Protein)
- Genome Indexes
 - Short read aligners

Databases

NCBI BLAST database indexes are updated on a bi-weekly basis and available to users with an account on the cluster. Genome indexes for short read aligners are updated on request. Older versions are not maintained but if you require a specific version of a database or one not listed here, please contact us.

BBC cluster database path is indicated in **black** text while the paths relevant to the **Xanadu cluster** database is indicated in **purple** text.

All databases listed are current as of 04/02/2018

NCBI BLAST Databases

Database	Type	Path on Server	Description
nt	Nucleic	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	Nucleotide sequence database, with entries from all traditional divisions of GenBank, EMBL, and DDBJ excluding bulk divisions (gss, sts, pat, est, and htg divisions). wgs entries are also excluded. Not non-redundant.
nr	Protein	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	Non-redundant protein sequence database with entries from GenPept, Swissprot, PIR, PDF, PDB and NCBI RefSeq
swissprot	Protein	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	Swiss-prot sequence databases (last major update), it's parent database is nr.
env_nt	Nucleic	<code>/common/blast/data</code> <code>/isg/shared/databases/Blast</code>	Nucleotide sequences for metagenome
env_nr	Protein	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	Nucleotide sequences for metagenome
tsa_nt	Nucleic	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	Sequences from the TSA division of GenBank, EMBL, and DDBJ
16SMicrobial	Nucleic	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	Bacterial and Archaeal 16S rRNA sequences from BioProjects 33175 and 33117
refseq_rna	Nucleic	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	NCBI Transcript reference sequences
refseq_protein	Protein	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	NCBI protein reference sequences.
refseq_genomic	Nucleic	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	NCBI genomic reference sequences
taxdb	Phylogeny	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	Taxonomy
cdd_delta	Nucleic	<code>/common/blast/data</code> <code>/isg/shared/databases/blast</code>	Conserved Domain Database sequences for use with stand alone deltablast

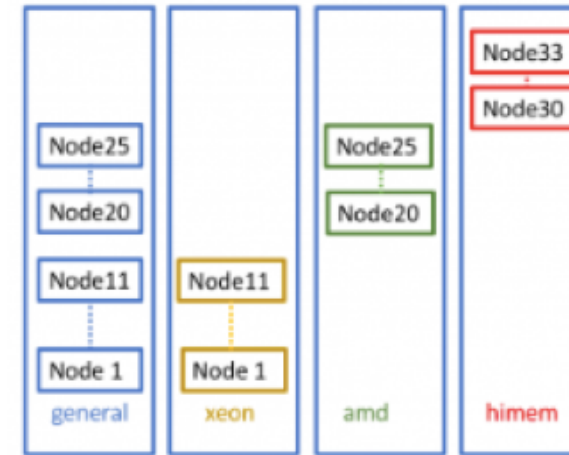


Accounts and Guides

- Request Account
 - Account needed for commercial software as well
- Xanadu Guide
 - What resources do I need?
 - Memory – fast access storage close to the CPU
 - CPU – processing units
 - Parallel applications?

HPC resources and limits

Xanadu cluster uses the Slurm, which is a highly scalable cluster management and job scheduling system for large and small Linux clusters. The nodes (individual nodes within the cluster) are divided into groups which are called partitions. Xanadu has several partitions available: general, xeon, amd, himem1, himem2, himem3, himem4.



Partitions on Xanadu

To look up the available partition information you can use 'sinfo -s' which will give you the current list:

```
$ sinfo -s
PARTITION AVAIL  TIMELIMIT  NODES(A/I/O/T)  NODELIST
general*   up        infinite   19/15/1/35     shangrila[01-18], xanadu-[01-11,20-25]
xeon       up        infinite   0/11/0/11     xanadu-[01-11]
amd        up        infinite   19/4/1/24     shangrila[01-18], xanadu-[20-25]
himem1     up        infinite   0/1/0/1       xanadu-30
himem2     up        infinite   0/1/0/1       xanadu-31
himem3     up        infinite   0/1/0/1       xanadu-32
himem4     up        infinite   0/1/0/1       xanadu-33
himem5     up        infinite   0/1/0/1       xanadu-29
```

In the above the **general** is the default partition for the users. Where NODES(A/I/O/T) are a count of a particular configuration by node state in the form of "Available / Idle / Other / Total".

Xanadu cluster is divided into six partitions:

- **general** partition
- **himem1** partition
- **himem2** partition
- **himem3** partition
- **himem4** partition
- **himem5** partition



Consultations

- CBC offers FREE consultations to members of the research community (UCH on Tuesdays!)
 - Connecting to the server for the first time?
 - Troubleshooting submission scripts
 - Configuring software
 - Basic shell scripting to assist with large jobs (moving over hundreds of files)
 - Recommending workflows and approaches
 - Installing software on the cluster
 - Troubleshooting commands and determining parameters
 - Indexing databases



Consultations

- Form available on the website to ***submit a request to meet*** or troubleshooting on specific scripts
- **SLACK channel**
 - Anyone with a uchc.edu or uconn.edu address can join!
 - #General channel allows community response and feedback
 - Good place to check on server status
 - Ask general bioinformatics questions to the community (recommended parameters)
 - Possibly not for troubleshooting more complex issues – best to send a request
- ***All requests go into the Help Desk as a ticketed/trackable item***



Request Support for your Project

- CBC approved a **rate plan** in January 2018
 - Community members interested in full/partail support for their project
 - Data storage – compress and archive
 - Data quality control and submission to Genbank
 - Alignment -> ChiP-Seq or RNA-Seq
 - Variant detection -> GBS or RAD-Seq
 - Genome assembly and annotation (microbial or eukaryotic)
 - Custom bioinformatics plans developed as well



CBC: Tutorials

CBC Provides access to step by step tutorials!

- RNA-Seq
- ChIP-Seq
- Genome Assembly
- Variant Detection (coming soon!)

Tutorials

Tutorial	Last Updated	Description
Server Access		
UConn Health Cluster (PBS)	June 2015	Understanding the UConn Health Cluster
BBC Cluster (SGE)	June 2015	Understanding the BBC Cluster
UNIX and R		
Unix Basics	November 2013	Introduction to Command Line Operations
VIM	December 2013	Introduction to VIM UNIX Editor
Unix Examples	November 2013	Basic Bioinformatic Exercises
Introduction to R	November 2013	Basic Analysis and Plots
RNA-Seq Guides		
Prokaryote RNA-Seq (EDGE-pro/DESeq2)	July 2015	EDGE-pro tutorial (with Listeria reference genome)
Model Plant RNA-Seq (STAR/DESeq2)	July 2015	RNASeq tutorial (with Glycine max reference genome)
Non-Model Plant RNA-Seq (Bowtie2/eXpress/DESeq2)	August 2015	RNA-Seq tutorial (with Picea rubens reference transcriptome)
Human RNA-Seq (no replicates) (STAR/DESeq2)	March 2015	Introduction to RNASeq
Model Insect RNASeq (Web-based Galaxy)	July 2015	RNASeq tutorial (with Drosophila reference genome)
Genome Assembly		
Genome Size Estimation Tutorial	January 2017	Genome Size Estimation Tutorial
Bacterial Genome Assembly Tutorial	September 2015	Genome Assembly tutorial



CBC: Tutorials

CBC Provides access to step by step tutorials!

- RNA-Seq
- ChIP-Seq
- Genome Assembly
- Variant Detection (coming soon!)

RNA-Seq: Reference Genome, Differential Expression, and Functional Annotation

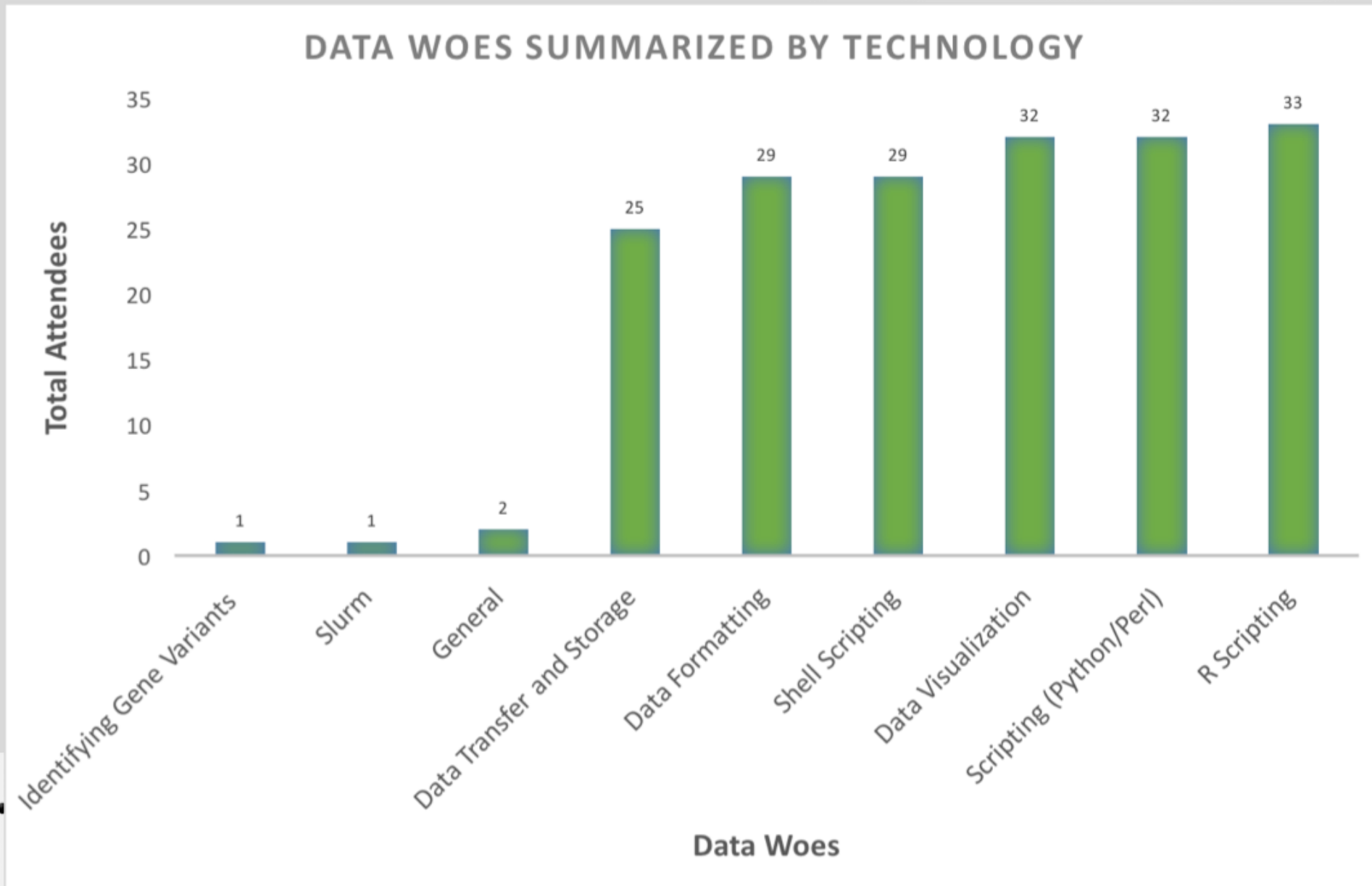
This repository is a usable, publicly available differential expression and functional annotation tutorial. All steps have been provided for the UConn CBC Xanadu cluster here with appropriate headers for the Slurm scheduler that can be modified simply to run. Commands should never be executed on the submit nodes of any HPC machine. If working on the Xanadu cluster, you should use `sbatch scriptname` after modifying the script for each stage. Basic editing of all scripts can be performed on the server with tools, such as `nano`, `vim`, or `emacs`. If you are new to Linux, please use [this](#) handy guide for the operating system commands. In this guide, you will be working with common bioinformatic file formats, such as [FASTA](#), [FASTQ](#), [SAM/BAM](#), and [GFF3/GTF](#). You can learn even more about each file format [here](#). If you do not have a Xanadu account and are an affiliate of UConn/UChC, please apply for one [here](#).

Contents

- [1 Overview and programs install](#)
- [2 Accessing the data using sra-toolkit](#)
- [3 Quality control using sickle](#)
- [4 Aligning reads to a genome using hisat2](#)
- [5 Generating total read counts from alignment using htseq-count](#)
- [6 Pairwise differential expression with counts in R with DESeq2](#)
- [7 EnTAP: Functional Annotation for Genomes](#)
- [8 Integrating the DE Results with the Annotation Results](#)
- [Citations](#)



Data Therapy Sessions



Data Therapy Sessions

- Bi-weekly discussion group at 10:30am in ESB 304
- Rotating Topics or Open Sessions
- Aimed at all levels
 - Encourage novice and advanced members for a productive discussion
- Metagenomics, Variant Detection, RNA-Seq

INSTITUTE FOR SYSTEMS GENOMICS
Computational Biology Core

Home People Hardware Software ▾ Databases Tutorials Resources ▾ **Data Therapy** Rate Plan

Data Therapy

These sessions will be held bi-weekly on Fridays.
We welcome beginners and advanced users to generate a productive discussion!
To get an idea of the group size (and your general interests) - we have a very quick survey.
This survey will also form an independent mailing list for this event so fill this out if you want to keep updated on our meetings!

CBC Events

- Data Therapy Week – 01 : RNA-Seq
- Data Therapy Week – 02 : RNA-Seq Discussion
- Data Therapy Week – 03 : Variant Detection Discussion
- Data Therapy Week – 04 : Variant Detection [WebEX Session]
- Data Therapy Week – 05 : HPC Resources and Package management [WebEX Session]
- Data Therapy Week – 06 : General Consultation [WebEX Session]
- Data Therapy Week – 07 : General Consultation [WebEX Session]



Education Mission

- Supporting Courses
- Tutorials in the Classroom
- Guest Lectures

Contact Us

Account and Support Requests

Please use this form to request an account, add software to either cluster, general bioinformatics/technical support, configure a virtual machine, or request additional cloud storage.

Inquiry Selection *

Account for Course (BBC) ▾



Workshop Offerings

CBC Workshops are offered as 2-3 day intensive sessions on topics of interest to researchers

- **RNA-Seq**
 - Experimental design
 - Quality control considerations
 - Alignment and de novo assembly strategies
- Next Workshop: **July 26-27, 2018 – Storrs Campus**



Software Carpentry



Teaching basic lab skills
for research computing

Our Lessons

Curriculum

Our lessons are developed collaboratively on [GitHub](#). You can check the status of each lesson on [our dashboard](#), or look at [older releases](#).

Availability

All of our lessons are freely available under the [Creative Commons - Attribution License](#). You may re-use and re-mix the material in any way you wish, without asking permission, provided you cite us as the original source (e.g., provide a link back to this website).

Contributing

If you have questions about contributing to particular lessons, please contact their maintainers (listed below). If you would like to develop new lessons, please see [our lesson incubation process](#).

Our lessons in English

Lesson	Site	Repository	Reference	Instructor Guide	Maintainer(s)
The Unix Shell					Gabriel Devenyi , Ashwin Srinath , Colin Morris, Will Pichters
Version Control with Git					Ivan Gonzalez , Daisie Huang , Nima Hejazi , Katherine Koziar, Madicken Munk
Version Control with Mercurial					Doug Latornell
Using Databases and SQL					Abigail Cabunoc Mayes , Jane Wyngaard, Sam Hames, Henry Senyondo
Programming with Python					Trevor Bekolay , Valentina Staneva , Anne Fouilloux, Maxim Belkin, Mike Trizna

- Open Source Curriculum
- Community Developed (and taught!)
 - Git
 - R
 - Python
 - SQL
- Frequently offered on campus
- Almost always free!

Proposal Preparation

Facilities and Other Resources

Please include the following text in research proposals to describe the computational resources available (As of June 2017):

Center for Genome Innovation

Laboratory: The Center for Genome Innovation (CGI) within the Institute for Systems Genomics on the University of Connecticut, Storrs campus, has established next generation library preparation and sequencing capacity including all necessary ancillary equipment for the following sequencing platforms: Applied Biosystems 5500xl with EZ Bead Enrichment System, two Illumina MiSeq and two Illumina NextSeq 500 systems, in addition to several Oxford Nanopore MiniONs. The laboratory is equipped with standard molecular biological equipment and resources, including thermal cyclers, standard and automated gel electrophoresis systems, centrifuges, micropipettors, water baths, incubators, refrigerators and freezers for sample and reagent storage. Ancillary equipment available for use include the: Beckman Allegra X-12R Refrigerated Centrifuge, Agilent Bioanalyzer 2100, Agilent TapeStation 2200, two PCR hoods, chemical and biosafety hoods, ABI 3500 DNA Analyzer, Affymetrix GeneAtlas and GeneChip hybridization and scanning systems, a Hydroshear, Pippin Prep, Fluidigm instruments including the C1, BioMark HD and various Access Arrays, and a Covaris S2. Genomics projects are further enhanced by the recent addition of the 10X Chromium Genomics System for long-range whole genome and exome library preparation, as well as single cell mRNA-Seq. The BioNano Irys System is also available for use in genome integrity studies looking at structural variation and haplotyping. The CGI offers both supervised and unsupervised access to instrumentation and training for use and implementation, along with an option to process samples through a fee-for-service structure. The CGI is also capable of assisting with protocol development.

The CGI is directed by Dr. Rachel O'Neill and operated by scientist, Dr. Bo Reese.

Computational Biology Core

The Computational Biology Core (CBC) within the Institute for Systems Genomics at the University of Connecticut supports bioinformatics research, teaching, and outreach. High Performance Computer (HPC) resources are housed at the Storrs campus (BBC) as well as UCHC (HPC1). All servers are freely available to members of the UConn research and their affiliates. The BBC cluster is running Centos (v6.3) and Rocks (v6.1 x64) with 1 Microway Navion Opteron Quadputers (16 x Quad-core 2.5 GHz AMD 6370P processors with 512 GB RAM), 1 Dell PowerEdge R710, 17 PowerEdge R410 nodes (2 x Quad-core 2.53 GHz Intel Xeon processors with 32 GB RAM), 4 PowerEdge R720 nodes (2 x 8-core 2.00 GHz Intel Xeon processors with 64 GB RAM), and 10 Asus RS161-E4 nodes (2 x Dual-core 2.8 GHz AMD processors with 16 GB RAM). The cluster is attached to a Dell PowerVault MD1000 disk array containing 18TB of storage (RAID 5). The cluster is connected via gigabit Ethernet internally.

The UConn Health High Performance Facility has two production compute clusters. The first has 18 nodes (Dell C6145 chassis with four 12-core AMD Opteron processors and 128 GB of RAM). It has a Dell C410x PCI expansion chassis with 16 NVIDIA M2075 GPGPU's attached to four of the compute nodes. The server is interconnected with ten-gigabit and has ten-gigabit interfaces to the public network.

The second compute cluster runs StackIQ (Rocks-) and has 14 nodes (Dell C6145 chassis with four 16-core AMD Opteron processors and either 128 GB RAM or 256 GB RAM). This cluster has gigabit interconnects and gigabit interfaces to the public network. Both clusters use a shared Isilon clustered file server of 282TB and 50 gigabit-per-second aggregate throughput.

An archival storage system is used by all clusters. It is a scalable Amplidata storage system with an Avere gateway and 1 PB in capacity with 30 gigabit-per-second aggregate throughput.

All clusters allow access to over 100 bioinformatic applications. These applications support analysis in phylogenetics, metagenomics, genome assembly, transcriptome assembly, sequence alignment, sequence annotation, variant detection, Chip-Seq, proteomics, and a variety of visualization tools.

The facility is directed by Dr. Jill Wegrzyn and operated by two scientists, Dr. Vijender Singh and Dr. Naranjan Perera. The hardware is supported by two system administrators that work across both campuses (Stephen King and Naranjan Perera).

Letter of Support

Request for Letter of Support

Facilities that you are requesting letters from *

- Computational Biology Core (CBC)
- Center for Genome Innovation (CGI)
- Proteomics Core

About You

Name *

First

Last

Email *

Institution

Department

Campus *

- Avery Point
- Farmington
- Stamford
- Storrs
- Other

Current Status *

- Undergraduate
- Graduate
- Postdoctoral
- Faculty



Institute for Systems Genomics:
Computational Biology Core
bioinformatics.uconn.edu

Stay Connected!

Online Support

Connect to the CBC team through



[#bioinformatics_help](#)

uconn-cbc.slack.com



[@Uconn_Bioinfo](#)



Bioinformatics-l@listserv.uconn.edu



Computational Biology Core

UCONN
UNIVERSITY OF CONNECTICUT