

Functional Genomics 2: Gene Annotation

Michael Schatz

March 4, 2019

Lecture 11: Applied Comparative Genomics



Assignment 4: Due Monday March 4

Assignment 4: Read mapping and variant calling

Assignment Date: Monday, Feb. 25, 2018

Due Date: Monday, March 4, 2018 @ 11:59pm

Assignment Overview

In this assignment, you will consider the algorithms and statistics to align reads to a reference genome to call SNPs and short indels. You will also perform an experiment to empirically determine the "mappability" of a genomic region. Finally, you will investigate some empirical behavior of the binomial test for heterozygous variant calling. As a reminder, any questions about the assignment should be posted to [Piazza](#). Don't forget to read the **Resources** section at the bottom of the page!

Question 1. Dynamic Programming [10 pts + 5pts]

- 1a. Compute the edit distance of (a portion of) the human hemoglobin alpha and beta subunits, showing the dynamic programming matrix and the aligned sequences. Assume a fixed unit cost to substitute one amino acid for another and a unit cost for an insertion or deletion. You are allowed to use the language of your choice, including spreadsheets (Excel, Google sheets, etc)

```
Alpha:  EALERMFLSFPTTKTYFPHFDLSHGSAQVK
Beta:   EALGRLLVYYPWTQRFFESFGDLSTPDAVMGNPKVK
```

- 1b. 5pt BONUS: Notice that the edit distance of GATTTACA and GATACA is 2, but there are multiple possible optimal alignments

```
GATTTACA      GATTTACA      GATTTACA
GAT--ACA      GA-T-ACA      GA--TACA
```

Print 5 optimal alignments between the alpha and beta sequences. If there are more than 5, just print the first 5 you find, although make sure they all have the same minimal edit distance. Hint: Instead of just following the pointers while backtracking, write a recursive depth first search to explore all the possible optimal alignments. The recursion should branch whenever there is a tie in the dynamic programming matrix.

Question 2. Small Variant Analysis [10 pts]

Download chromosome 22 from build 38 of the human genome from here:

Assignment 5: Due Monday March 11

Assignment 5: RNA-seq and differential expression

Assignment Date: Monday, March 4, 2019

Due Date: Monday, March 11, 2019 @ 11:59pm

Assignment Overview

In this assignment, you will analyze gene expression data and learn how to make several kinds of plots in the environment of your choice. (We suggest Python or R.) Make sure to show your work/code in your writeup! As before, any questions about the assignment should be posted to [Piazza](#).

Question 1. Gene Annotation Preliminaries [10 pts]

Download the annotation of build 38 of the human genome from here: ftp://ftp.ensembl.org/pub/release-87/gtf/homo_sapiens/Homo_sapiens.GRCh38.87.gtf.gz

- Question 1a. How many annotated protein coding genes are on each autosome of the human genome? [Hint: Protein coding genes will have "gene" in the 3rd column, and contain the following text: gene_biotype "protein_coding"]
- Question 1b. What is the maximum, minimum, mean, and standard deviation of the span of protein coding genes? [Hint: use the genes identified in 1a]
- Question 1c. What is the maximum, minimum, mean, and standard deviation in the number of exons for protein coding genes? [Hint: you should separately consider each isoform for each protein coding gene]

Question 2. Time Series [10 pts]

[This file](#) contains pre-normalized expression values for 100 genes over 10 time points. Most genes have a stable background expression level, but some special genes show increased expression over the timecourse and some show decreased expression.

- a. Cluster the genes using an algorithm of your choice. Which genes show increasing expression and which genes show decreasing expression, and how did you determine this? What is the background expression level (numerical value) and how did you determine this? [Hint: K-means and hierarchical clustering are common clustering algorithms you could try.]
- b. Calculate the first two principal components of the expression matrix. Show the plot and color the points based on their cluster from part (a). Does the PC1 axis, PC2 axis, neither, or both correspond to the clustering?
- c. Create a heatmap of the expression matrix. Order the genes by cluster, but keep the time points in numerical order.

Question 3. Sampling Simulation [10 pts]

A typical human cell has ~250,000 transcripts, and a typical bulk RNA-seq experiment may involve millions of cells. Consequently in an RNAseq experiment you may start with trillions of RNA molecules, although your sequencer will only give a few million to billions of reads. Therefore your RNAseq experiment will be a small sampling of the full composition. We hope the sequences will be a representative sample of the total population, but if your sample is very unlucky or biased it may not represent the true distribution. We will explore this concept by sampling a small subset of transcripts (1000 to 5000) out of a much larger set (1M) so that you can evaluate this bias.

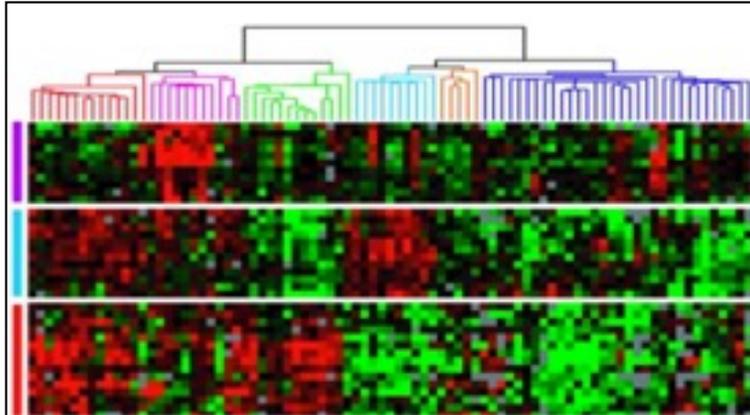
In [data1.txt](#) with 1,000,000 lines we provide an abstraction of RNA-seq data where normalization has been performed and the number of times a gene name occurs corresponds to the number of transcripts sequenced.

- a. Randomly sample 1000 rows. Do this simulation 10 times and record the relative abundance of each of the 15 genes. Plot the mean vs. variance.
- b. Do the same sampling experiment but sample 5000 rows each time. Again plot the mean vs. variance.
- c. Is the variance greater in (a) or (b)?, and explain why. What is the relationship between abundance and variance?

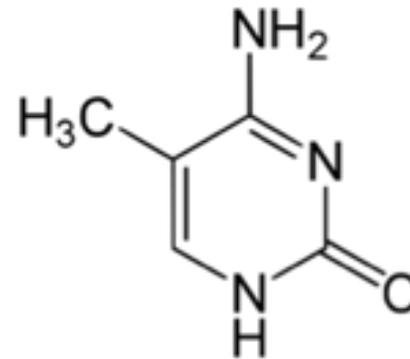


*-seq in 4 short vignettes

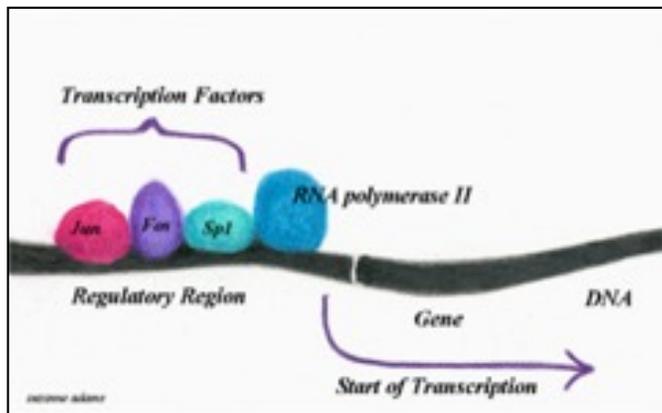
RNA-seq



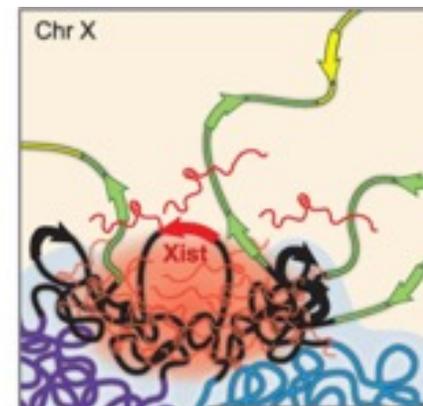
Methyl-seq



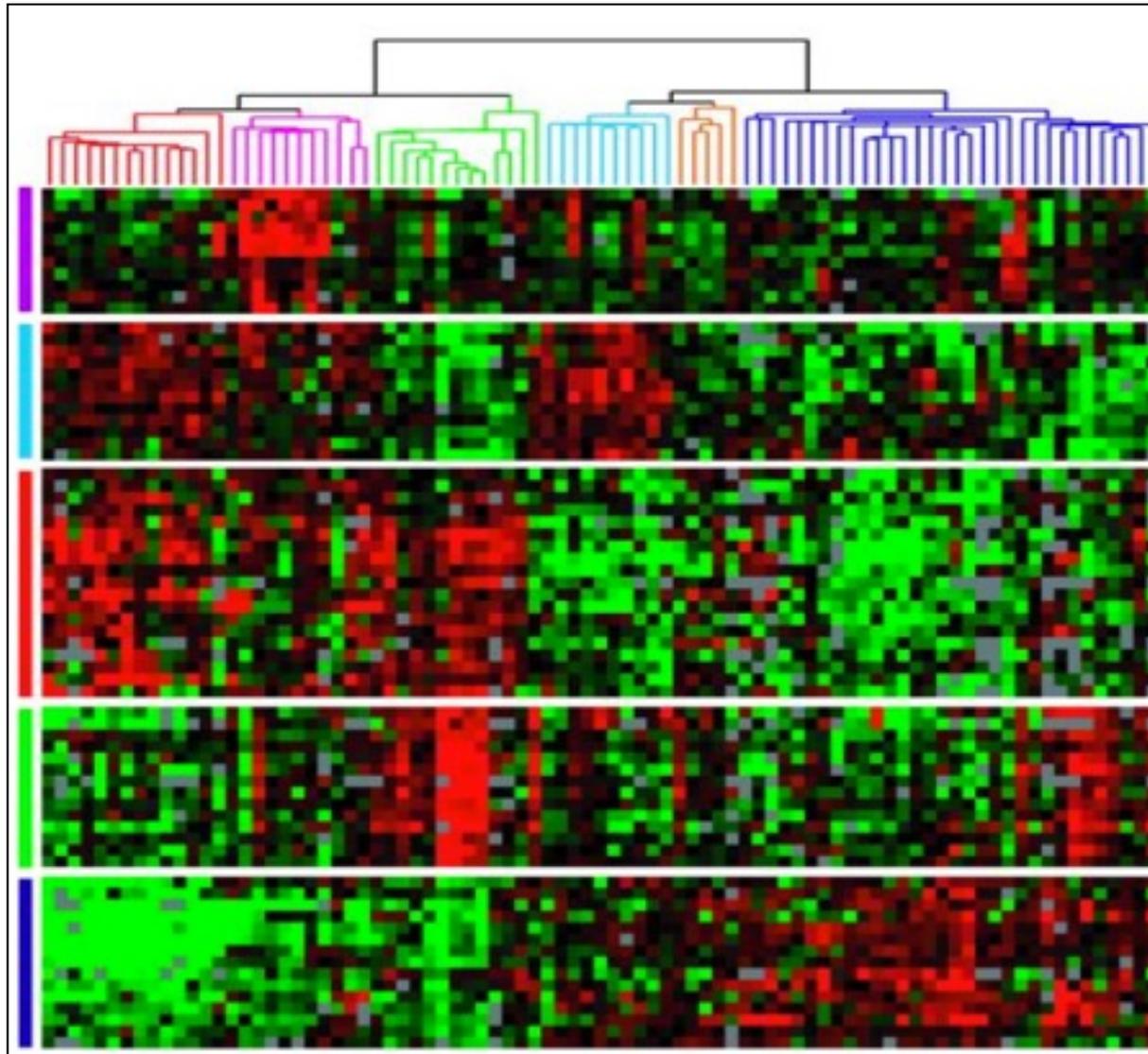
ChIP-seq



Hi-C

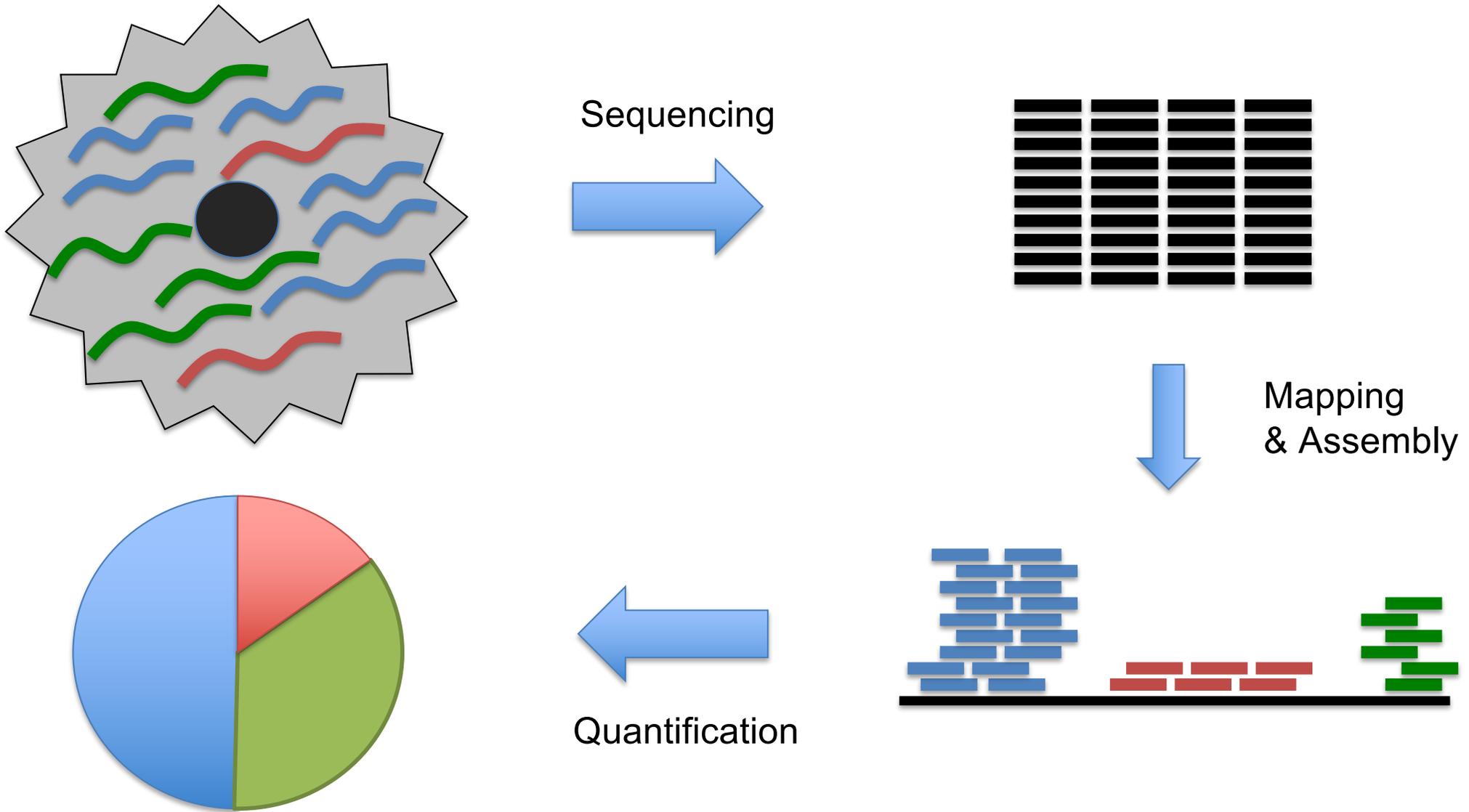


RNA-seq

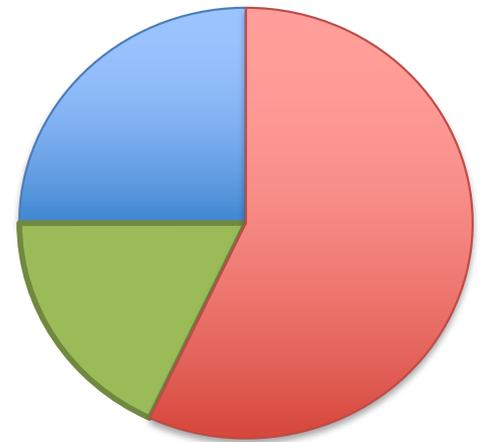
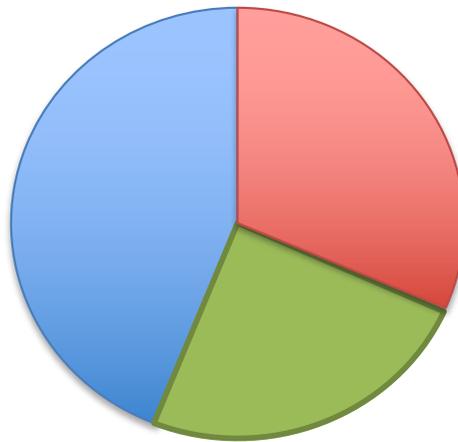
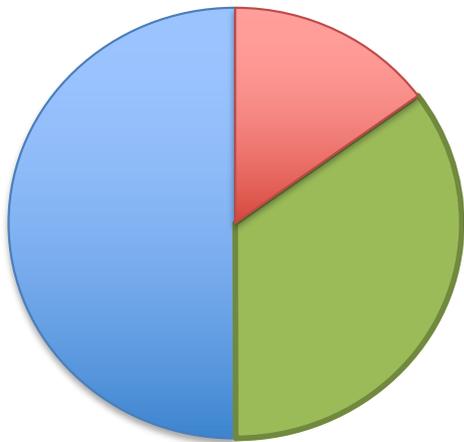
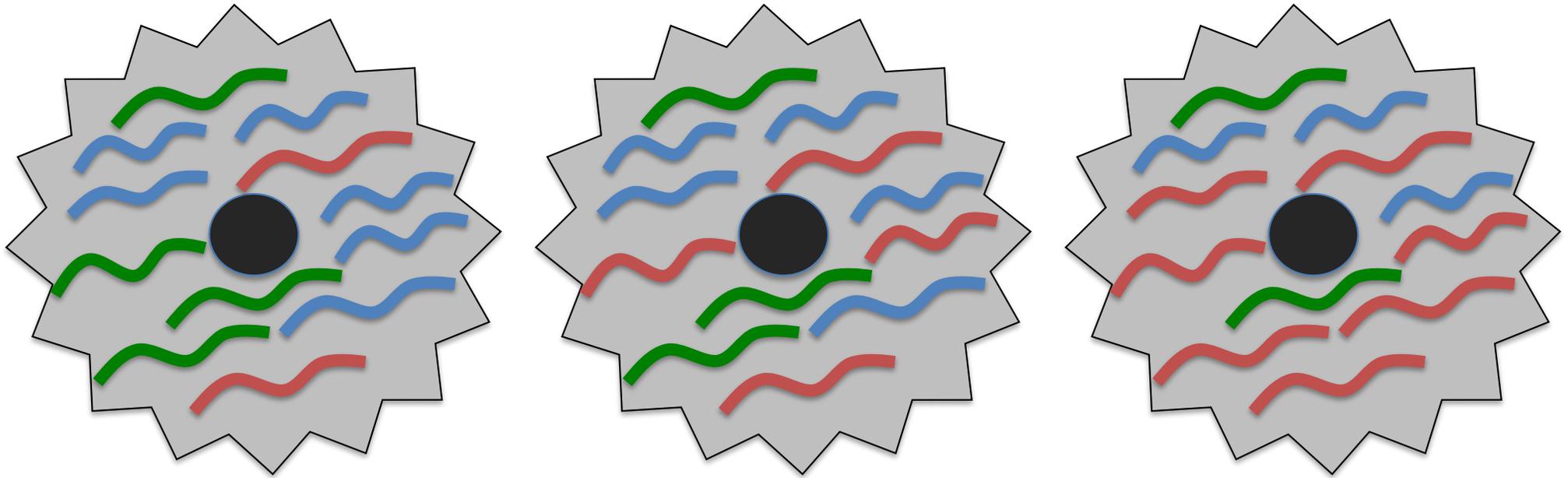


Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications.
Sørli et al (2001) *PNAS*. 98(19):10869-74.

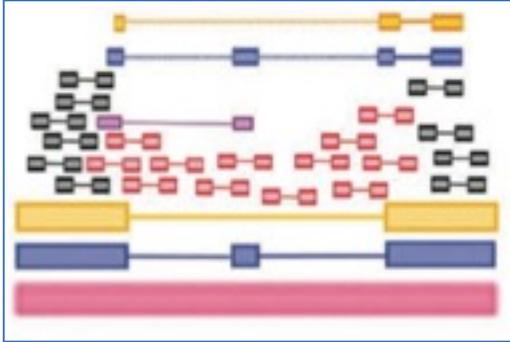
RNA-seq Overview



RNA-seq Overview



RNA-seq Challenges

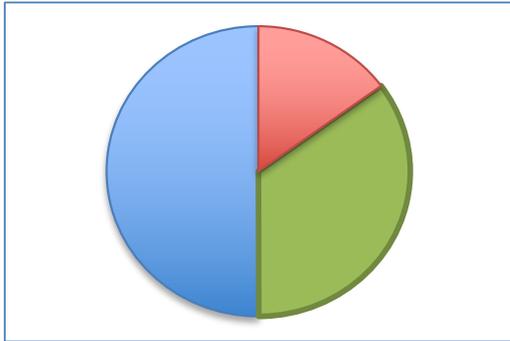


Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

TopHat: discovering spliced junctions with RNA-Seq.

Trapnell et al (2009) *Bioinformatics*. 25:0 1105-1111

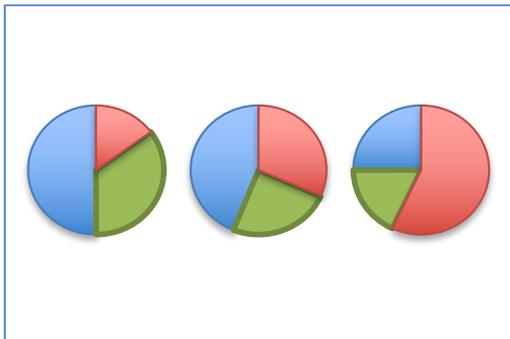


Challenge 2: Read Count \neq Transcript abundance

Solution: Infer underlying abundances (e.g. TPM)

Transcript assembly and quantification by RNA-seq

Trapnell et al (2010) *Nat. Biotech.* 25(5): 511-515



Challenge 3: Transcript abundances are stochastic

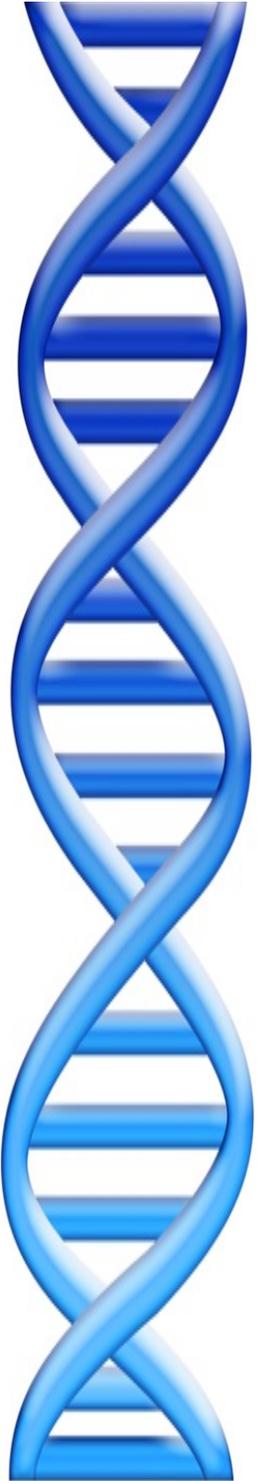
Solution: Replicates, replicates, and more replicates

RNA-seq differential expression studies: more sequence or more replication?

Liu et al (2013) *Bioinformatics*. doi:10.1093/bioinformatics/btt688

Outline

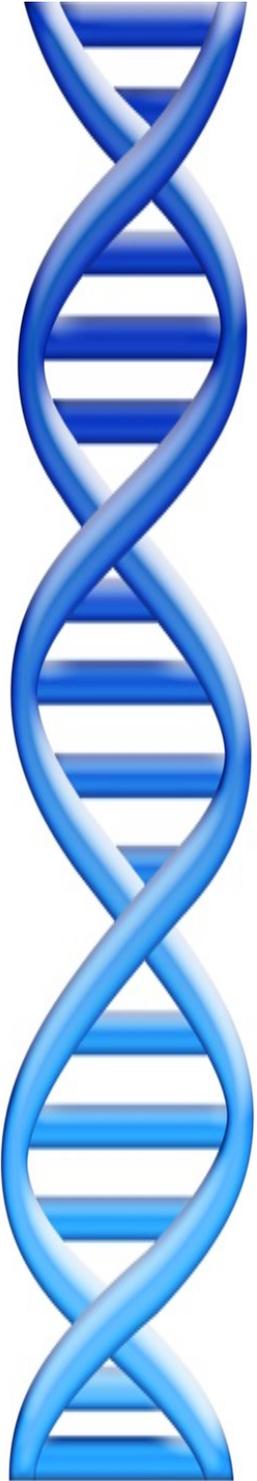
1. Experimental: RNAseq
2. Homology: Alignment to other genomes
3. Prediction: “Gene Finding”



Outline

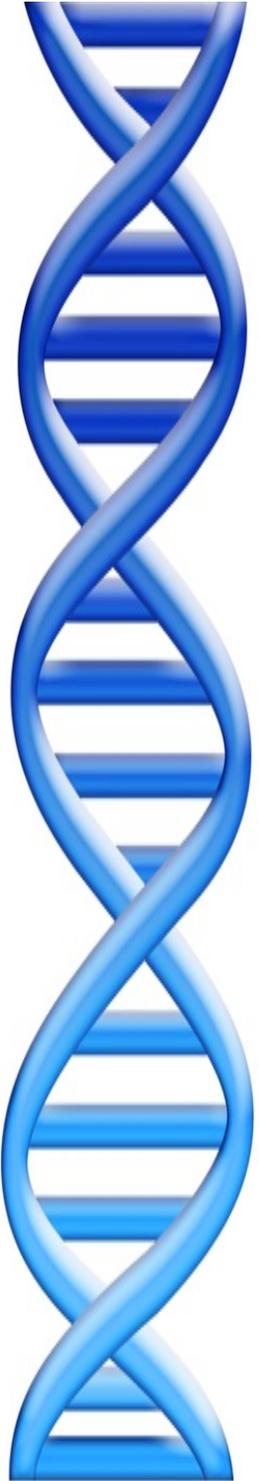
I. Experimental: RNAseq

- 😊 Direct evidence for expression!
 - Including novel genes within a species
- 😞 Typical tissues only express 25% to 50% of genes
 - Many genes are restricted to very particular cell types, developmental stages, or stress conditions
 - Our knowledge of alternative splicing is very incomplete
- 😞 Can resolve gene structure, but nothing about gene function
 - Co-expression is sometimes a clue, but often incomplete



Outline

1. Experimental: RNAseq
2. Homology: Alignment to other genomes
3. Prediction: “Gene Finding”



Basic Local Alignment Search Tool

- Rapidly compare a sequence Q to a database to find all sequences in the database with an score above some cutoff S .
 - Which protein is most similar to a newly sequenced one?
 - Where does this sequence of DNA originate?
- Speed achieved by using a procedure that typically finds “most” matches with scores $> S$.
 - Tradeoff between sensitivity and specificity/speed
 - Sensitivity – ability to find all related sequences
 - Specificity – ability to reject unrelated sequences

Seed and Extend

FAKDFLAGGVAAAI**SKTAVAPIERVKLLLVQV**HASKQITADKQYKGIIDCVVRI PKEQGV
FLIDLASGGTAAAV**SKTAVAPIERVKLLLVQV**DASKAIAVDKRYKGIMDVLIRVPKEQGV

- Homologous sequences are likely to contain a **short high scoring word pair**, a seed.
 - Smaller seed sizes make the sense more sensitive, but also (much) slower
 - Typically do a fast search for prototypes, but then most sensitive for final result
- BLAST then tries to extend high scoring word pairs to compute **high scoring segment pairs (HSPs)**.
 - Significance of the alignment reported via an e-value

BLAST E-values

E-value = the number of HSPs having alignment score **S** (or higher) expected to occur **by chance**.

→ Smaller E-value, more significant in statistics

→ Bigger E-value, less significant

→ Over 1.0 means expect this totally by chance
(not significant at all!)

The expected number of HSPs with the score at least **S** is :

$$E = K * n * m * e^{-\lambda S}$$

K, λ are constant depending on model

n, m are the length of query and sequence

E-values quickly drop off for better alignment bits scores

Very Similar Sequences

Query: HBA_HUMAN Hemoglobin alpha subunit

Sbjct: HBB_HUMAN Hemoglobin beta subunit

Score = 114 bits (285), Expect = 1e-26

Identities = 61/145 (42%), Positives = 86/145 (59%), Gaps = 8/145 (5%)

Query 2 LSPADKTNVKAAWGKVGHAHAGEYGAELERMFLSFPTTKTYFPHF-----DLSHGSAQV 55
L+P +K+ V A WGKV + E G EAL R+ + +P T+ +F F D G+ +V

Sbjct 3 LTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV 60

Query 56 KGHGKKVADALTNAVAHVDDMPNALSALSDDLHAHKL RVD PVNFKLLSHCLLVTLAAHLPA 115
K HGKKV A ++ +AH+D++ + LS+LH KL VDP NF+LL + L+ LA H

Sbjct 61 KAHGKKVLGAFSDGLAHL DNLKGT FATLSELHCDKLHVDPENFRL LGNVLVCVLAH HFGK 120

Query 116 EFTP AVHASLDKFLASVSTVLTSKY 140

EFTP V A+ K +A V+ L KY

Sbjct 121 EFTPPVQAAYQKVVAGVANALAHKY 145

Quite Similar Sequences

Query: HBA_HUMAN Hemoglobin alpha subunit

Sbjct: MYG_HUMAN Myoglobin

Score = 51.2 bits (121), Expect = 1e-07,

Identities = 38/146 (26%), Positives = 58/146 (39%), Gaps = 6/146 (4%)

```
Query 2  LSPADKTNVKAAWGKVGHAHAGEYGAELERMFLSEPTTKTYFPHF-----DLSHGSAQV 55
      LS  +   V   WGKV A   +G E L R+F   P T   F F       D   S   +
Sbjct 3  LSDGEWQLVLNVWGKVEADIPGHGQEVLRIRLFKGGHPETLEKFDKFKHLKSEDEMKAEDL 62
```

```
Query 56 KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKL RVD PVNFKLLSHCLLVTLAAHLPA 115
      K HG V AL   +           + L+ HA K ++       + +S C++ L + P
Sbjct 63 KKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPG 122
```

```
Query 116 EFTP AVHASLDKFLASVSTVLT SKYR 141
      +F      +++K L      + S Y+
Sbjct 123 DFGADAQGAMNKALELFRKDMASNYK 148
```

Not similar sequences

Query: HBA_HUMAN Hemoglobin alpha subunit

Sbjct: SPAC869.02c [Schizosaccharomyces pombe]

Score = 33.1 bits (74), Expect = 0.24

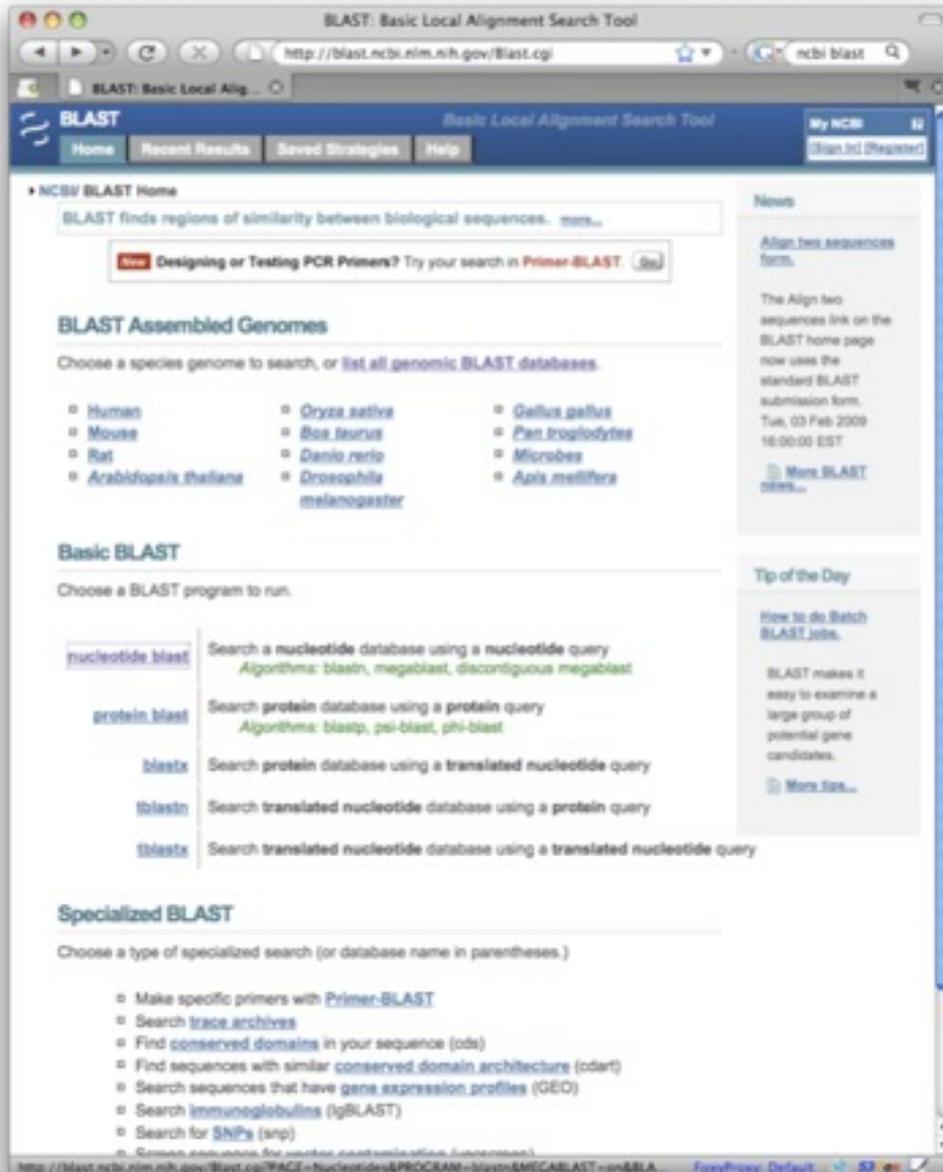
Identities = 27/95 (28%), Positives = 50/95 (52%), Gaps = 10/95 (10%)

```
Query 30  ERMFLSFPTTKTYFPHFDSLHGSAQVKGHGKKVADALTNAVAHVDDMPNALSALS  
      ++M  ++P      P+F+ +H  +      + +A AL N  ++DD+  +LSA  D  
Sbjct 59  QKMLGNYPEV---LPYFNKAHQISL--SQPRILAFALLNYAKNIDDL-TLSAFMDQIVV 112  
  
Query 90  K---LRVDPVNFKLLSHCLLVTLAAHLPAEF-TPA 120  
      K  L++  ++ ++ HCLL T+  LP++  TPA  
Sbjct 113  KHVGLQIKAEHYPIVGHCLLSTMQELLPSDVATPA 147
```

Blast Versions

Program	Database	Query
BLASTN	Nucleotide	Nucleotide
BLASTP	Protein	Protein
BLASTX	Protein	Nucleotide translated into protein
TBLASTN	Nucleotide translated into protein	Protein
TBLASTX	Nucleotide translated into protein	Nucleotide translated into protein

NCBI Blast



- Nucleotide Databases
 - nr:All Genbank
 - refseq: Reference organisms
 - wgs:All reads
- Protein Databases
 - nr:All non-redundant sequences
 - Refseq: Reference proteins

Genetic Code

		1st base								
		U		C		A		G		
2nd base	U	UUU	Phenylalanine	UCU	Serine	UAU	Tyrosine	UGU	Cysteine	U
		UUC	Phenylalanine	UCC	Serine	UAC	Tyrosine	UGC	Cysteine	C
		UUA	Leucine	UCA	Serine	UAA	Stop	UGA	Stop	A
		UUG	Leucine	UCG	Serine	UAG	Stop	UGG	Tryptophan	G
	C	CUU	Leucine	CCU	Proline	CAU	Histidine	CGU	Arginine	U
		CUC	Leucine	CCC	Proline	CAC	Histidine	CGC	Arginine	C
		CUA	Leucine	CCA	Proline	CAA	Glutamine	CGA	Arginine	A
		CUG	Leucine	CCG	Proline	CAG	Glutamine	CGG	Arginine	G
	A	AUU	Isoleucine	ACU	Threonine	AAU	Asparagine	AGU	Serine	U
		AUC	Isoleucine	ACC	Threonine	AAC	Asparagine	AGC	Serine	C
		AUA	Isoleucine	ACA	Threonine	AAA	Lysine	AGA	Arginine	A
		AUG	Methionine (Start)	ACG	Threonine	AAG	Lysine	AGG	Arginine	G
	G	GUU	Valine	GCU	Alanine	GAU	Aspartic Acid	GGU	Glycine	U
		GUC	Valine	GCC	Alanine	GAC	Aspartic Acid	GGC	Glycine	C
		GUA	Valine	GCA	Alanine	GAA	Glutamic Acid	GGA	Glycine	A
		GUG	Valine	GCG	Alanine	GAG	Glutamic Acid	GGG	Glycine	G

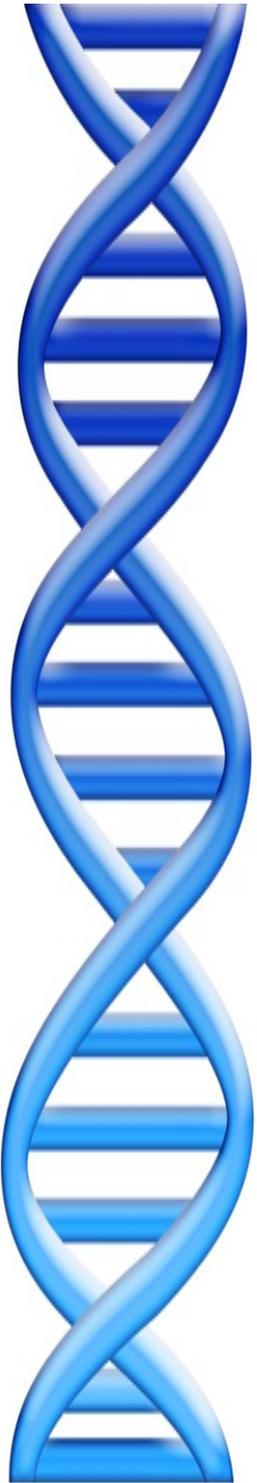
3rd base

Nonpolar, aliphatic Polar, uncharged Aromatic Positively charged Negatively charged

Outline

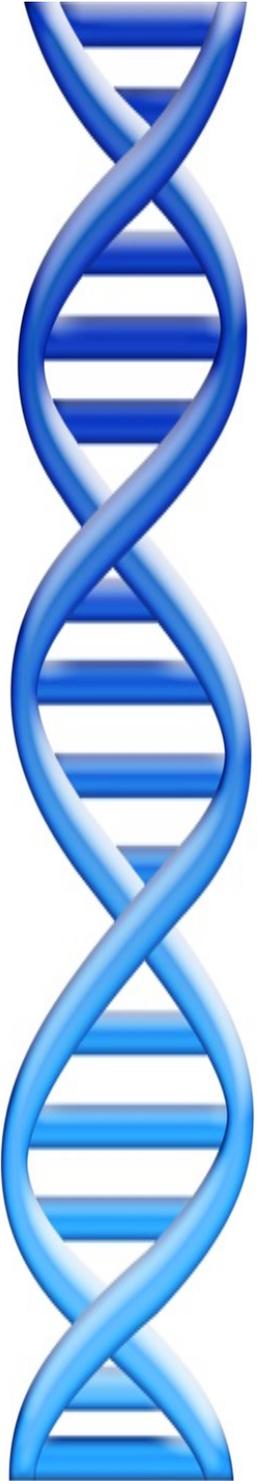
2. Homology: Alignment to other genomes

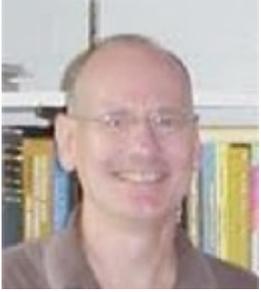
- :-/ Indirect evidence for expression
 - Works well for familiar species, but more limited for unexplored clades
 - Relatively few false positives, but many false negatives
- 😊 Universal across tissues (and species)
 - Proteins often have highly conserved domains, whereas genome/transcript may have many mutations (especially “wobble” base)
- :-/ Transfer gene function across species
 - Reciprocal best blast hit a widely used heuristic
 - Often works, but examples where single base change leads to opposite function



Outline

1. Experimental: RNAseq
2. Homology: Alignment to other genomes
3. Prediction: “Gene Finding”





Bacterial Gene Finding and Glimmer

(also Archaeal and viral gene finding)

Arthur L. Delcher and Steven Salzberg
Center for Bioinformatics and Computational Biology
Johns Hopkins University School of Medicine

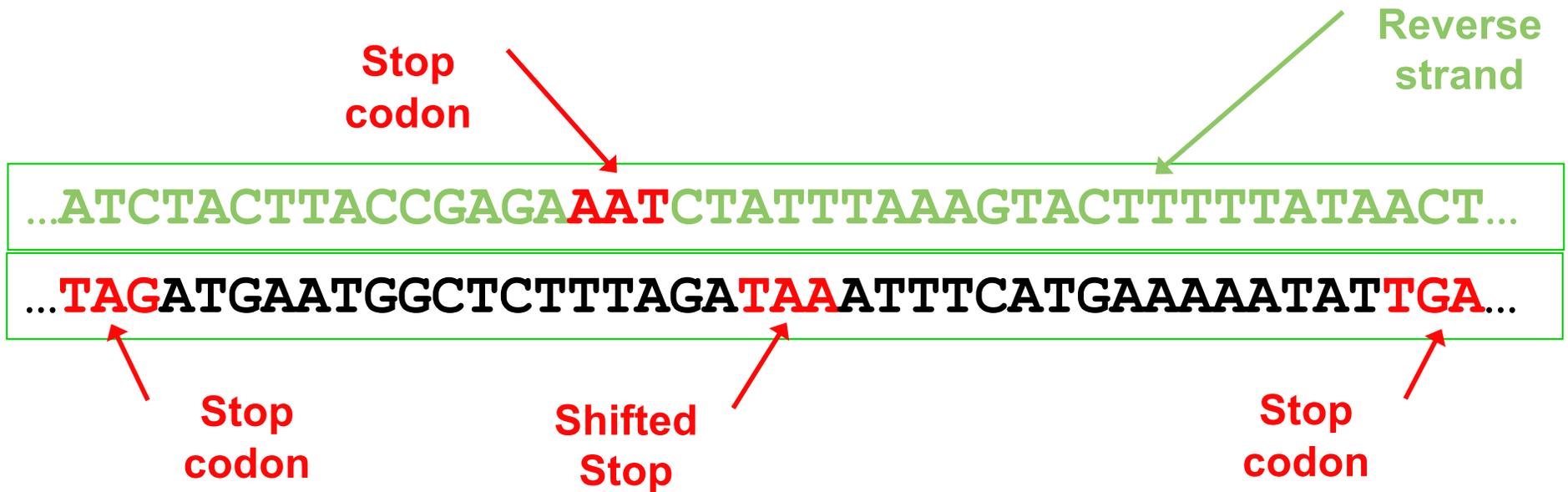
Step One

- Find open reading frames (ORFs).



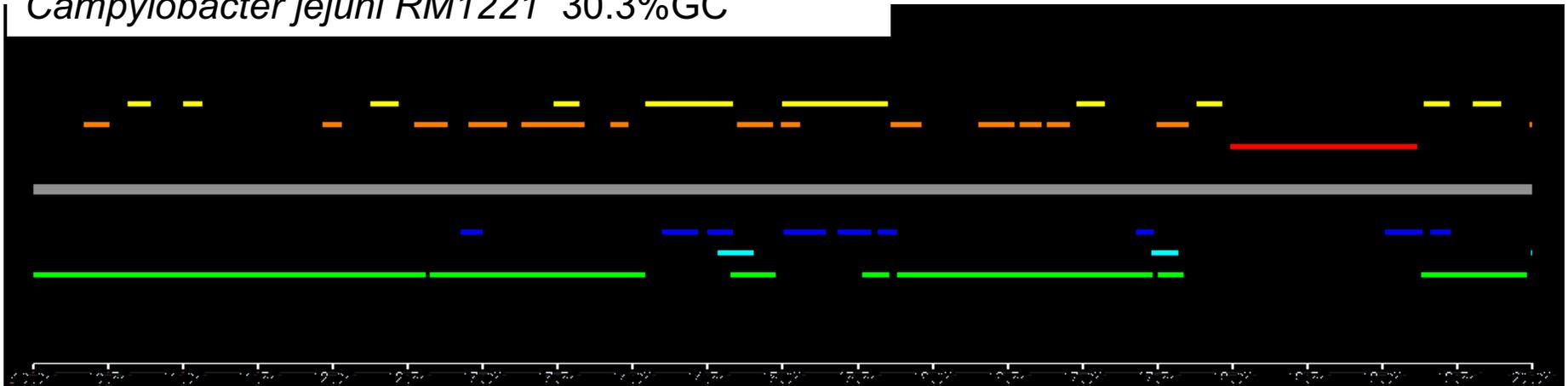
Step One

- Find open reading frames (ORFs).



- But ORFs generally overlap ...

Campylobacter jejuni RM1221 30.3%GC



All ORFs longer than 100bp on both strands shown
- color indicates reading frame

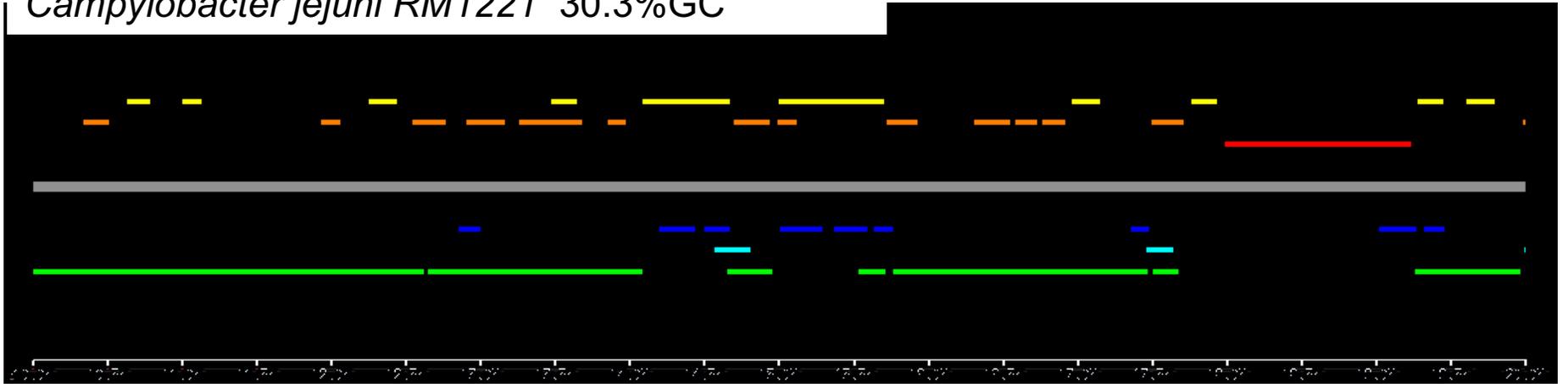
Note the low GC content

- many A+T -> many stop codons (TAA/TAG/TGA)

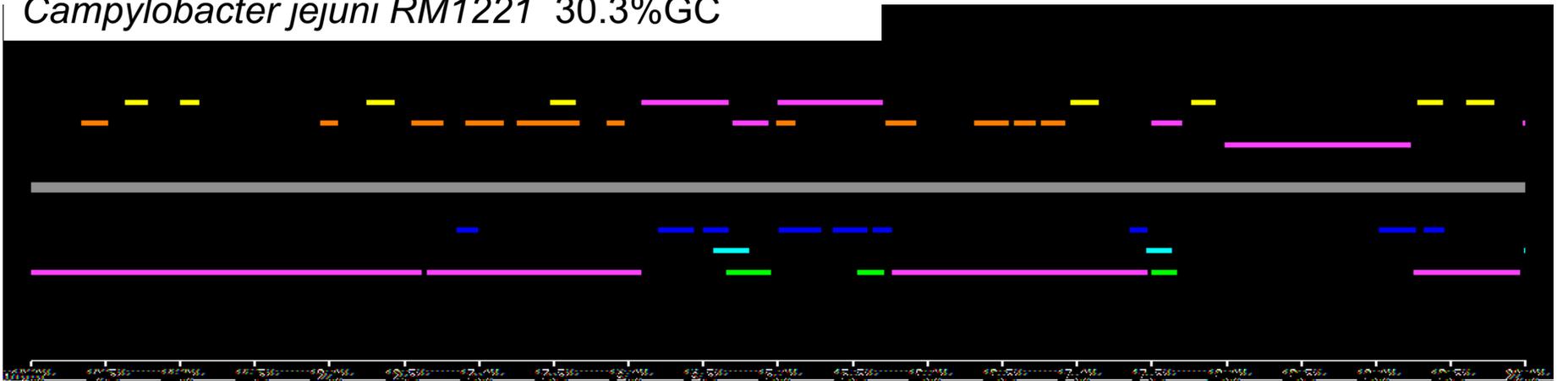
All genes are ORFs but not all ORFs are genes

- Longest ORFs likely to be protein-coding genes

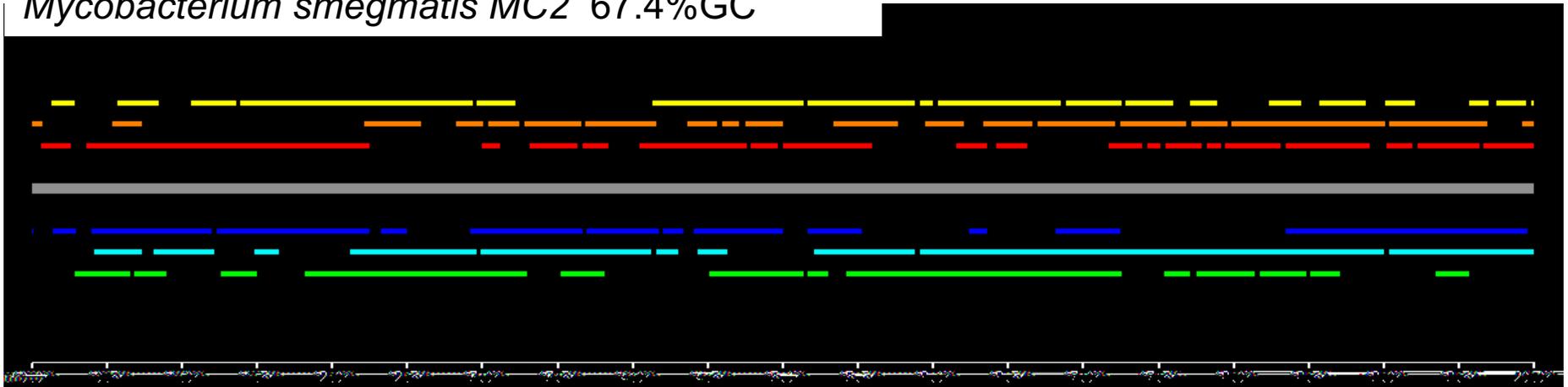
Campylobacter jejuni RM1221 30.3%GC



Campylobacter jejuni RM1221 30.3%GC

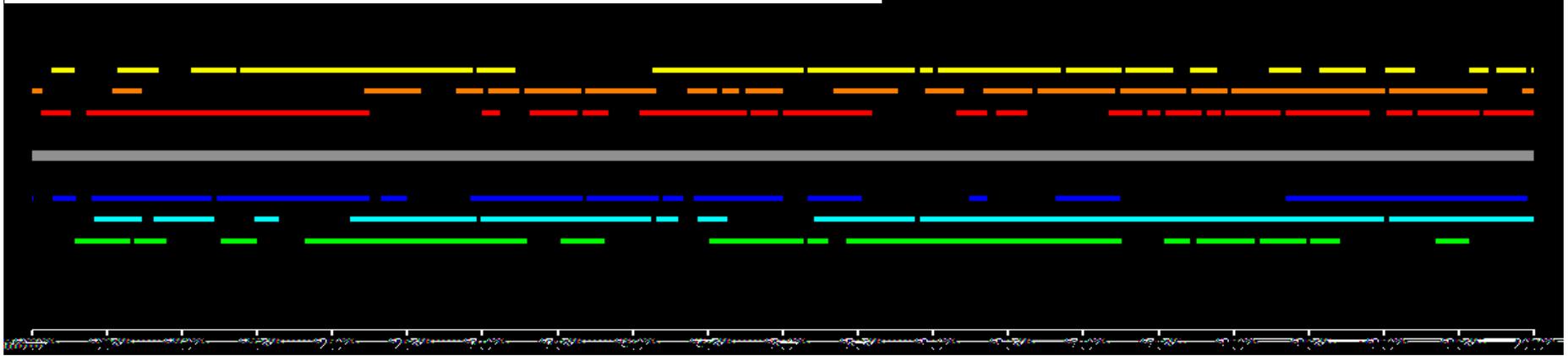


Mycobacterium smegmatis MC2 67.4%GC

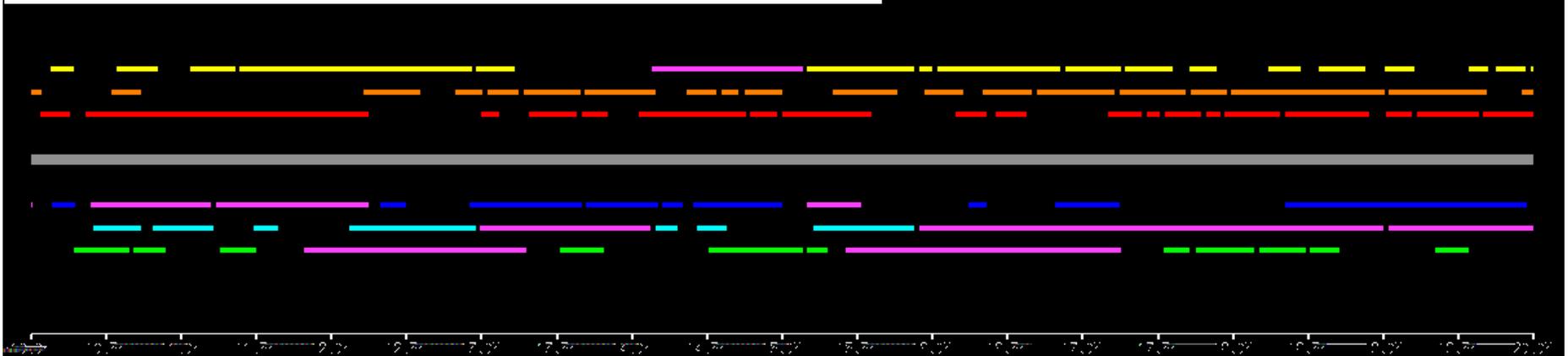


Note what happens in a high-GC genome

Mycobacterium smegmatis MC2 67.4%GC



Mycobacterium smegmatis MC2 67.4%GC



Probabilistic Methods

- Create models that have a probability of generating any given sequence.
 - Evaluate gene/non-genome models against a sequence
- Train the models using examples of the types of sequences to generate.
 - Use RNA sequencing, homology, or “obvious” genes
- The “score” of an orf is the probability of the model generating it.
 - Most basic technique is to count how kmers occur in known genes versus intergenic sequences
 - More sophisticated methods consider variable length contexts, “wobble” bases, other statistical clues

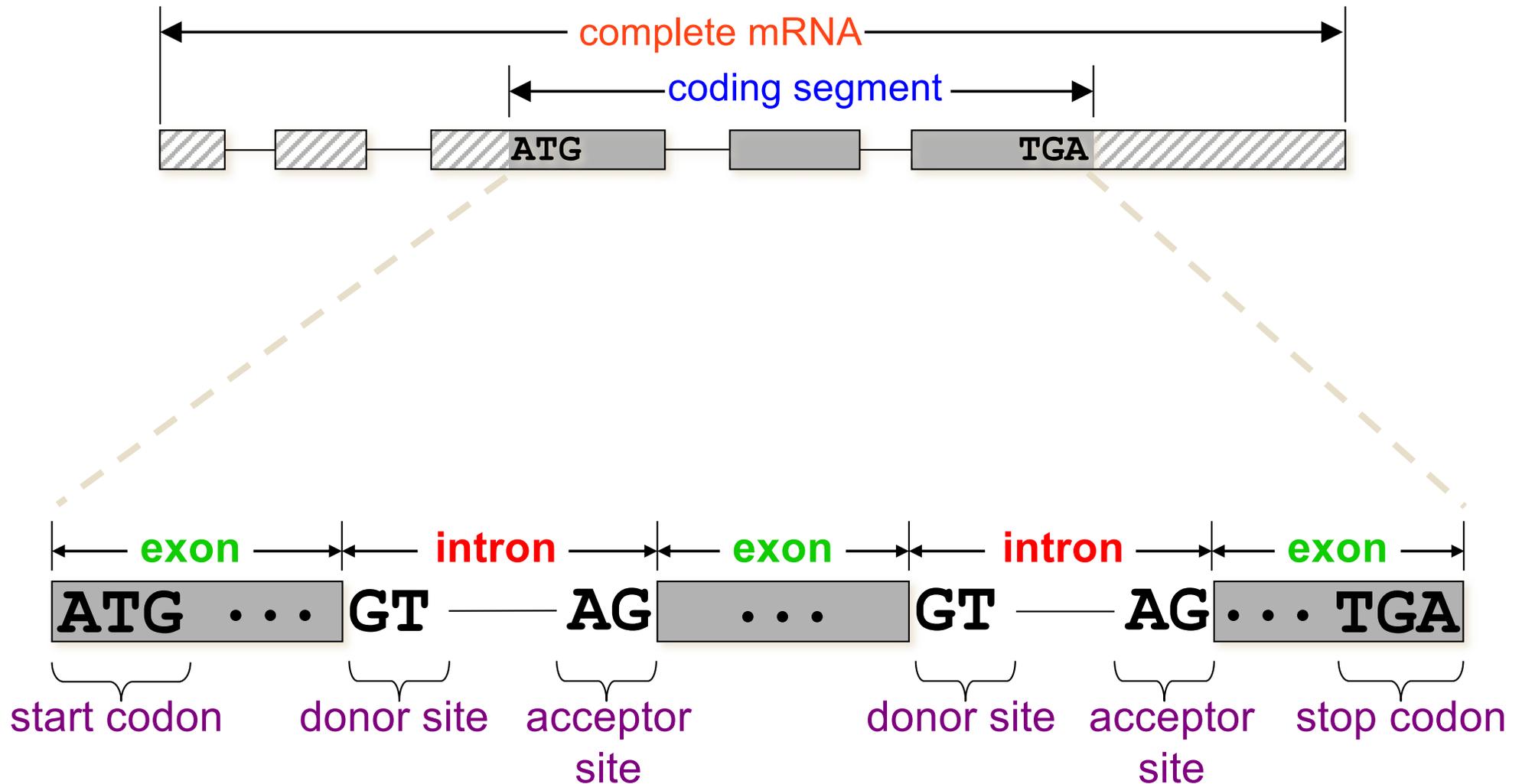


Overview of Eukaryotic Gene Prediction

CBB 231 / COMPSCI 261

W.H. Majoros

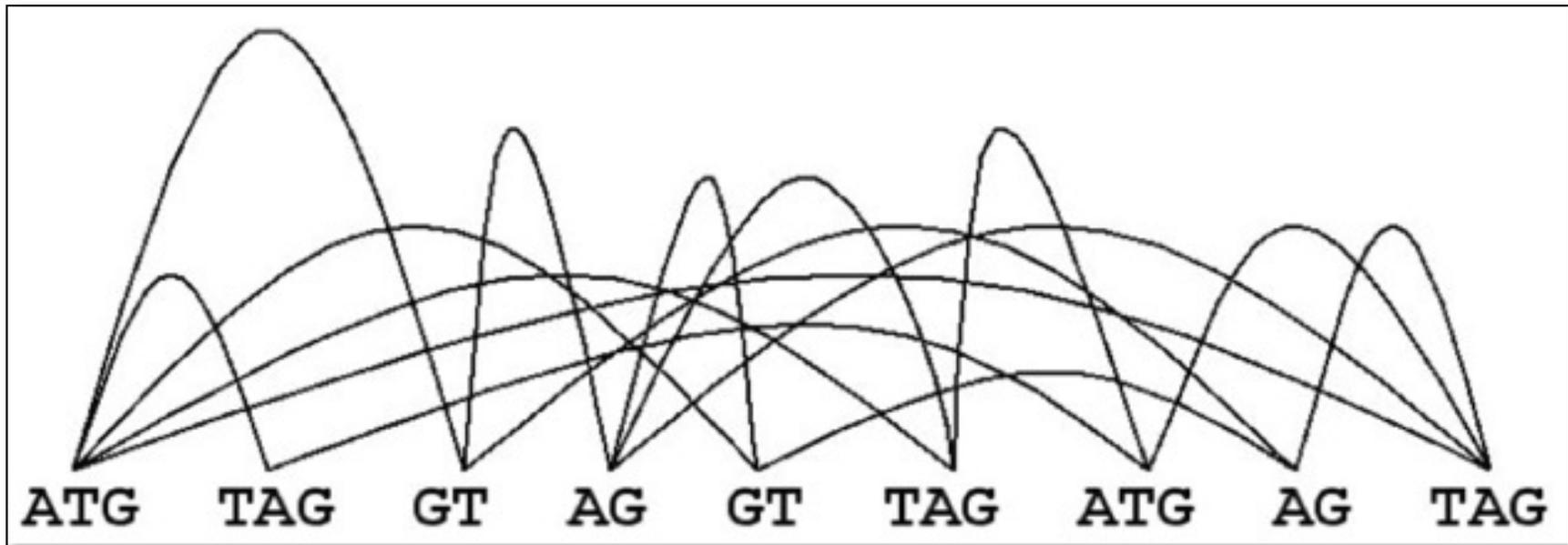
Eukaryotic Gene Syntax



Regions of the gene outside of the CDS are called *UTR*'s (*untranslated regions*), and are mostly ignored by gene finders, though they are important for regulatory functions.

Representing Gene Syntax with ORF Graphs

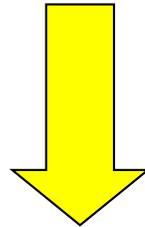
After identifying the most promising (i.e., highest-scoring) signals in an input sequence, we can apply the gene syntax rules to connect these into an *ORF graph*:



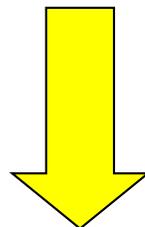
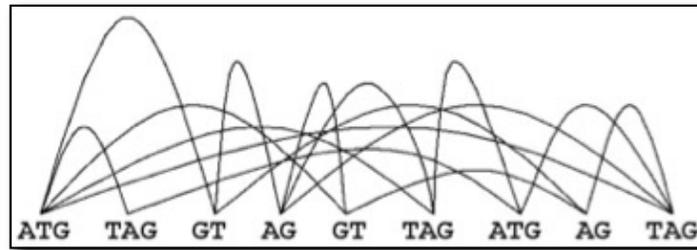
An ORF graph represents all possible *gene parses* (and their scores) for a given set of putative signals. A *path* through the graph represents a single gene parse.

Conceptual Gene-finding Framework

```
TATTCCGATCGATCGATCTCTCTAGCGTCTACG  
CTATCATCGCTCTCTATTATCGCGGATCGTCG  
ATCGCGGAGAGTATGCTACGTCGATCGAATTG
```



identify most promising signals, score signals and content regions between them; induce an ORF graph on the signals



find highest-scoring path through ORF graph; interpret path as a gene parse = gene structure



What is an HMM?

- Dynamic Bayesian Network

- A set of states

- {Fair, Biased} for coin tossing
- {Gene, Not Gene} for Bacterial Gene
- {Intergenic, Exon, Intron} for Eukaryotic Gene

- A set of emission characters

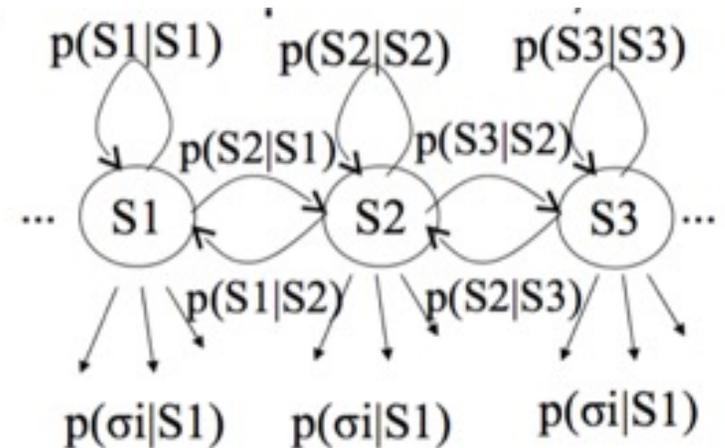
- $E = \{H, T\}$ for coin tossing
- $E = \{1, 2, 3, 4, 5, 6\}$ for dice tossing
- $E = \{A, C, G, T\}$ for DNA

- State-specific emission probabilities

- $P(H | \text{Fair}) = .5, P(T | \text{Fair}) = .5, P(H | \text{Biased}) = .9, P(T | \text{Biased}) = .1$
- $P(A | \text{Gene}) = .9, P(A | \text{Not Gene}) = .1 \dots$

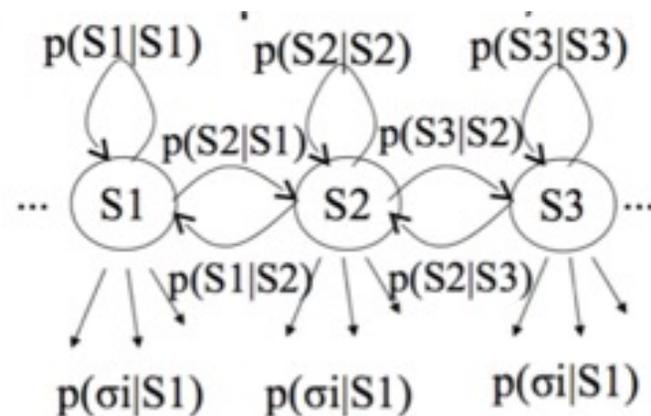
- A probability of taking a transition

- $P(s_i = \text{Fair} | s_{i-1} = \text{Fair}) = .9, P(s_i = \text{Bias} | s_{i-1} = \text{Fair}) = .1$
- $P(s_i = \text{Exon} | s_{i-1} = \text{Intergenic}), \dots$



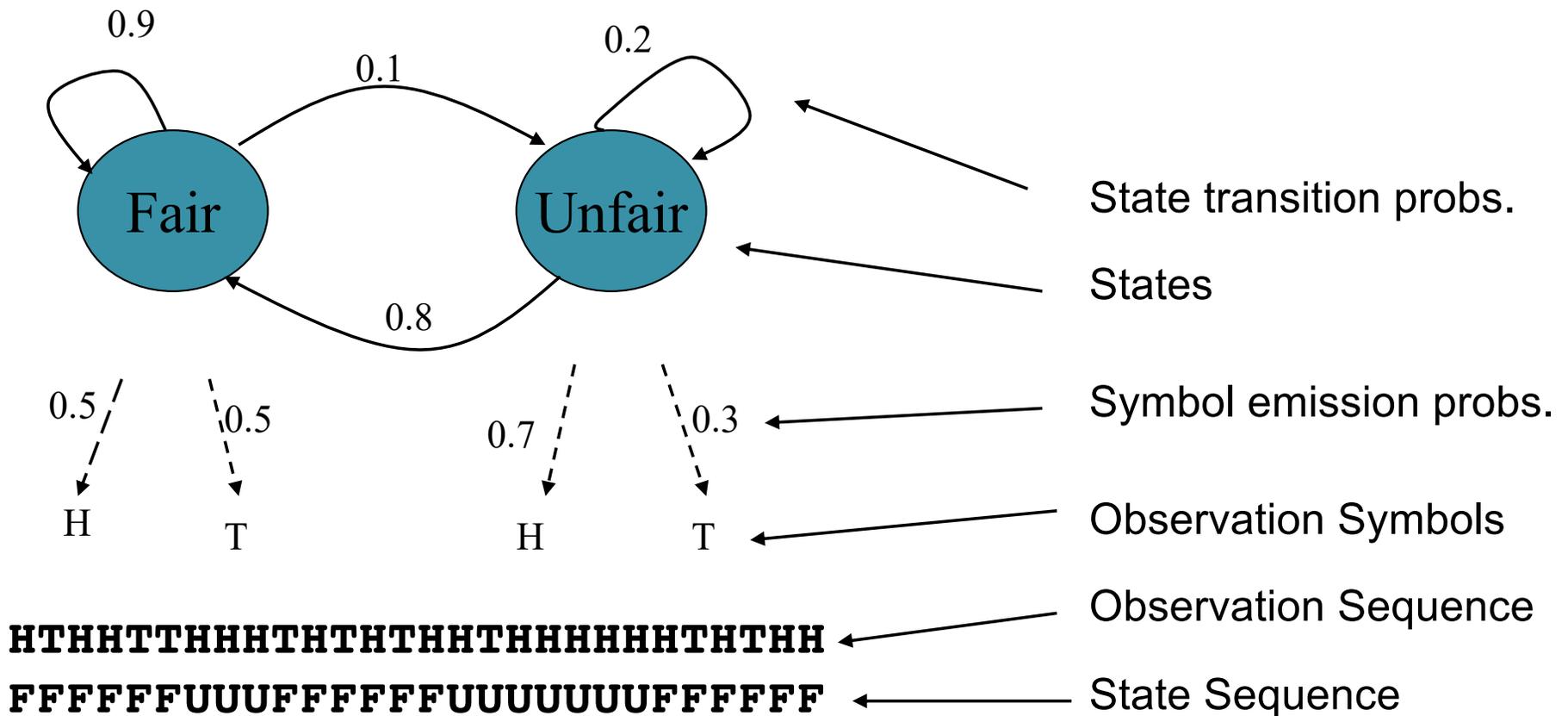
Why Hidden?

- Similar to Markov models used for prokaryotic gene finding, but system may transition between multiple models called states (gene/non-gene, intergenic/exon/intron)
- Observers can see the emitted symbols of an HMM (i.e., nucleotides) but have no ability to know which state the HMM is currently in.
 - But we can *infer* the most likely hidden states of an HMM based on the given sequence of emitted symbols.



AAAGCATGCATTTAACGTGAGCACAAATAGATTACA

HMM Example - Casino Coin



Motivation: Given a sequence of H & Ts, can you tell at what times the casino cheated?