# BGGN 213

## Foundations of Bioinformatics

**Barry Grant**

UC San Diego

http://thegrantlab.org/bggn213

**05:00**

# Introduce Yourself!

Your preferred name,
Place you identify with,
Major area of study/research,
Favorite joke (optional)!

# Today's Menu

| | |
|---|---|
| **Course Logistics** | Website, screencasts, survey, ethics, assessment and grading. |
| **Learning Objectives** | What you need to learn to succeed in this course. |
| **Course Structure** | Major lecture topics and specific leaning goals. |
| **Introduction to Bioinformatis** | **Introducing the *what*, *why* and *how* of bioinformatics?** |
| **Bioinformatics Database** | **Hands-on** exploration of several major databases and their associated tools. |

# http://thegrantlab.org/bggn213/

# http://thegrantlab.org/bggn213/

# What essential concepts and skills should YOU attain from this course?



## Learning Goals

At the end of this course students will:

- Understand the increasing necessity for computation in modern life sciences research.
- Be able to use and evaluate online bioinformatics resources including major biomolecular and genomic databases, search and analysis tools, genome browsers, structure viewers, and select quality control and analysis tools to solve problems in the biological sciences.
- Be able to use the UNIX command line and the R environment to analyze bioinformatics data at scale.
- Understand the process by which genomes are currently sequenced and the bioinformatics processing and analysis required for their interpretation.
- Be familiar with the research objectives of the bioinformatics related sub-disciplines of Genomics, Transcriptomics and Structural bioinformatics.

In short, students will develop a solid foundational knowledge of bioinformatics and be able to evaluate new biomolecular and genomic information using existing bioinformatic tools and resources.

### BGGN 213

A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the Division of Biological Sciences, UCSD.

Overview

Lectures

Computer Setup

Learning Goals

Assignments & Grading

Ethics Code

Screen Cast Videos

# At the end of this course students will:

- Understand the increasing necessity for computation in modern life sciences research.

- Be able to use and evaluate online bioinformatics resources and analysis tools to solve problems in the biological sciences.

- Be able to use the R environment to analyze bioinformatics data at scale.

- Be familiar with the research objectives of the bioinformatics related sub-disciplines of Genome informatics, Transcriptomics and Structural informatics.

In short, you will develop a solid foundational knowledge of **bioinformatics** and be able to evaluate new biomolecular and genomic information using **existing bioinformatic tools and resources**.

# Specific Learning Goals....
## What I want you to know by course end!

# Course Structure
## Derived from specific learning goals

# Course Structure
## Derived from specific learning goals

# Class Details

## Goals, Class material, Screencasts & **Homework**

# Homework

## Goals, Class material, Screencasts & **Homework**

# Homework

## Goals, Class material, Screencasts & **Homework**

# Homework

## Goals, Class material, Screencasts & **Homework**

# Homework

## Goals, Class material, Screencasts & **Homework**

# Projects
## Week long **mini-projects** (x2), and 1 five week main project



bioboot.github.io/bggn213_W19/lectures/#9

Home  Gmail  Gcal  GitHub  BIMM143  BGGN213  Atmosphere  BIMM194  Blink  News ∨  + ∨

## UC San Diego

### BGGN 213

A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the **Division of Biological Sciences, UCSD** ⬀.

**Overview**

**Lectures**

**Computer Setup**

**Learning Goals**

**Assignments & Grading**

## 9: Unsupervised learning mini-project

**Topics**: Longer hands-on session with unsupervised learning analysis of cancer cells, Practical considerations and best practices for the analysis and visualization of high dimensional datasets.

**Goals**:

- Be able to import data and prepare for unsupervised learning analysis.
- Be able to apply and test combinations of PCA, k-means and hierarchical clustering to high dimensional datasets and critically review results.

**Material**:

- Lecture Slides: **To Update** Large PDF ⬀, Small PDF ⬀
- Lab: Hands-on Worksheet ⬀
- Data file: WisconsinCancer.csv ⬀, new_samples.csv ⬀.
- Bio3D PCA App: http://bio3d.ucsd.edu/pca-app/ ⬀.
- Feedback: Muddy-Point-Assesment ⬀

# Projects

## Week long **mini-projects** (x2), and 1 five week main project

**UC San Diego**

## BGGN 213

A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the Division of Biological Sciences, UCSD ↗.

**Overview**

**Lectures**

**Computer Setup**

**Learning Goals**

### 18: Cancer genomics

**Topics**: Cancer genomics resources and bioinformatics tools for investigating the molecular basis of cancer. Large scale cancer sequencing projects; NCI Genomic Data Commons; What has been learned from genome sequencing of cancer? **Immunoinformatics, immunotherapy and cancer**; Using genomics and bioinformatics to harness a patient's own immune system to fight cancer. Implications for the development of personalized medicine.

**N.B.** Find a gene assignment due before next class!

**Material:**

- Lecture Slides: Large PDF ↗, Small PDF ↗
- Lab: **T0 UPDATE** Hands-on Worksheet Part 1. ↗
- Lab: **T0 UPDATE** Hands-on Worksheet Part 2. ↗
- Data files:
  - lecture18_sequences.fa ↗,

# Projects

Week long mini-projects (x2),
and 1 five week **main project**



bioboot.github.io/bggn213_W19/lectures/#9

Home  Gmail  Gcal  GitHub  BIMM143  BGGN213  Atmosphere  BIMM194  Blink  News ▾  + ▾

## 10: Project: Find a gene assignment (Part 1)

The **find-a-gene project** ↗ is a required assignment for BIMM-143. The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.

You may wish to consult the scoring rubric at the end of the above linked project description and the **example report** ↗ for format and content guidance.

Your responses to questions Q1-Q4 are due at the beginning of class **Fri Feb 22nd** (02/22/19)).

The complete assignment, including responses to all questions, is due at the beginning of class **Wed March 13th** (03/13/19).

Late responses will not be accepted under any circumstances.

# Why Projects?

- Projects allow you to practice your new Bioinformatics skills in a less guided environment.

- In Projects, we provide datasets and ask you questions about them; just like a research project.

- Projects help build a personal portfolio and showcase your new skills, as well as help put what we have learned into practice.

# Online portfolio of **your** bioinformatics work!

# Online portfolio of **your** bioinformatics work!



jasonpbennett.github.io/bimm143/class13/NGS.html

class13     Bioinformatics Class 5

# class13

*Jason Patrick Bennett*

*May 15, 2018*

## Identifying SNP's in a Population

Lets analyze SNP's from the Mexican-American population in Los Angeles:

```
genotype <- read.csv("373531-SampleGenotypes-Homo_sapiens_Variation_Sample_rs8067378.csv")
```

Now lets look at a table of the data:

```
table(genotype)
```

```
## , , Population.s. = ALL, AMR, MXL, Father = -, Mother = -
##
##                              Genotype..forward.strand.
## Sample..Male.Female.Unknown. A|A A|G G|A G|G
##               NA19648 (F)     1   0   0   0
##               NA19649 (M)     0   0   0   1
##               NA19651 (F)     1   0   0   0
##               NA19652 (M)     0   0   0   1
##               NA19654 (F)     0   0   0   1
##               NA19655 (M)     0   1   0   0
##               NA19657 (F)     0   1   0   0
##               NA19658 (M)     1   0   0   0
##               NA19661 (M)     0   1   0   0
##               NA19663 (F)     1   0   0   0
##               NA19664 (M)     0   0   1   0
```

# Online portfolio of **your** bioinformatics work!



And finally, the fanciest graph!

```
ggplot(expr, aes(geno, exp, fill=geno)) +
  geom_boxplot(notch=TRUE, outlier.shape = NA) +
  geom_jitter(shape=16, position=position_jitter(0.2), alpha=0.4)
```

# Bonus:

## Bioinformatics & Genomics in industry

bioboot.github.io/bggn213_W19/lectures/#21

Home   Gmail   Gcal   GitHub   BIMM143   BGGN213   Atmosphere   BIMM194   Blink   News ⌄   + ⌄

## UC San Diego

### BGGN 213

A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the **Division of Biological Sciences, UCSD** ↗.

**Overview**

**Lectures**

**Computer Setup**

**Learning Goals**

**Assignments & Grading**

**Ethics Code**

### 21: Bonus: Bioinformatics & Genomics in industry

Friday March 15th at 1pm come and enjoy a set of short open ended guest lectures from leading genomic scientists at Illumina Inc., Synthetic Genomics Inc., Samumed and the La Jolla Institute for Allergy and Immunology. Come prepared for networking and to have your questions about industry careers in Bioinformatics and Genomics answered.

# Side Note: **Why stick with this course?**

**Provides a hands-on practical introduction to major bioinformatics concepts and resources.**

Covers modern hot topics and the intimate coupling of informatics with biology - highlighting the impact of computing advances and 'big data' on biology!

Designed for graduates in the biosciences with no programing experience or high level math skills.

Provides a hook for increasing computational and data science competencies in the biosciences - valuable high demand translational skills!

# Side Note: **Why stick with this course?**

**Provides a hands-on practical introduction to major bioinformatics concepts and resources.**

Covers modern hot topics and the intimate coupling of informatics with biology - highlighting the impact of computing advances and 'big data' on biology!

Designed for graduates in the biosciences with no programing experience or high level math skills.

Provides a hook for increasing computational and data science competencies in the biosciences - valuable high demand translational skills!

# BGGN-213 Learning Goals....
## Advanced UNIX and R based learning goals

# BGGN-213 Learning Goals....
## Delve deeper into "real-world" bioinformatics

# These support a major learning objective

**At the end of this course students will:**

- Understand the increasing necessity for computation in modern life sciences research.

- Be able to use and evaluate online bioinformatics resources and analysis tools to solve problems in the biological sciences.

- Be able to use UNIX and the R environment to analyze bioinformatics data at scale.

- Be familiar with the research objectives of the bioinformatics related sub-disciplines of Genome informatics, Transcriptomics and Structural informatics.

# Why use R?

Productivity
Flexibility
Genomic data analysis

# IEEE 2016 Top Programming Languages

| Language Rank | Types | Spectrum Ranking |
|---|---|---|
| 1. C | | 100.0 |
| 2. Java | | 98.1 |
| 3. Python | | 98.0 |
| 4. C++ | | 95.9 |
| 5. R | | 87.9 |
| 6. C# | | 86.7 |
| 7. PHP | | 82.8 |
| 8. JavaScript | | 82.2 |
| 9. Ruby | | 74.5 |
| 10. Go | | 71.9 |

http://spectrum.ieee.org/computing/software/the-2016-top-programming-languages

# R and Python: The Numbers

## Popularity Rankings

R and Pythons popularity between 2013 and February 2015 (Tiobe Index)



Redmonk ranking, comparing the relative performance of programming languages on GitHub and Stack Overflow (September 2012 and January 2013, 2014, 2015)



| | 2012 | 2013 | 2014 | 2015 |
|---|---|---|---|---|
| Python | 4 | 4 | 5 | 4 |
| R | 17 | 17 | 15 | 13 |

## Jobs And Salary?

2014 Dice Tech Salary Survey:
Average Salary For High Paying Skills and Experience

R $ 115,531

Python $94,139

# R is designed specifically for data analysis

- Large friendly user and developer community.

  - As of Jan 6th 2019 there are 15,352 add on **R packages** on **CRAN** and 1,823 on **Bioconductor** - much more on these later!

- Virtually every statistical technique is either already built into R, or available as a free package.

- Unparalleled data analysis environment for **high-throughput genomic data**.

# < https://www.datacamp.com/ >

# < https://www.datacamp.com/ >

# < https://www.datacamp.com/ >

# < https://www.datacamp.com/ >

## **Homework** assignments will be via DataCamp

# < https://www.datacamp.com/ >

# Today's Menu

| | |
|---|---|
| **Course Logistics** | Website, screencasts, survey, ethics, assessment and grading. |
| **Learning Objectives** | What you need to learn to succeed in this course. |
| **Course Structure** | Major lecture topics and specific leaning goals. |
| **Introduction to Bioinformatis** | Introducing the *what*, *why* and *how* of bioinformatics? |
| **Computer Setup** | Ensuring your laptop is all set for future sections of this course. |

# "What is Bioinformatics?"

"*Bioinformatics is the application of <u>computers</u> to the collection, archiving, organization, and analysis of <u>biological data</u>.*"

… A hybrid of biology and computer science

*"Bioinformatics is the application of computers to the collection, archiving, organization, and analysis of biological data."*

**Bioinformatics is computer aided biology!**

*"Bioinformatics is the application of <u>computers</u> to the collection, archiving, organization, and analysis of <u>biological data</u>."*

**Bioinformatics is computer aided biology!**

**Goal: Data to Knowledge**

# There are many useful definitions...

- "Computer based management and analysis of biological and biomedical data with useful applications in many disciplines, particularly genomics, proteomics, metabolomics, and related fields."
  (BGGN-213)

- "Bioinformatics is conceptualizing biology in terms of macromolecules and then applying "informatics" techniques (derived from disciplines such as applied maths, computer science, and statistics) to understand and organize the information associated with these molecules, on a large-scale."
  (Luscombe *et al.* 2001)

- "Bioinformatics is research, development, or application of computational approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize and analyze such data ...<cut>..."
  (National Institutes of Health:  http://tinyurl.com/l3gxr6b)

# There are many useful definitions...

- "Computer based management and analysis of biological and biomedical data with useful applications in many disciplines, particularly genomics, proteomics, metabolomics, and related fields."
  (BGGN-213)

- "Bioinformatics is conceptualizing biology in terms of macromolecules and then applying "informatics" techniques (derived from disciplines such as applied math, computer science, and statistics) to understand and organize the information associated with these molecules on a large-scale."
  (Luscombe et al. 2001)

- "Bioinformatics is research, development, or application of computational approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize and analyze such data ...<cut>..."
  (National Institutes of Health: http://tinyurl.com/l3gxr6b)

**Key Point:** Bioinformatics is Computer Aided Biology

# Major types of Bioinformatics Data

# Major types of Bioinformatics Data



**Genomes**

**Literature and ontologies**

**Gene expression**

**DNA & RNA sequence**

**Protein sequence**

**Proteom...**

**Goal:** Integrate sequence, 3D structure, expression patterns, interaction and function of biomolecules to gain a deeper understanding of biological mechanisms, processes and systems.

**Protein interactions**

**Pathways**

**Systems**

# Major types of Bioinformatics Data



Literature and ontologies

Gene expression

Genomes

Protein sequence

DNA & RNA sequence

Proteomes

Chemical entities

**Bioinformatics aims to bridge the gap between data and knowledge.**

Protein interactions

Pathways

Systems

# How do we do Bioinformatics?

- A "*bioinformatics approach*" involves the application of **computer algorithms**, **computer models** and **computer databases** with the broad goal of understanding the action of both individual genes, transcripts, proteins and <u>large collections</u> of these entities.

| DNA | → | RNA | → | Protein |
|---|---|---|---|---|
| **Genome** | → | **Transcriptome** | → | **Proteome** |

x 1,000

x 100,000

# How do we *actually* do Bioinformatics?

**Pre-packaged tools and databases**

- Many online
- Most are free to use
- Time consuming methods require downloading…

**Advanced tool application & development**

- Mostly on a **UNIX** environment
- Knowledge of programing languages frequently required (*e.g.* **R**, Python, Perl, C, Java, Fortran)
- May require specialized high performance computing…

# How do we *actually* do Bioinformatics?

**Pre-packaged tools and databases**

- Many online
- Most are free to use
- Time consuming methods require downloading…

**Advanced tool application & development**

- Mostly on a <u>UNIX</u> environment
- Knowledge of programing languages frequently required (*e.g.* **<u>R</u>**, Python, Perl, C, Java, Fortran)
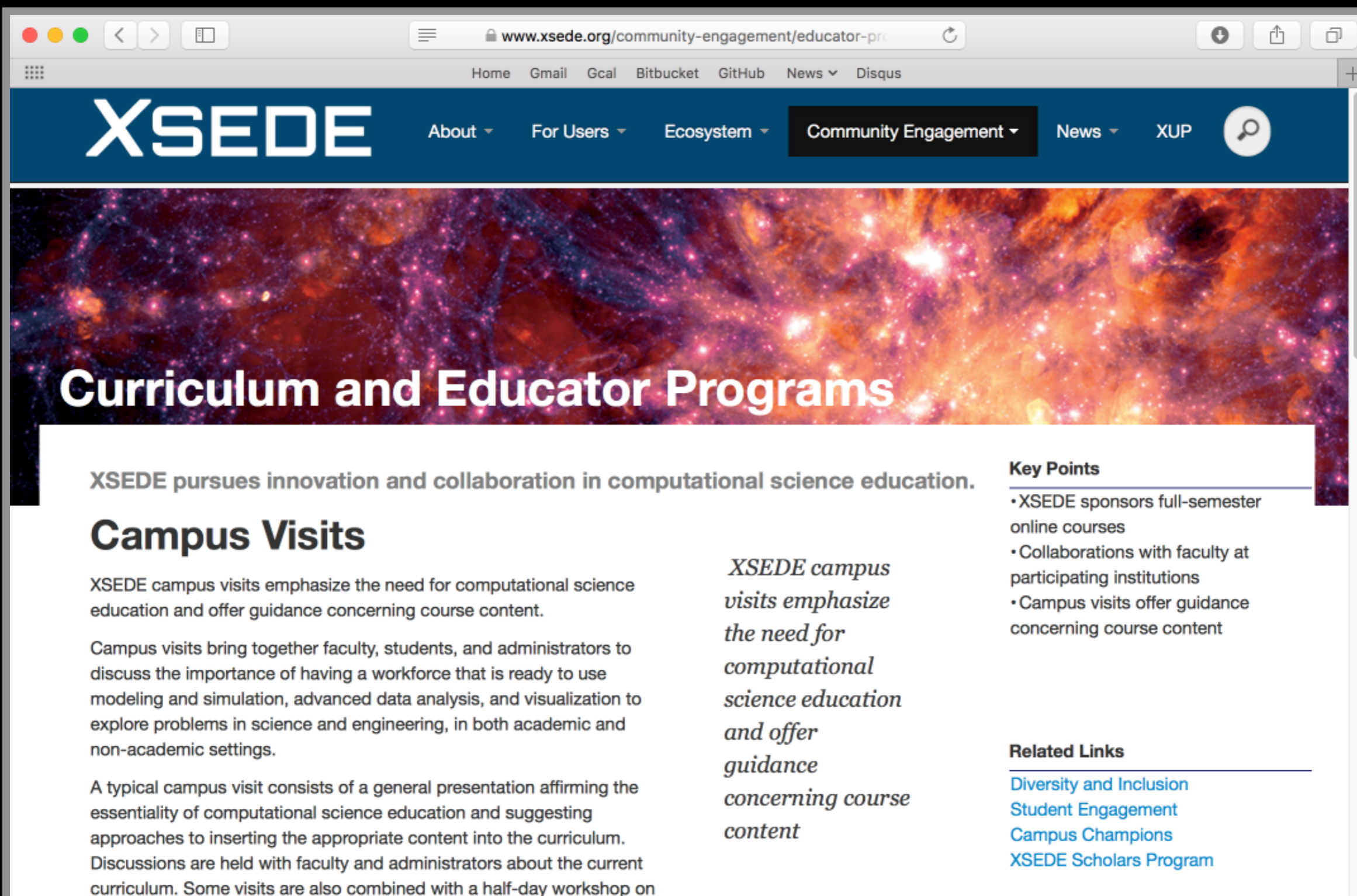- May require specialized <u>high performance computing</u>…

# NSF Extreme Science and Engineering Discovery Environment (XSEDE)

# What is *Jetstream*?

- A new cloud computing environment based at Indiana University and the Texas Advanced Computing Center (TACC) providing on-demand access to interactive computing and data analysis resources.

# Jetstream tutorials

Developed *user friendly* labs for Jetstream basics

# Jetstream tutorials

Developed *user friendly* labs for Jetstream basics

# Jetstream tutorials
## Developed *user friendly* labs for Jetstream basics

# Skepticism & Bioinformatics

We have to approach computational results the same way we do wet-lab results:

- Do they make sense?

- Is it what we expected?

- Do we have adequate controls, and how did they come out?

- Modeling is modeling, but biology is different...

  *What does this model actually contribute?*

- Avoid the miss-use of 'black boxes'

# Skepticism & Bioinformatics

Gunnar von Heijne in "*Sequence Analysis in Molecular Biology; Treasure Trove or Trivial Pursuit*" states:

➡ "Think about what you're doing; use your knowledge of the molecular system involved to guide both your interpretation of results and your direction of inquiry; use as much information as possible; and do not blindly accept everything the computer offers you".

Key-Point: **Avoid the miss-use of 'black boxes'!**

# Common problems with Bioinformatics

Confusing multitude of tools available
‣ Each with many options and settable parameters

Most tools and databases are written by and for nerds
‣ Same is true of documentation - if any exists!

Most are developed independently

Notable exceptions are found at the:
• **EBI** (European Bioinformatics Institute) and
• **NCBI** (National Center for Biotechnology Information)

Protein BLAST: search protein databases using a protein query

blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&BLAST_PROGRAMS=blastp&PAGE_TYPE=BlastSearch&SHOW_DEFAULTS=on&LINK_LOC=blasthome

### General Parameters

**Max target sequences**    500
Select the maximum number of aligned sequences to display

**Short queries**    ☑ Automatically adjust parameters for short input sequences

**Expect threshold**    10

**Word size**    3

**Max matches in a query range**    0

### Scoring Parameters

**Matrix**    BLOSUM62

**Gap Costs**    Existence: 11 Extension: 1

**Compositional adjustments**    Conditional compositional sco...

### Filters and Masking

**Filter**    ☐ Low complexity regions

**Mask**    ☐ Mask for lookup table only
☐ Mask lower case letters

### PSI/PHI/DELTA BLAST

**Upload PSSM**    Choose File    no file selected
Optional

**PSI-BLAST Threshold**    0.005

**Pseudocount**    0

Even Blast has many settable parameters

Related tools with different terminology

**STEP 3 - Set you...**

PROGRAM
FASTA

| MATRIX | GAP OPEN | GAP EXTEND | KTUP | EXPECTATION UPPER VALUE | EXPECTATION LOWER VALUE |
|---|---|---|---|---|---|
| BLOSUM50 | −10 | −2 | 2 | 10 | 0 (default) |

| DNA STRAND | HISTOGRAM | FILTER | STATISTICAL ESTIMATES |
|---|---|---|---|
| N/A | no | none | Regress |

| SCORES | ALIGNMENTS | SEQUENCE RANGE | DATABASE RANGE | MULTI HSPs |
|---|---|---|---|---|
| 50 | 50 | START-END | START-END | no |

SCORE FORMAT
Default

# Key Online Bioinformatics Resources: NCBI & EBI

The NCBI and EBI are invaluable, publicly available resources for biomedical research



http://www.ncbi.nlm.nih.gov

https://www.ebi.ac.uk

# National Center for Biotechnology Information (NCBI)

- Created in 1988 as a part of the National Library of Medicine (NLM) at the National Institutes of Health

- NCBI's mission includes:
  - Establish **public databases**
  - Develop **software tools**
  - **Education** on and dissemination of biomedical information

Bethesda, MD

- We will cover a number of core NCBI databases and software tools in the lecture

# http://www.ncbi.nlm.nih.gov

# http://www.ncbi.nlm.nih.gov

# http://www.ncbi.nlm.nih.gov



Notable NCBI databases include:
**GenBank**, **RefSeq**, PubMed, dbSNP

and the search tools **ENTREZ** and **BLAST**

# Key Online Bioinformatics Resources: NCBI & EBI

The NCBI and EBI are invaluable, publicly available resources for biomedical research



http://www.ncbi.nlm.nih.gov

https://www.ebi.ac.uk

# The EBI maintains a number of high quality curated **secondary databases** and associated tools

The EBI also provides a growing selection of **online tutorials** on EBI databases and tools



Notable EBI databases include:
ENA, **UniProt**, **Ensembl**

and the tools FASTA, BLAST, InterProScan, **MUSCLE**, DALI, **HMMER**

# Bioinformatics Databases

AATDB, AceDb, ACUTS, ADB, AFDB, AGIS, AMSdb, ARR, AsDb, BBDB, BCGD, Beanref, BioImage, BioMagResBank, BIOMDB, BLOCKS, BovGBASE, BOVMAP, BSORF, BTKbase, CANSITE, CarbBank, CARBHYD, CATH, CAZY, CCDC, CD4OLbase, CGAP, ChickGBASE, Colibri, COPE, CottonDB, CSNDB, CUTG, CyanoBase, dbCFC, dbEST, dbSTS, DDBJ, DGP, DictyDb, Picty_cDB, DIP, DOGS, DOMO, DPD, DPInteract, ECDC, ECGC, EC02DBASE, EcoCyc, EcoGene, EMBL, EMD db, ENZYME, EPD, EpoDB, ESTHER, FlyBase, FlyView, GCRDB, GDB, GENATLAS, Genbank, GeneCards, Genlilesne, GenLink, GENOTK, GenProtEC, GIFTS, GPCRDB, GRAP, GRBase, gRNAsdb, GRR, GSDB, HAEMB, HAMSTERS, HEART-2DPAGE, HEXAdb, HGMD, HIDB, HIDC, HIVdb, HotMolecBase, HOVERGEN, HPDB,  HSC-2DPAGE, ICN, ICTVDB, IL2RGbase, IMGT, Kabat,    KDNA, KEGG, Klotho, LGIC, MAD, MaizeDb, MDB, Medline, Mendel, MEROPS, MGDB, MGI, MHCPEP5  Micado, MitoDat, MITOMAP, MJDB, MmtDB, Mol-R-Us, MPDB, MRR, MutBase, MycDB, NDB, NRSub, 0-lycBase,  OMIA, OMIM, OPD, ORDB, OWL, PAHdb, PatBase, PDB, PDD, Pfam, PhosphoBase, PigBASE, PIR, PKR, PMD,  PPDB, PRESAGE, PRINTS, ProDom, Prolysis, PROSITE, PROTOMAP, RatMAP, RDP, REBASE, RGP, SBASE,   SCOP, SeqAnaiRef, SGD, SGP, SheepMap, Soybase, SPAD, SRNA db, SRPDB, STACK, StyGene, Sub2D, SubtiList, SWISS-2DPAGE, SWISS-3DIMAGE, SWISS-  MODEL Repository, SWISS-PROT, TelDB, TGN, tmRDB, TOPS, TRANSFAC, TRR, UniGene, URNADB, V BASE,   VDRR, VectorDB, WDCM, WIT, WormPep, etc ……………… !!!!

# Bioinformatics Databases

AATDB, AceDb, ACUTS, ADB, AFDB, AGIS, AMSdb, ARR, AsDb, BBDB, BCGD, Beanref, BioImage, BioMagResBank, BIOMDB, BLOCKS, BovGBASE, BOVM... ...TKbase, CANSITE, CarbBank, CARBHYD, CATH, CAZY... ...AP, ChickGBASE, Colibri, COPE, CottonDB... ...bEST, dbSTS, DDBJ, DGP, DictyDb... ...CDC, ECGC, EC02DBASE... ...THER, FlyBase... ...Link, G... ...DB, HAEMB, H... ...Hivdb, HotMolecBase, H... ...E2RGbase, IMGT, Kabat, KDNA, K... ...DB, Medline, Mendel, MEROPS, MGDB, MGI, MHC... ...OMAP, MJDB, MmtDB, Mol-R-Us, MPDB, MRR, MutBase, My... ...ub, 0-lycBase, OMIA, OMIM, OPD, ORDB, OWL, PAHdb, PatBase, PD... DD, Pfam, PhosphoBase, PigBASE, PIR, PKR, PMD, PPDB, PRESAGE, PRINTS, ProDom, Prolysis, PROSITE, PROTOMAP, RatMAP, RDP, REBASE, RGP, SBASE, SCOP, SeqAnaiRef, SGD, SGP, SheepMap, Soybase, SPAD, SRNA db, SRPDB, STACK, StyGene, Sub2D, SubtiList, SWISS-2DPAGE, SWISS-3DIMAGE, SWISS- MODEL Repository, SWISS-PROT, TelDB, TGN, tmRDB, TOPS, TRANSFAC, TRR, UniGene, URNADB, V BASE, VDRR, VectorDB, WDCM, WIT, WormPep, etc .................. !!!!

**There are lots of Bioinformatics Databases**

For a annotated listing of major bioinformatics databases please see the online handout

< **Major_Databases.pdf** >

# Side-note: **Databases come in all shapes and sizes**

Databases can be of variable quality and often there are multiple databases with overlapping content.

# Today's Menu

| | |
|---|---|
| **Course Logistics** | Website, screencasts, survey, ethics, assessment and grading. |
| **Learning Objectives** | What you need to learn to succeed in this course. |
| **Course Structure** | Major lecture topics and specific leaning goals. |
| **Introduction to Bioinformatis** | Introducing the *what*, *why* and *how* of bioinformatics? |
| **Bioinformatics Database** | **Hands-on** exploration of several major databases and their associated tools. |

# Hands-on section

http://thegrantlab.org/bggn213/

**BGGN-213: FOUNDATIONS OF BIOINFORMATICS (Lecture 1)**

**<u>Bioinformatics Databases and Key Online Resources</u>**
<u>https://bioboot.github.io/bggn213_S18/lectures/#1</u>
Dr. Barry Grant


**<u>Overview</u>:** The purpose of this lab session is to introduce a range of bioinformatics databases and associated services available on the Web whilst investigating the molecular basis of a common human disease.

Sections 1 and 2 deal with querying and searching GenBank, GENE and OMIM databases at NCBI. Sections 3 and 4 provide exposure to EBI resources for comparing proteins and visualizing protein structures. Finally, section 5 provides an opportunity to explore these and other databases further with additional examples.

**Side-note:** The Web is a dynamic environment, where information is constantly added and removed. Servers "go down", links change without warning, etc. This can lead to "broken" links and results not being returned from services. Don't give up - give it a second go and try a search engine using terms related to the page you are trying to access.


**<u>Section 1</u>**
The following transcript was found to be abundant in a human patient's blood sample.

```
>example1
ATGGTGCATCTGACTCCTGTGGAGAAGTCTGCCGTTACTGCCCTGTGGGGCAAGGTGAACGTGGATGAAG
TTGGTGGTGAGGCCCTGGGCAGGCTGCTGGTGGTCTACCCTTGGACCCAGAGGTTCTTTGAGTCCTTTGG
GGATCTGTCCACTCCTGATGCAGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAGAAAGTGCTCGGT
GCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGTGAGCTGCACT
GTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGCTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCA
TCACTTTGGCAAAGAATTCACCCCACCAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAAT
GCCCTGGCCCACAAGTATCACTAAGCTCGCTTTCTTGCTGTCCAATTT
```


The only information you are given is the above sequence so you must begin your investigation with a sequence search - for this example we will use NCBI's **BLAST** service at: <u>http://blast.ncbi.nlm.nih.gov/</u>

*Note that there are several different "basic BLAST" programs available at NCBI (including nucleotide BLAST, protein BLAST, and BLASTx).*

# YOUR TURN!

- There are five major hands-on sections including:

    1. BLAST, GenBank and OMIM @ **NCBI**    [~35 mins]
    2. GENE database @ **NCBI**    [~15 mins]
        — BREAK —
    3. UniProt & Muscle @ **EBI**    [~25 mins]
    4. PFAM, PDB & NGL    [~30 mins]
        — BREAK —
    5. Extension exercises    [~30 mins]

    ‣ Please do answer the last review question (**Q19**).

    ‣ We encourage discussion and exploration!

# YOUR TURN!

- There are five major hands-on sections including:

  End times:

  1. BLAST, GenBank and OMIM @ **NCBI**  [2:35 pm]
  2. GENE database @ **NCBI**  [2:55 pm]
     — BREAK —  — 3:10 pm —
  3. UniProt & Muscle @ **EBI**  [3:30 pm]
  4. PFAM, PDB & NGL  [4:00 pm]
     — BREAK —  — 4:10 pm —
  5. Extension exercises  [4:40 pm]

  ‣ Please do answer the last review question (**Q19**).
  ‣ We encourage discussion and exploration!

# SUMMARY

- Bioinformatics is computer aided biology.

- Bioinformatics deals with the collection, archiving, organization, and interpretation of a wide range of biological data.

- The NCBI and EBI are major online bioinformatics service providers.

- Introduced Gene, UniProt, PDB databases as well as a number of 'boutique' databases including PFAM and OMIM.

# HOMEWORK

☑ Complete the initial **course questionnaire**:

☑ Check out the "**background reading**" material online:

☑ Complete the lecture 1 **homework questions**: