# CPBS 7711 Course Review

Presented by Evan Stene

# Topics Covered

- Databases
- Genetics
- Machine Learning
- Protein Folding
- Phylogenetic Trees
- Visualization
- And More…

# MolBio Databases

- Biological data is very heterogeneous

- Data is scattered among thousands of databases

- Much of the data stored in redundant

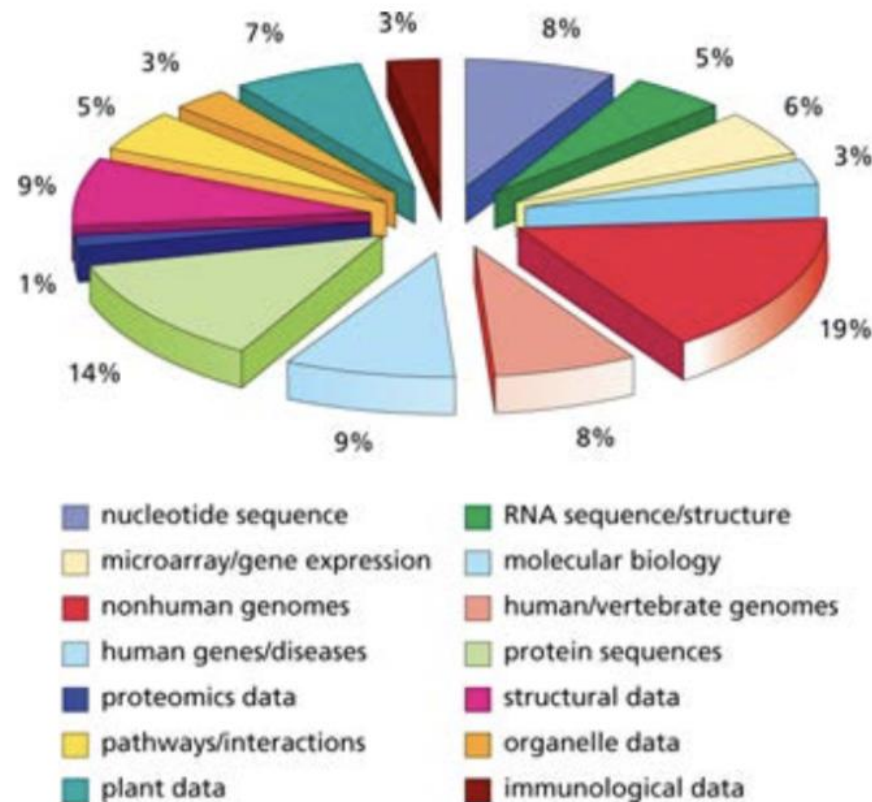- Slides in this section from: Dr. Robin Dowell

# Some statistics

- More than 1000 different bio-databases
- Generally accessible through the web

(useful link: www.expasy.ch/alinks.html)

- Variable size: <100Kb to >10Gb
  - DNA: > 10 Gb
  - Protein: 1 Gb
  - 3D structure: 5 Gb
  - Other: smaller
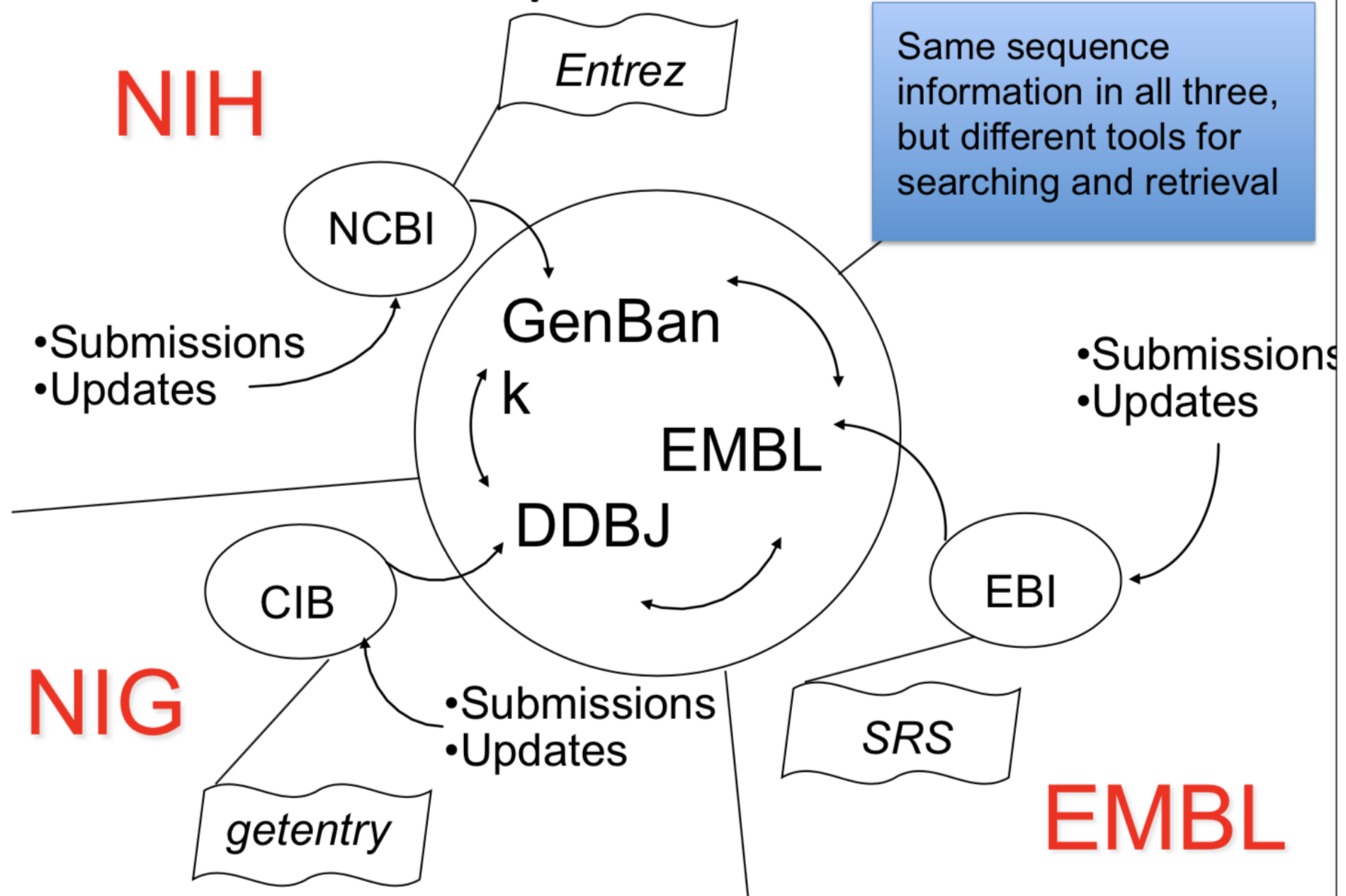
- Update frequency: daily to annually

# NAR Database Issue

- Online collection of biological databases:

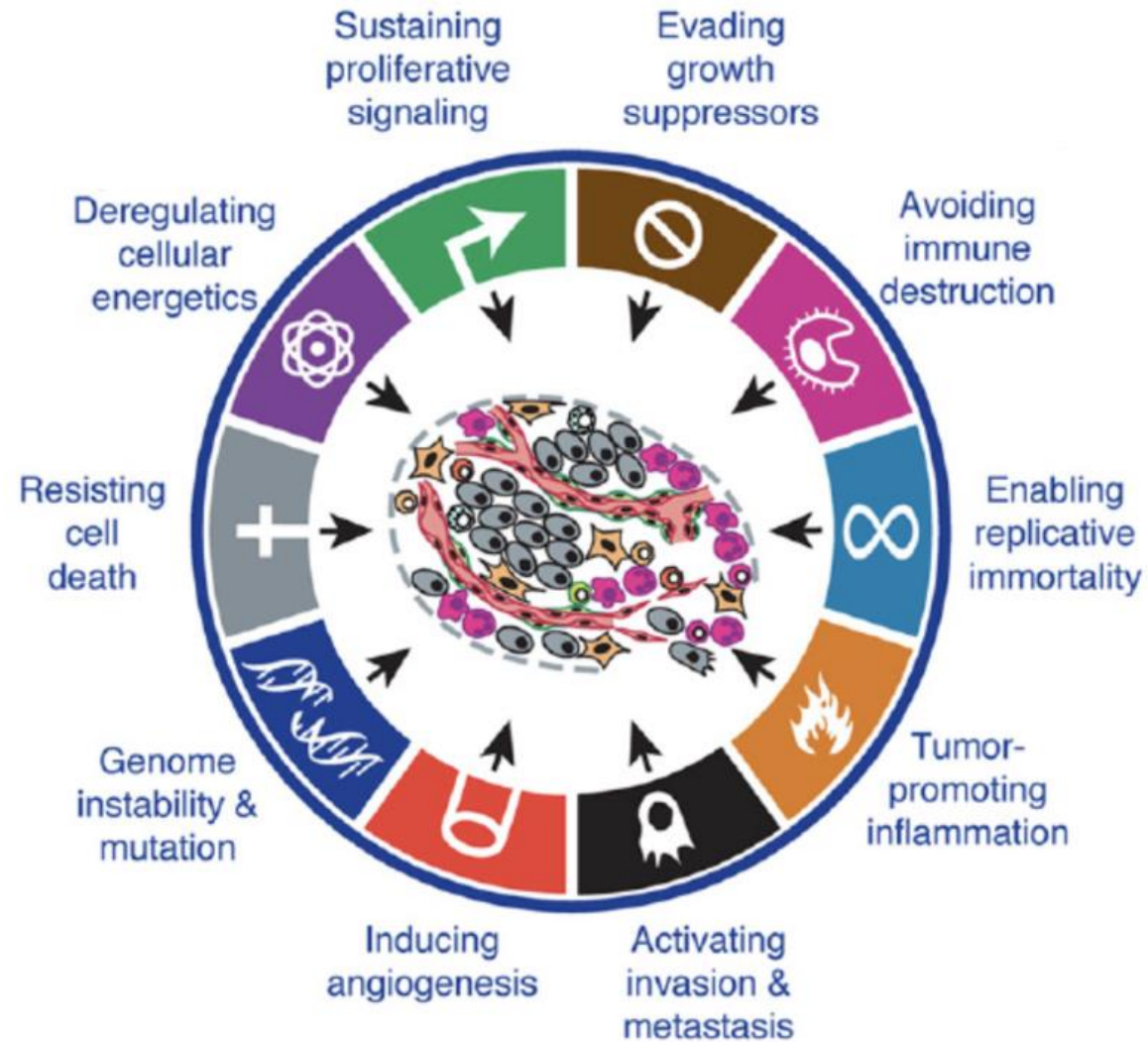http://www.oxfordjournals.org/nar/database/c/

# Public Sequence Databases
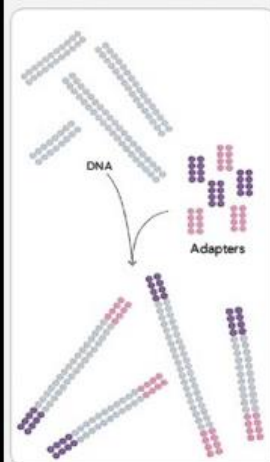
# Cancer Genomics

- Humans are adept at survival
  - Errors (mutations) in DNA are repaired (usually)
  - Redundant systems take over when errors occur
  - Cancerous cells are usually killed before becoming damaging
- Cancer is a complex disease
  - Many mutations must accumulate to cause cancer
  - There is no single way to treat cancer


- Slides in this section from Dr. James Costello
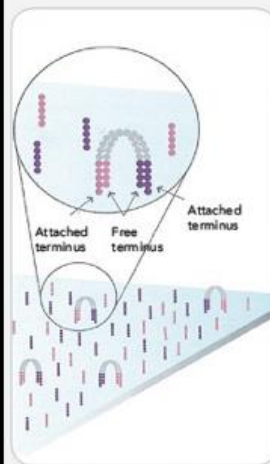
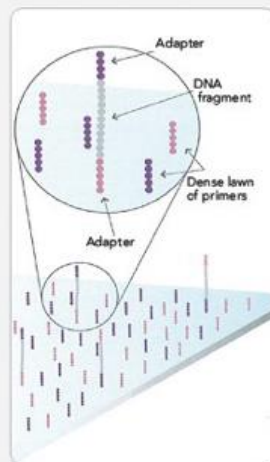# Hallmarks of Cancer

# Illumina Sequencing

# Mutation Calling

# Mutation Calling Pipeline



GATK Best Practices

# Synthetic Lethality

| Gene *A* | Gene *B* | |
|----------|----------|--------|
| *A* | *B* | Viable |
| *A* | *b* | Viable |
| *a* | *B* | Viable |
| *a* | *b* | Lethal |

Kaelin (2005) Nature Reviews Cancer. 5:689-98.

# Gene Expression Analysis

- Studying the transcriptome
  - Parts of the DNA that act like blueprints
  - Why cells with the same DNA can be entirely different

- Slides in this section from Dr. Aik Choon Tan

# The Central Dogma of Molecular Biology



(From Wikipedia)

# cDNA microarray schema

# Gene Expression Profile

$m$ samples

| Geneid | Condition 1 | Condition 2 | ... | Condition $m$ |
|--------|-------------|-------------|-----|---------------|
| Gene1 | 103.02 | 58.79 | ... | 101.54 |
| Gene2 | 40.55 | 1246.87 | ... | 1432.12 |
| ... | ... | ... | ... | ... |
| Gene $n$ | 78.13 | 66.25 | ... | 823.09 |

$n$ genes

# Gene expression data analysis



(Ramaswamy and Golub 2002)

# Gene Expression Clustering



(Cited by 8782)

# Next Generation Sequencing (NGS)

- High throughput sequencers that produce many short sequences

- Lower accuracy (than single sequence Sanger sequencing)

- Requires processing data before use in analysis pipelines

- Slides in this section from Dr. Katerina Kechris

Sequencing the genome of a species

Cataloging variation between individuals in a species

Characterizing differences between cells within an individual

Describing the underlying cellular mechanisms

*Shendure & Aiden Nature Biotechnology 30, 1084–1094 (2012)*

# Base-Calling

- Converts the fluorescence signals of four nucleotides for each cycle into sequence data

- Methods differ on correcting for cross-talk, phasing/prephasing, signal decay

- Return sequence read and quality score for each base

TGCTACGAT...

# Quality Score Per Base

- Assess reliability of base call
- For each base, converts error probability ($p$) to integer score (rounded)

$$q = -10\log_{10}(p)$$

| $q$ | $p$ | Probability called base is correct |
|-----|--------|------------------------------------|
| 10  | 0.1    | 0.9    |
| 20  | 0.01   | 0.99   |
| 30  | 0.001  | 0.999  |
| 40  | 0.0001 | 0.9999 |

# Output: FASTQ Format

```
@SEQ_ID
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>>>CCCCCCC65
```

- ASCII text
- 4 lines per sequence
- Line 1 begins with the @ character, a sequence ID, and an optional description
- Line 2 is the sequence for the read
- Line 3 begins with the + character, followed by the same sequence ID, and another optional description
- Line 4 encodes quality values in hexadecimal format for the sequence letters in line 2
  - Must contain the same number of characters as the sequence in line 2
  - Hexadecimal – use single ASCII to represent up to 92 numeric value (save space)

# Mapping & Reconstruction

1. Mapping: Align reads to reference genome (read mapping)

2. Reconstruct transcriptome (using reference or *de novo*)

# Mapping Options

Align and then assemble

vs

Assemble and then align


Align to genome

vs

Transcriptome


Non-spliced vs Spliced

# Challenges for Mapping

- Multiple hits
- Allowing mismatches
  - High error rate
- Paired-end reads
- Span exon-exon junctions
  - Align with gaps
- Very large number of reads
  - Storage & processing power

# Pharmacogenomics & Personalized Medicine

- Use patient data (DNA) to determine treatments on an individual level
- Requires understanding of the disease and drug(s)
- Benefits from tracking if the patient took the prescribed drug (adherance)

- Slides in this section from: Dr. Debashis Ghosh

# What is Personalized Medicine?

- Also known as "precision medicine."
- Provide "the right patient with the right drug at the right dose at the right time."
- Tailor medical treatment to the individual characteristics, needs, and preferences of a patient during all stages of care.
- Identify <u>genetic</u> or <u>protein</u> biomarkers to predict drug response and side effects.

# Pharmacogenetics Terminology

- A genetic variant is a difference in the DNA sequence compared with a reference sequence.
  - <u>Polymorphism</u>:  A genetic variant that is common, often defined as ≥ 1% in the population.
  - <u>Mutation</u>:  A genetic variant that is rare, often defined as <1% in the population.

| Examples of Types of Genetic Variants | Definition |
|---|---|
| Single nucleotide polymorphism (SNP) | Difference in one nucleotide (base pair) |
| Insertion or deletion (indel) | Insertion or deletion of multiple consecutive nucleotides |
| Repeat polymorphism | Variable number of nucleotides that are repeated |
| Copy number variation (CNV) | Abnormal number of copies of one or more DNA regions (e.g., gene duplication or deletion) |

# *Realities* of Drug Therapy

Physician Prescribes a Medication → Anything goes

A → B C D E F X Y Z #%!

```
MD Prescribes Medication
    → Pt Takes Drug
        → + Effect
        → Side Effect
        → No Effect
        → + and Side Effect
    → Pt Doesn't Take Drug
        → No Effect
        → + / - Effect

+ Effect / Side Effect / No Effect / + and Side Effect →
    Drug – Drug Interaction
    Prescription Changes
    Diagnoses Changes

+ / - Effect ----> Diagnoses Changes
```

```
Medication Prescribed → Medication Taken → Biological Effect
```

# Clopidogrel (Plavix)

- Antiplatelet agent used to reduce the risk of atherosclerotic events.

- Metabolized to 2-oxo-clopidogrel and an active thiol metabolite by numerous CYP enzymes.

- CYP2C19 is the key enzyme involved in the formation of the active metabolite.

- The active metabolite binds to and inhibits P2Y12 receptors on platelets, thereby inhibiting platelet activation and aggregation.

www.pharmgkb.org

# Computational Phylogeny

- Determine evolutionary relationships using genetic information
- Similar sequence => similar function

- Slides in this section from: Dr. David Pollock

# Tetrapods



Birds

Crocodiles

Turtles

Mammals

Amphibians

Lizards 'n' things

Cytochrome C Oxidase I Bayesian consensus tree

# Why Phylogenetics?

- Resolve evolutionary history
  - Important for comparative analysis to account for correlations due to relatedness
- Disease origins, paths of infection
  - Influenza, HIV
- Origin of genes, systems, functions

# Standard Empirical Models of Sequence Evolution

# Topology Space

$$\prod_{i=3}^{T}(2i-5)$$

Times branch lengths

Need tree search algorithm

# Heuristic, Optimal, Posterior

- NP hard, so need tricks
- Distance
  - Estimate of amount of change separating two sequences (species)
  - Calculate analytically (limited), or ML
  - Requires a reversible model of evolution
- Parsimony
  - Minimal number of changes
  - Poorly specified model (but there is one)
  - Easy to calculate
- ML, Bayes
  - Model based, don't toss the data

# Multiple Sequence Alignment

- Compare multiple sequences to identify similarities
- Parts of a sequence that are often preserved are likely important
- Combines fields of Biochemistry, Statistics, & Comp Sci

- Slides in this section from: Dr. Scott Walmsley

# Why MSA?

"Whether the ultimate aim is a **_phylogenetic_** analysis of several orthologues, the identification of a **_pattern_** for particular feature or motif, or the basis for **_structural modelling_**, multiple sequence alignments allow the researcher to gather more biological information than a single sequence can offer"

"The importance of a residue for maintaining the structure and function of a protein can usually be inferred from how conserved it appears in a multiple sequence alignment of that protein and its homologues"

Valdar WS. Scoring residue conservation. Proteins. 2002 Aug 1;48(2):227-41. Review

# Pre-requisite knowledge

*Computational / Math / Statistics & Biochemistry*

- Alphabets

DNA (n= 4)

RNA (n = 4)

Amino Acids (n = 20)

**CODON (n=64)**

Table 3 - Codon usage of the *Arapaima gigas* mtDNA.

| Amino acid (anticodon) | Codon group | Usage of codon ending in A | C | G | T | Total | % |
|---|---|---|---|---|---|---|---|
| Ala (UGC) | GCN | 94 | 108 | 4 | 88 | 294 | 7.53 |
| Cys (GCA) | TGY | 0 | 19 | 0 | 11 | 30 | 0.77 |
| Asp (GUC) | GAY | 0 | 37 | 0 | 39 | 76 | 1.95 |
| Glu (UUC) | GAR | 90 | 0 | 6 | 0 | 96 | 2.46 |
| Phe (GAA) | TTY | 0 | 120 | 0 | 127 | 247 | 6.33 |
| Gly (UCC) | GGN | 92 | 63 | 35 | 46 | 236 | 6.05 |
| His (GUG) | CAY | 0 | 66 | 0 | 45 | 111 | 2.84 |
| Ile (GAU) | ATY | 0 | 110 | 0 | 201 | 311 | 7.97 |
| Lys (UUU) | AAR | 82 | 0 | 5 | 0 | 87 | 2.23 |
| Leu (UAG) | CTN+TTR | 367 | 116 | 44 | 107 | 634 | 16.24 |
| Met (CAU) | ATR | 146 | 0 | 40 | 0 | 186 | 4.77 |
| Asn (GUU) | AAY | 0 | 70 | 0 | 64 | 134 | 3.43 |
| Pro (UGG) | CCN | 122 | 36 | 7 | 43 | 208 | 5.33 |
| Gln (UUG) | CAR | 94 | 0 | 94 | 0 | 188 | 4.82 |
| Arg (UCG) | CGN | 44 | 12 | 4 | 12 | 72 | 1.84 |
| Ser (UGA) | TCN+AGY | 89 | 99 | 4 | 60 | 252 | 6.46 |
| Thr (UGU) | ACN | 139 | 86 | 8 | 81 | 314 | 8.05 |
| Val (UAC) | GTN | 86 | 35 | 12 | 54 | 187 | 4.79 |
| Trp (UCA) | TGR | 111 | 0 | 9 | 0 | 120 | 3.07 |
| Tyr (GUA) | TAY | 0 | 49 | 0 | 64 | 113 | 2.90 |
| Stop (UUA) | TAR | 5 | 0 | 1 | 0 | 6 | 0.15 |
| Stop (UCA) | TGR | 1 | 0 | 0 | 0 | 1 | 0.03 |
| Total | | 1562 | 1026 | 273 | 1042 | 3903 | 100.00 |

# Pre-requisite knowledge

## *Biochemistry / Molecular Biology*

- Mutation rates drive evolution

- Biophysical mechanisms produce mutation rates:

  DNA / RNA Polymerase

- Insertion /Deletion :  frameshift  → altered CODON

Wikipedia

Primary structure
amino acid sequence

Secondary structure
regular sub-structures

alpha helix

beta sheet

hemoglobin

Tertiary structure
three-dimensional structure

P13 protein

Quaternary structure
complex of protein molecules

Wikipedia

# Anatomy of a MSA

Procter JB, Thompson J, Letunic I, Creevey C, Jossinet F, Barton GJ. Visualization of multiple alignments, phylogenies and gene family evolution. Nat Methods. 2010 Mar;7(3 Suppl):S16-25.

# Protein Structure and Prediction

- Given a sequence of DNA can we compute:
  - The translated protein sequence?
  - The protein's shape?
  - The protein's functionality?

- Slides in this section from: Dr. Michael Strong

# Why do we care about protein structures

**Combining Structure and Genomic Information**
**- Help us understand implications of mutations**

## Drug Resistance

**Tuberculosis**



Rifampin

rpoB drug target

# Most Proteins Spontaneously Fold

Important to Computational Biologists, because this suggests that all information relating to the correct folding of a protein is contained in it's primary amino acid sequence, but ...



©1999 GARLAND PUBLISHING INC.
A member of the Taylor & Francis Group

# Experimental Methods of Structure Determination

## X-ray crystallography
## High resolution structure determination

- Intensities and phases of all reflections are combined in a Fourier transform to provide maps of electron density

Phases determined by using heavy metals or selenomethionine (MAD)



Experimental electron density
Phe or Tyr
A break

Fitted electron density
Y110
F64
F121
F62
Y80

# Experimental Methods of Structure Determination

## NMR – Nuclear Magnetic Resonance
## High resolution structure determination

- Smaller Proteins than X-ray

- Distances between pairs of hydrogen atoms

- Lots of information about dynamics

- Requires soluble, non-aggregating material

- Assignment sometimes difficult



NOESY

NOE cross-peak if they are within 5.0 Å

# Experimental Methods of Structure Determination
## Cryo Electron Microscopy
## Low to medium resolution structure determination

- Medium resolution typically ~10-15Å (up to 3.8Å in some special cases)

- Limited information about dynamics

- Can be used for very large molecules and complexes





GroEL crystal structure at 30 A resolution

GroEL    GroEL-ADP    GroEL-AMPPNP    GroEL-ATP

Rosetta structure prediction

2 phases

1.Low-resolution phase – statistical scoring function and fragment assembly

       A. local structure conformations using info from PDB (3 and 9mer stretches)

       B. multiple fragment substitution simulated annealing – to find best arrangement of the fragments (Monte Carlo Search)

       C. low resolution ensemble of decoy conformations

2. Atomic refinement phase using rotamers and small backbone angle moves (in populated regions of Ramachandran plot)

       A. Refinement

       B. Then structures clustered based on RMSD

       C. Center of the Largest Clusters chosen as representative folds (likely to be correct fold)

# The energy model

- Proposed by Linus Pauling in the 1930s
- Bond angles and lengths are almost always the same
- Energy model broken up into two parts:

  Covalent terms
  - Bond distances
  - Bond angles
  - Dihedral angles

  Non-covalent terms
  - Forces at a distance between all non-bonded atoms

# Microbiome & Metagenomics

- Microbiome: "The ecological community of commensal, symbiotic, and pathogenic microorganisms that share our body space"

- Metagenomics: In essence, using DNA sequences from a microbial sample to identify the species present

- Slides in this section from: Dr. Catherine Lozupone

# What do we want to understand?

- What does a healthy microbiome look like?
  - How diverse is it?
  - What types of bacteria are there?
  - What is their function?
- How variable is the microbiome?
  - Over time within an individual?
  - Across individuals?
  - Functionally?
- What are driving factors of variability?
  - Age, culture, physiological state (pregnancy)
- How do changes affect disease?
  - What properties (taxa, amount of diversity) change with disease?
  - Cause or affect?
  - Functional consequences of dysbiosis
- Host Interactions
  - Evolution/adaptation to the host over time.
  - Immune system

Gut microbiota has simple composition at the phylum level

Legend:
- Verrucomicrobia
- Tenericutes
- Proteobacteria
- Fusobacteria
- Firmicutes
- Bacteroidetes
- Actinobacteria

A  B  C  D

Different phyla: Animals and plants

Data from:  Yatsunenko *et. al.* 2012.  Nature.

# Small Subunit Ribosomal RNA

- Present in all known life forms
- Highly conserved
- Resistant to horizontal transfer events
- Big database!



Figure 1 | **Variable regions of the 16S ribosomal RNA.** Secondary structure of the 16S rRNA of *Escherichia coli*, as generated using the xrna program (see Further information). For our analysis, six R fragments of ~250 nucleotides were designed according to the known V regions (Supplementary information S2 (box)). In red, fragment R1 including regions V1 and V2; in orange, fragment R2 including region V3; in yellow, fragment R3 including region V4; in green, fragment R4 including regions V5 and V6; in blue, fragment R5 including regions V7 and V8; and in purple, fragment R6 including region V9.

J. Gregory Caporaso, et al., Nature Methods 2010.
QIIME allows analysis of high- throughput community sequencing data

# BioViz

- How should one visualize all of this complex bio data?
  - Sequences/genes
  - Pathways
  - Multiple –omics
- Why visualize?
  - Help the patient understand their treatment
  - Help clinicians make critical decisions
  - Understand the data and transformations as a bioinformatisist
- Information in this section from: Dr. Carsten Göerg

From Barsky, Aaron et. al. (2008)

# Knowledge-based Analysis

- Doing a bio study is costly
  - Patient consent and the IRB (prepare to wait)
  - Chemicals (reagents) used in some experiments are expensive
  - How to avoid running redundant experiments and use existing data?
- Many databases exist for bio data
  - Many are small and specialized
- Many papers are submitted to PubMed each year
  - Difficult to parse all of them, even in a narrow field

- Slides in this section from Dr. Larry Hunter

# Knowledge-based Analysis

- Build a database that crawls the existing works and builds a graph of **knowledge** about a certain keyword/phrase

- Aggregate the many data sources into one

- Link data through knowledge parsed from PubMed articles

- (Similar to W3's semantic web)

# Analysis is the hard part

- "We are close to having a $1,000 genome sequence, but this may be accompanied by a $1,000,000 interpretation."

  - Bruce Korf, president American College of Medical Genetics

- Not only is the cost of sequencing essentially free, but big computers and big storage are cheap, too. What will keep us busy for the next 50 years is understanding the data"

  - Russ Altman, chair of Biomedical Engineering at Stanford

# Hanalyzer PoC

- Goal: Mixed human/computer approach to knowledge-based approach to analyzing genome-scale datasets.

- Uses graphs to align experimental results with knowledge about genes extracted from many databases (and, increasingly, directly from the biomedical literature).

- [Leach, et al., PLoS Comp Bio 2009]
http://hanalyzer.sourceforge.org
Search YouTube, for "Hanalyzer"