# BIMM 143

## Introduction to Bioinformatics

**Barry Grant**
UC San Diego

http://thegrantlab.org/bimm143

---

**HELLO** my name is

*BARRY*

bjgrant@ucsd.edu

**Office Hours:**
SignUp

**Location:**
*TATA, #2501*

**HELLO** ~~HER~~ my name is

*YAN*

yal069@ucsd.edu

**HELLO** ~~HER~~ my name is

*NHIEN*

n7nguyen@ucsd.edu

---

# 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/bimm143/

**Bioinformatics (BIMM 143, Fall 2018)**

Course Director
Prof. Barry J. Grant (Email: bjgrant@ucsd.edu)

Instructional Assistant
Chao Shi (Email: bioshichao@gmail.com)

Course Syllabus
Fall 2018 (PDF)

**Overview**

Bioinformatics - the application of computational and analytical methods to biological problems - is a rapidly maturing field that is driving the collection, analysis, and interpretation of the avalanche of data in modern life sciences and medical research.

This upper division 4-unit course is designed for biology majors and provides an introduction to the principles and practical approaches of bioinformatics as applied to genes and proteins.

---

http://thegrantlab.org/bimm143/

(Same content repeated with "Learning Goals" highlighted)

---

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 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.

**Specific 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 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

# Homework (35% of course grade)

## Goals, Class material, Screencasts & Homework

BIMM143 Lecture 1 Homework

Please answer the following questions ... email address and UCSD PID number so you can ...

**Homework is due before the next weeks class!**

Email address *

Your email

UCSD PID number (exam number)

Your answer

Which of the following operating systems is most frequently used for bioinformatics tool development

1 point

---

# Projects

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

### UC San Diego

**BIMM 143**

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

### 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: Large PDF, Small PDF,
- Lab: Hands-on section worksheet for PCA
- Data file: WisconsinCancer.csv, new_samples.csv.
- Bio3D PCA App: http://bio3d.ucsd.edu/pca-app/.
- Feedback: Muddy point assessment.
- Bonus: Kevin's StackExchange Link on PCA.

---

# Projects

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

### UC San Diego

**BIMM 143**

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

### Designing a personalized cancer vaccine

**BIMM-143 Lecture 18:**
Barry Grant < http://thegrantlab.org >
Date: 2018-03-07 (15:24:21 PST on Wed, Mar 07)

**Notes**: To identify somatic mutations in a tumor, DNA from the tumor is sequenced and compared to DNA from normal tissue in the same individual using *variant calling algorithms*.

Comparison of tumor sequences to those from normal tissue (rather than 'the human genome') is important to ensure that the detected differences are not germline mutations.

To identify which of the somatic mutations leads to the production of aberrant proteins, the location of the mutation in the genome is inspected to identify non-

---

# Projects (20% of course grade)

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

### UC San Diego

**BIMM 143**

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

### 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.

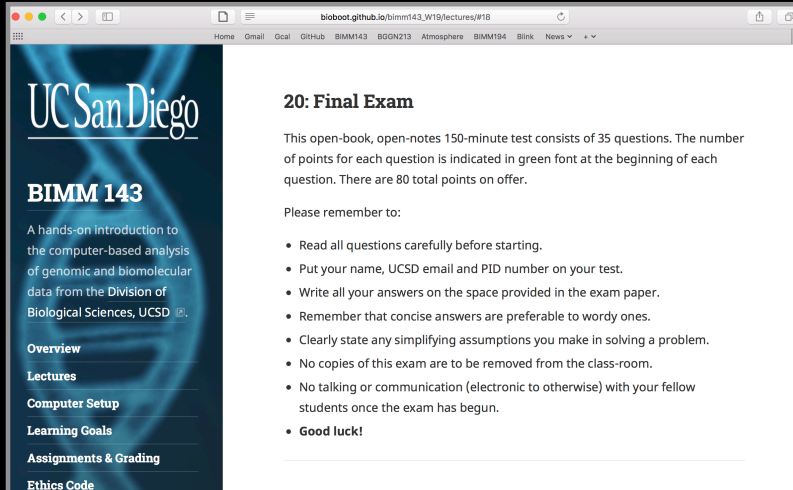Your responses to questions Q1-Q4 are due at the beginning of class **Thursday Nov 15th** (11/15/18).

The complete assignment, including responses to all questions, is due at the beginning of class **Thursday Dec 4th** (12/04/18).

Late responses will not be accepted under any circumstances.
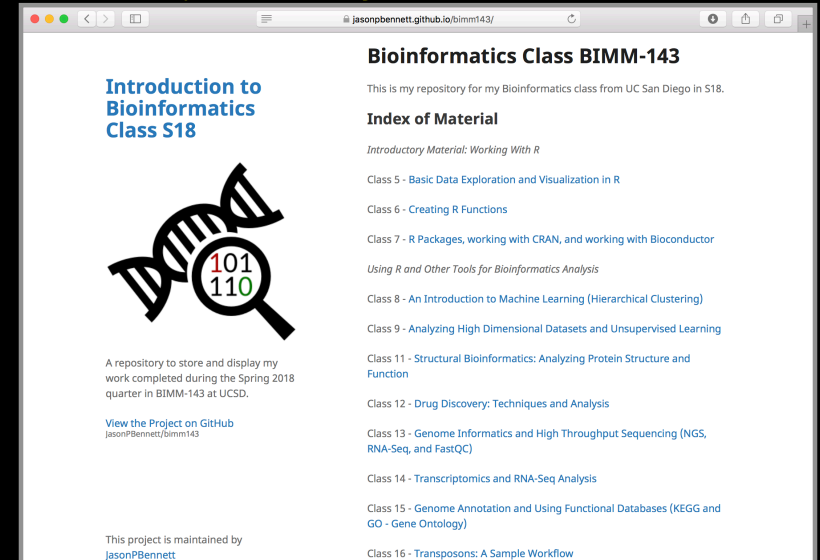
# Final Exam

## Open-book, open-notes 150-minute test

### (45% of course grade)

UC San Diego

**BIMM 143**

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

### 20: Final Exam

This open-book, open-notes 150-minute test consists of 35 questions. The number of points for each question is indicated in green font at the beginning of each question. There are 80 total points on offer.

Please remember to:

- Read all questions carefully before starting.
- Put your name, UCSD email and PID number on your test.
- Write all your answers on the space provided in the exam paper.
- Remember that concise answers are preferable to wordy ones.
- Clearly state any simplifying assumptions you make in solving a problem.
- No copies of this exam are to be removed from the class-room.
- No talking or communication (electronic to otherwise) with your fellow students once the exam has begun.
- **Good luck!**

---

# Bonus:

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

### Introduction to Bioinformatics Class S18

A repository to store and display my work completed during the Spring 2018 quarter in BIMM-143 at UCSD.

View the Project on GitHub
JasonPBennett/bimm143

This project is maintained by JasonPBennett

### Bioinformatics Class BIMM-143

This is my repository for my Bioinformatics class from UC San Diego in S18.

#### Index of Material

*Introductory Material: Working With R*

Class 5 - Basic Data Exploration and Visualization in R

Class 6 - Creating R Functions

Class 7 - R Packages, working with CRAN, and working with Bioconductor

*Using R and Other Tools for Bioinformatics Analysis*

Class 8 - An Introduction to Machine Learning (Hierarchical Clustering)

Class 9 - Analyzing High Dimensional Datasets and Unsupervised Learning

Class 11 - Structural Bioinformatics: Analyzing Protein Structure and Function

Class 12 - Drug Discovery: Techniques and Analysis

Class 13 - Genome Informatics and High Throughput Sequencing (NGS, RNA-Seq, and FastQC)

Class 14 - Transcriptomics and RNA-Seq Analysis

Class 15 - Genome Annotation and Using Functional Databases (KEGG and GO - Gene Ontology)

Class 16 - Transposons: A Sample Workflow

---

# Bonus:

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

### 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)   0   1   0   0
##               NA19664 (M)   0   0   1   0
```

---

# Bonus:

## 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)
```

geno
A/A
A/G
G/G

**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 biology majors 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 biology majors 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!

---

# BIMM-143 Learning Goals….
## Data science R based learning goals



---

# BIMM-143 Learning Goals….
## Delve deeper into "real-world" bioinformatics

## Slide 1

**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 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.

## Slide 2

# Why use R?

Productivity
Flexibility
Genomic data analysis

## Slide 3

### 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

## Slide 4

### R and Python: The Numbers

**Popularity Rankings**

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

Python

R

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

http://www.kdnuggets.com/2015/05/r-vs-python-data-science.html?utm_medium=email&utm_source=flipboard

## Slide 1

- R is the "lingua franca" of data science in industry and academia and was designed specifically for data analysis.

- Large friendly user and developer community.
  - As of Jan 6th 2019 there are 13,645 add on **R packages** on **CRAN** and 1,649 on **Bioconductor** - 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**.

## Slide 2

# Past Student Opinions...



## Slide 3

# Past Student Opinions...



## Slide 4

# Past Student Opinions...

## 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. |

---

**Q.** What is Bioinformatics?

---

**Q.** What is Bioinformatics?

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

… Bioinformatics is a hybrid of biology and computer science

---

**Q.** What is Bioinformatics?

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

… Bioinformatics is a hybrid of biology and computer science
… **Bioinformatics is computer aided biology!**

## Slide 1

**Q.** What is Bioinformatics?

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

… Bioinformatics is a hybrid of biology and computer science

… **Bioinformatics is computer aided biology!**

Computer based management and analysis of biological and biomedical data with useful applications in many disciplines, particularly genomics, proteomics, metabolomics, etc...

## Slide 2

## MORE DEFINITIONS

▸ "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 NM, *et al.* Methods Inf Med. 2001;40:346.

▸ "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." National Institutes of Health (NIH)  ( http://tinyurl.com/l3gxr6b )

## Slide 3

## MORE DEFINITIONS

▸ "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 analyze the information associated with these molecules, on a **large-scale**. Luscombe NM, *et al.* Methods Inf Med. 2001;40:346.

▸ "Bioinformatics 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." National Institutes of Health (NIH)  ( http://tinyurl.com/l3gxr6b )

*Key Point: Bioinformatics is Computer Aided Biology*

## Slide 4

# Major types of Bioinformatics Data

Literature and ontologies

Gene expression

Genomes

Protein sequence

DNA & RNA sequence

Protein structure

Proteomes

Chemical entities

Protein families, motifs and domains

Protein interactions

Pathways

Systems

## Major types of Bioinformatics Data



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

Labels: Literature and ontologies, Genomes, Gene expression, DNA & RNA sequence, Protein sequence, Proteomes, Protein interactions, Pathways, Systems

## Major types of Bioinformatics Data



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

Labels: Literature and ontologies, Genomes, Gene expression, DNA & RNA sequence, Protein sequence, Proteomes, chemical entities, Protein interactions, Pathways, Systems

## BIOINFORMATICS RESEARCH AREAS

Include but are not limited to:

- Organization, classification, dissemination and analysis of biological and biomedical data (particularly '-omics' data).
- Biological sequence analysis and phylogenetics.
- Genome organization and evolution.
- Regulation of gene expression and epigenetics.
- Biological pathways and networks in healthy & disease states.
- Protein structure prediction from sequence.
- Modeling and prediction of the biophysical properties of biomolecules for binding prediction and drug design.
- Design of biomolecular structure and function.

With applications to Biology, Medicine, Agriculture and Industry

## Where did bioinformatics come from?

Bioinformatics arose as molecular biology began to be transformed by the emergence of molecular sequence and structural data

**Recap: The key dogmas of molecular biology**
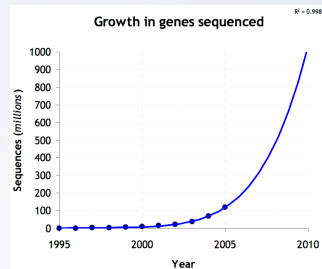
- *DNA sequence* determines *protein sequence*.
- *Protein sequence* determines *protein structure*.
- *Protein structure* determines *protein function*.
- *Regulatory mechanisms* (e.g. gene expression) determine the amount of a particular *function in space and time*.

Bioinformatics is *now* essential for the archiving, organization and analysis of data related to all these processes.

## Why do we need Bioinformatics?

Bioinformatics is necessitated by the rapidly expanding quantities and complexity of biomolecular data

- Bioinformatics provides methods for the efficient:
  ‣ **storage**
  ‣ **annotation**
  ‣ **search** and **retrieval**
  ‣ data **integration**
  ‣ data **mining** and **analysis**
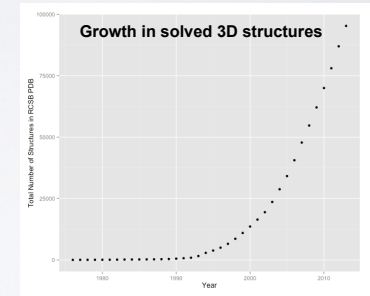


Growth in genes sequenced

E.G. data from sequencing, structural genomics, proteomics, new high throughput assays, *etc…*

---

## Why do we need Bioinformatics?

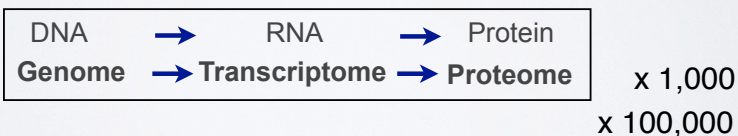Bioinformatics is necessitated by the rapidly expanding quantities and complexity of biomolecular data

- Bioinformatics provides methods for the efficient:
  ‣ **storage**
  ‣ **annotation**
  ‣ **search** and **retrieval**
  ‣ data **integration**
  ‣ data **mining** and **analysis**



Growth in solved 3D structures

E.G. data from sequencing, structural genomics, proteomics, new high throughput assays, *etc…*

---

## 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 large collections 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…
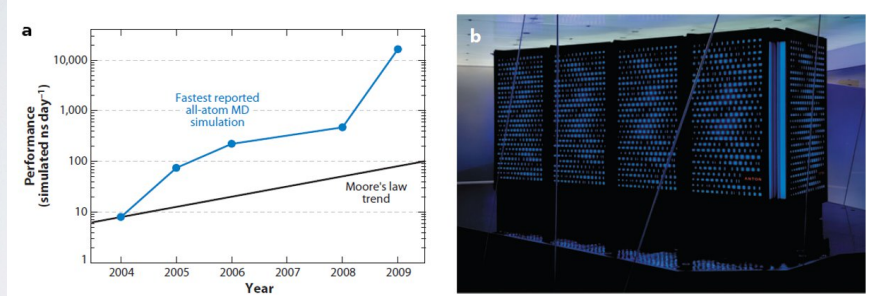
## How do we *actually* do Bioinformatics?

**Pre-packaged tools and databases**

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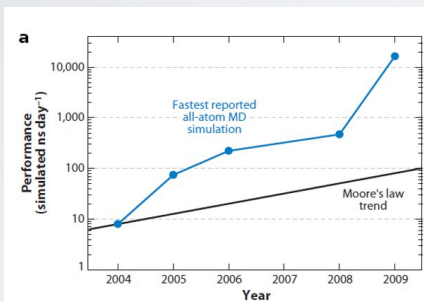- Time consuming methods require downloading…

**Advanced tool application & development**

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

---

## SIDE-NOTE: SUPERCOMPUTERS AND GPUS



---

## SIDE-NOTE: SUPERCOMPUTERS AND GPUS



---

## 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"* 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)



Even Blast has many settable parameters

Related tools with different terminology

# 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**

**Slide 1:**

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

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

**Slide 2:**

# 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

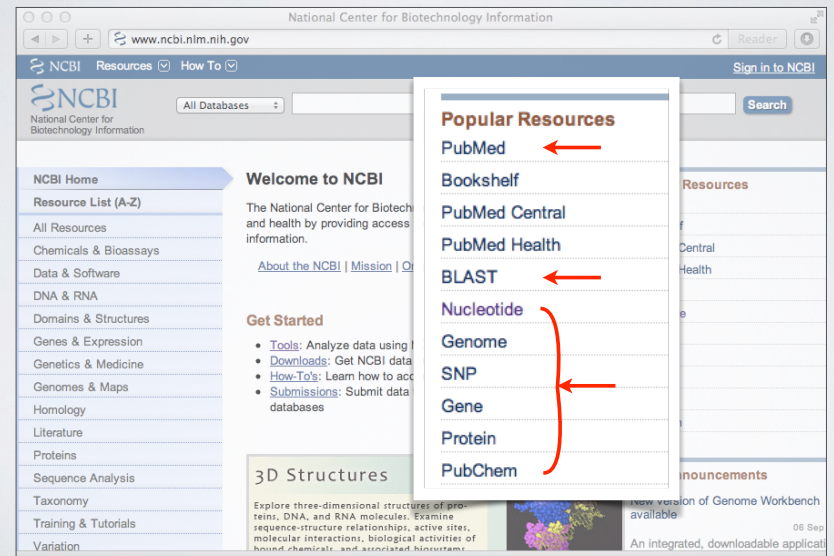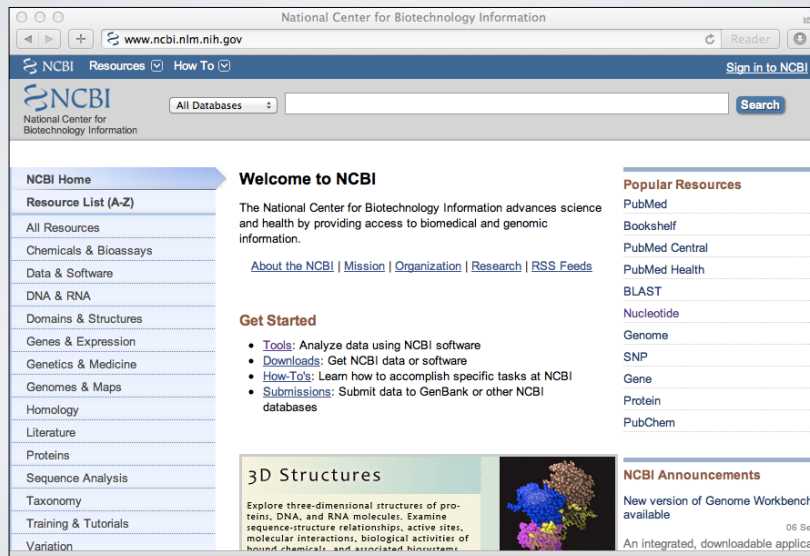**Slide 3:**

# European Bioinformatics Institute (EBI)

- Created in 1997 as a part of the European Molecular Biology Laboratory (EMBL)

- EBI's mission includes:
  ‣ providing freely available **data** and **bioinformatics services**
  ‣ and providing advanced **bioinformatics training**

Hinxton, UK

- We will briefly cover several EBI databases and tools that have advantages over those offered at NCBI

**Slide 4:**

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

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



**https://www.ebi.ac.uk**

The EBI makes available a wider variety of **online tools** than NCBI



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



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

## Slide 1 (top-left)

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**

## Slide 2 (top-right)

# 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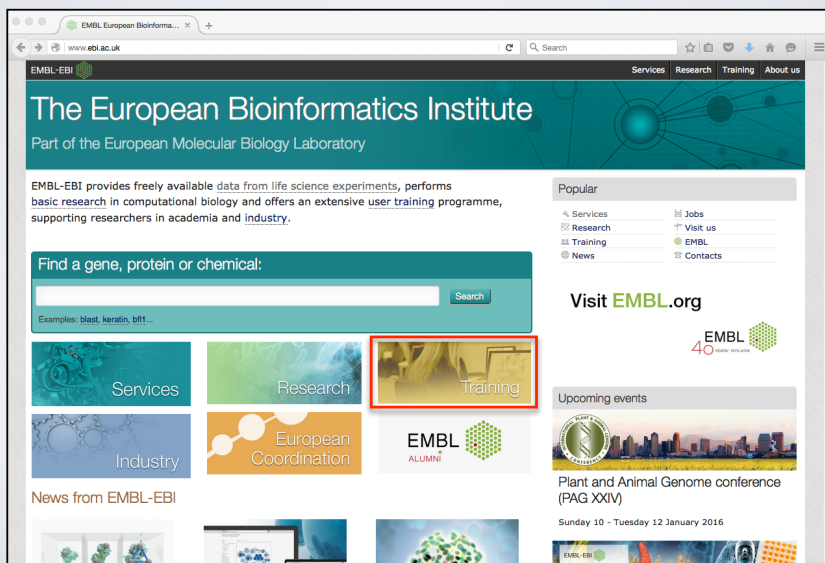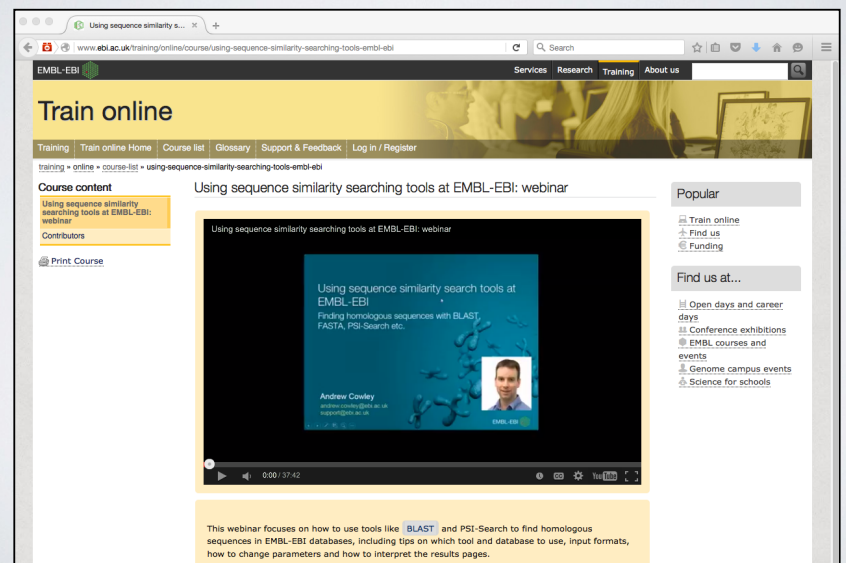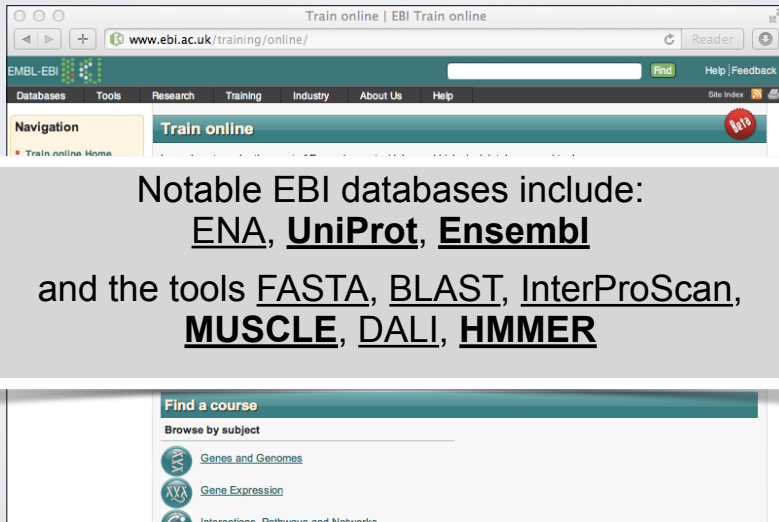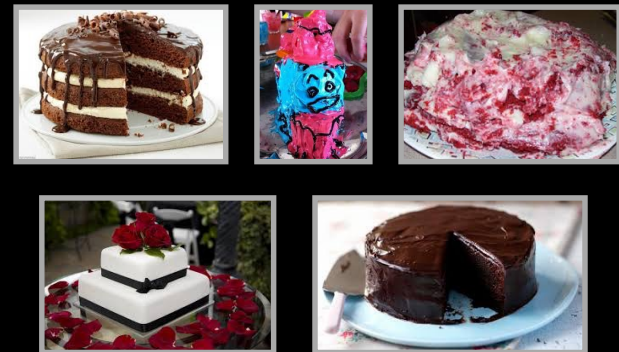, EPD, EpoDB, ESTHER, FlyBase, FlyView, GCRDB, GDB, GENATLAS, Genbank, GeneCards, Genlilesne, GenLink, GENOTK, GenProtEC, GIFTS, GPCRDB, GRAP, GRBase, gRNAsdb, GRR, GSDB, HAEMB, HAMSTERS, HEART-2DPAGE, HEXAdb, HGMD, HIDB, HIDC, HIVdb, HotMolecBase, HOVERGEN, HPDB, HSC-2DPAGE, ICN, ICTVDB, IL2RGbase, 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 ................. !!!!

## Slide 3 (bottom-left)

# Bioinformatics Databases

AATDB, AceDb, ACUTS, ADB, AFDB, AGIS, AMSdb, ARR, AsDb,BBDB, BCGB, 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, Fl... Link, G... b, HAEMB, H... ivdb, HotMolecBase, H... 2RGbase, IMGT, Kabat, KDNA, K... Medline, Mendel, MEROPS, MGDB, MGI, MH... OMAP, MJDB, MmtDB, Mol-R-Us, MPDB, MRR, MutBase, Myc... 0-lycBase, OMIA, OMIM, OPD, ORDB, OWL, PAHdb, PatBase, PD... 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 >

## Slide 4 (bottom-right)

# 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. |

---

## Your Turn!

https://bioboot.github.io/bimm143_S19/lectures/#1

**1: Welcome to Foundations of Bioinformatics**

**Topics:**
Course introduction, Leaning goals & expectations, Biology is an information science, History of Bioinformatics, Types of data, Application areas and introduction to upcoming course segments, Student 30-second introductions, Student computer setup.

**Goals:**
- Understand course scope, expectations, logistics and ethics code.
- Understand the increasing necessity for computation in modern life sciences research.
- Get introduced to how bioinformatics is practiced.
- Complete the pre-course questionnaire.
- Setup your laptop computer for this course.

**Material:**
- Lecture Slides: Large PDF, Small PDF.
- Lab: Hands-on section worksheet
- Feedback: Muddy Point Assessment

- Handout: Class Syllabus
- Computer Setup Instructions.

**BIMM 143**
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

---

**BIMM-143: INTRODUCTION TO BIOINFORMATICS (Lecture 1)**

**Bioinformatics Databases and Key Online Resources**
https://bioboot.github.io/bimm143_W18/lectures/#1
Dr. Barry Grant
Jan 2018

**Overview:** 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.

**Section 1**
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: http://blast.ncbi.nlm.nih.gov/

*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** | [11:05 am] |
| 2. | GENE database @ **NCBI** | [11:25 am] |
|    | — BREAK — | — 11:35 am — |
| 3. | UniProt & Muscle @ **EBI** | [12:00 am] |
| 4. | PFAM, PDB & NGL | [12:30 pm] |
|    | — BREAK — | |
| 5. | Extension exercises | |

- ‣ 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

https://bioboot.github.io/bimm143_S19/lectures/#1

- ☑ Complete the **initial course questionnaire**:

- ☑ Check out the "**Background Reading**" material online:

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