

Lecture 13. Gene Finding & RNAseq

Michael Schatz

March 8, 2021

Applied Comparative Genomics



Assignment 5: Due Wed Mar 17

The screenshot shows a GitHub repository page for "Assignment 5: BWT and RNA-seq". The page has a light gray header with standard browser controls. Below the header, the repository details are shown: "148 Lines (89 sloc) 8.82 KB". On the right side of the header, there are buttons for "Raw", "Blame", "Copy", and "Edit". The main content area has a title "Assignment 5: BWT and RNA-seq" and a subtitle "Assignment Date: Wednesday, Mar. 3, 2021" and "Due Date: Wednesday, Mar. 17, 2021 @ 11:59pm". A section titled "Assignment Overview" contains text about the assignment requirements and a note to show work/code in the writeup. It also mentions Piazza for questions. A question section for BWT Encoding is described with pseudo code.

Assignment 5: BWT and RNA-seq

Assignment Date: Wednesday, Mar. 3, 2021
Due Date: Wednesday, Mar. 17, 2021 @ 11:59pm

Assignment Overview

In this assignment you will write a simple BWT encoder and decoder, and explore a couple of aspects of RNA-seq (with a small introduction to clustering). For this assignment, you will have to generate some visualizations - we recommend R or Python, but use a language you are comfortable with!

Make sure to show your work/code in your writeup!

As a reminder, any questions about the assignment should be posted to [Piazza](#).

Question 1. BWT Encoding [10 pts]

In the language of your choice, implement a BWT encoder and encode the string below. Linear time methods exist for computing the BWT, although for this assignment you can use the simple method based on standard sorting techniques. Your solution does not need to be an optimal algorithm and can use $O(n^2)$ space and $O(n^2 \lg n)$ time.

Here is the recommended pseudo code (make sure to submit your code as well as the encoded string):

```
computeBwt(string s)
    ## add the magic end-of-string character
    s = s + "$"

    ## build up the BWT from the cyclic permutations
    ## note the i-th cyclic permutation is just "s[i..n] + s[0..i]"
    StringList rows = []
    for (i = 0; i < length(s); i++)
        rows.append(cyclic_permutation(s, i))

    ## last use the builtin sort command
```

Project Proposal: Due March 15

A screenshot of a web browser window showing a GitHub project proposal page. The title bar says "github.com". The page content includes:

Project Proposal

Assignment Date: Monday March 8, 2021
Due Date: Monday, March 15 2021 @ 11:59pm

Review the [Project Ideas](#) page

Work solo or form a team for your class project of no more than 3 people.

The proposal should have the following components:

- Name of your team
- List of team members and email addresses
- Short title for your proposal
- 1 paragraph description of what you hope to do and how you will do it
- References to 2 to 3 relevant papers
- References/URLs to datasets that you will be studying (Note you can also use simulated data)
- Please add a note if you need me to sponsor you for a MARCC account (high RAM, GPUs, many cores, etc)

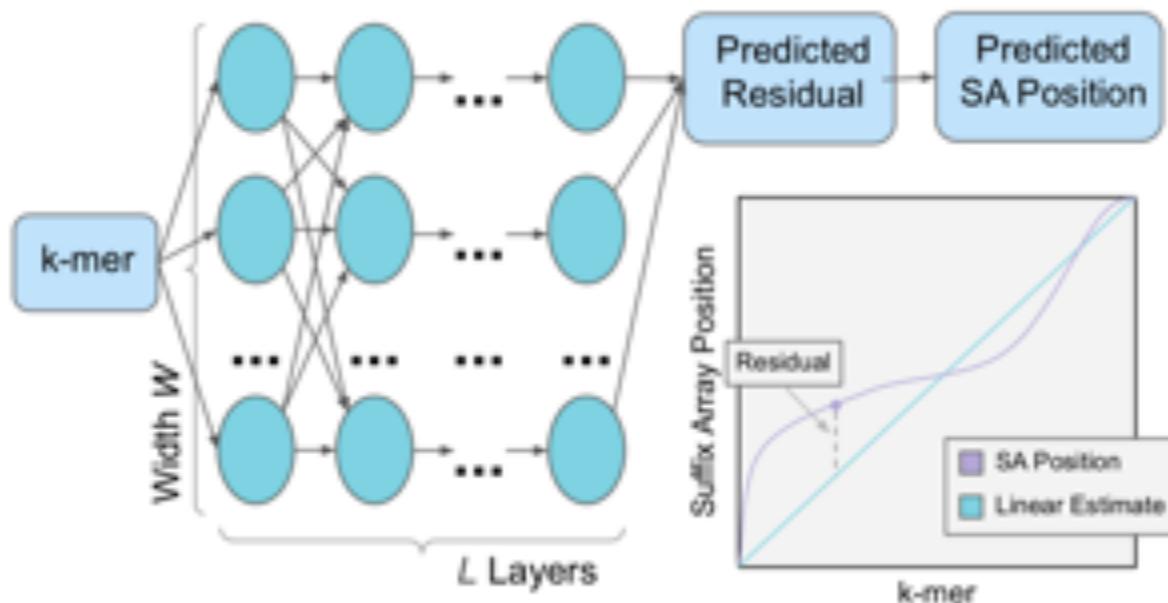
Submit the proposal as a 1 to 2 page PDF on GradeScope (each team member should submit the same PDF). After submitting your proposal, we can schedule a time to discuss your proposal, especially to ensure you have access to the data that you need. The sooner that you submit your proposal, the sooner we can schedule the meeting. No late days can be used for the project.

Later, you will present your project in class during the last week of class. You will also submit a written report (5-7 pages) of your project, formatting as a Bioinformatics article (Intro, Methods, Results, Discussion, References). Word and LaTeX templates are available at https://academic.oup.com/bioinformatics/pages/submission_online

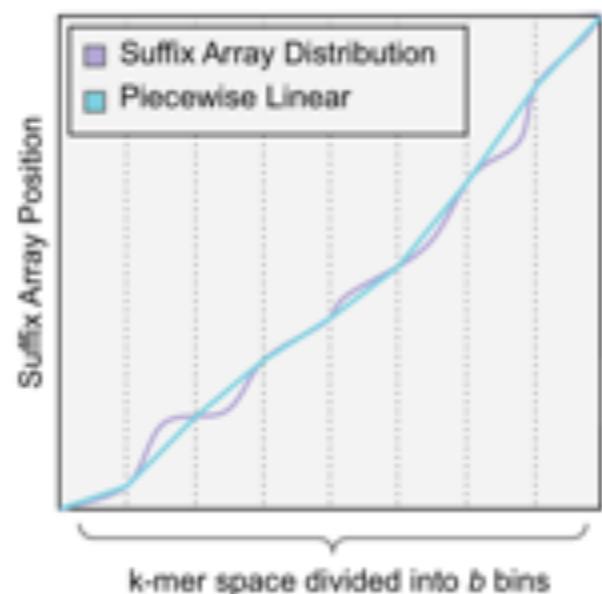
Please use Piazza to coordinate proposal plans!

Modeling the function

a) ANN Architecture



b) Piecewise Linear Architecture



Goal: Genome Annotations

aatgcatgcggctatgctaattgcattgcggctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgcaaggcttggatccgatgactatgctaagctggatccgatgacaatgcattgcggctatgctaattgaatggtcttggatttaccttggaaatgctaagctggatccgatgacaatgcattgcggctatgctaattgaaatggtcttggatttaccttggaaatatgctaattgcattgcggctatgctaagctggatccgatgacaatgcattgcg gctatgctaattgcattgcggctatgcaaggctggatccgatgactatgctaagctggatccgatgacaatgcattgcg gctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgctaattgcattgcg gctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgctaattgcattgcg gctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgctaattgcattgcg gctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgctaattgcattgcg atgactatgctaagctgcggctatgctaattgcattgcggctatgctaagctcatgcggctatgctaagctggaaatgcatgcggctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgctaattgcattgcg gcatgcggctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgctaattgcattgcg ggatccgatgactatgctaagctgcggctatgctaattgcattgcggctatgctaagctcggtatgctaattgaatgtcttggatttaccttggaaatgctaagctggatccgatgacaatgcattgcggctatgctaattgaatggtcttggatttaccttggaaatatgctaattgcattgcggctatgctaagctggaaatgcatgcggctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcg gctatgctaattgcattgcggctatgctaagctcatgcggctatgctaattgcattgcggctatgctaattgcattgcg

Goal: Genome Annotations

aatgcatgcggctatgctaattgcattgcggctatgctaagctggatccgatgacaatgcattgcggctatgctaatt
gcatgcggctatgcaaggctggatccgatgactatgctaagctggatccgatgacaatgcattgcggctatgctaatt
aatgaatggtcttggattttaccttggaaatgtctaagctggatccgatgacaatgcattgcggctatgctaatt
tggtcttggattttaccttggaaatgtctaattgcattgcggctatgctaagctggatccgatgacaatgcattgcg
gctatgctaattgcattgcggctatgcaaggctggatccgatgactatgctaagctgcggctatgctaattgcattgcg
gctatgctaagctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgctaagctggatcc
gctatgctaattgcattgcggctatgctaattgcattgcggctatgctaagctggatccgatgacaatgcattgcgg
atgctaattgcattgcggctatgctaattgcattgcggctatgctaagctggatccgatgacaatgcattgcgg
gctatgctaagctggatccgatgacaatgcattgcggctatgctaagctggatccgatgacaatgcattgcgg
atgactatgctaagctgcggctatgctaattgcattgcggctatgctaagctcatgcggctatgctaagctgg
gcatgcggctatgctaaggctggatccgatgacaatgcattgcggctatgctaattgcattgcggctatgctaagct
ggatccgatgactatgctaagctgcggctatgctaattgcattgcggctatgctaagctcggtatgctaatt
gtcttggattttaccttggaaatgtctaagctggatccgatgacaatgcattgcggctatgctaatt
gatttaccttggaaatgtctaattgcattgcggctatgctaagctggatgcattgcggctatgctaagctgg
cgatgacaatgcattgcggctatgctaattgcattgcggctatgctaagctggatccgatgactatgctaagctgc
gctatgctaattgcattgcggctatgctaagctcatgcgg

Gene!



Outline

1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. Experimental & Functional Assays



Outline

1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. Experimental & Functional Assays

Basic Local Alignment Search Tool

- Rapidly compare a sequence Q to a database to find all sequences in the database with a 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

FAKDFLAGGVAAAISKTAVAPIERVKLLLQVQHASKQITADKQYKGIIDCVVRIPKEQGV

FLIDLASGGTAAAVSKTAVAPIERVKLLLQVQDASKAIAVDKRYKGIMDVLIRVPKEQGV

- 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

Seed and Extend

FAKDFLAGGVAAAISKTAVAPIERVKLLLQVQHASKQITADKQYKGIIDCVVRIPKEQGV
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
FLIDLASGGTAAAVSKTAVAPIERVKLLLQVQDASKAIAVDKRYKGIMDVLIRVPKEQGV

- 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 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 LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-----DLSHGSAQV 55

L+P +K+ V A WGKV + E G EAL R+ + +P T+ +F F D G+ +V

Sbjct 3 LTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPKV 60

Query 56 KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRDPVNFKLLSHCLLVTLAHLPA 115

K HGKKV A ++ +AH+D++ + LS+LH KL VDP NF+LL + L+ LA H

Sbjct 61 KAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGK 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	LSPADKTNVKAAGKVGAGAHEYGAELERMFLSFPTTKTYFPF-----DLSHGSAQV	55
		LS + V WGKV A +G E L R+F P T F F D S +	
Sbjct	3	LSDGEWQLVVLNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKEKHLKSEDEMKAEDL	62

Query	56	KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPA	115
		K HG V AL + + L+ HA K ++ + +S C++ L + P	
Sbjct	63	KKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPG	122

Query	116	EFTPAVHASLDKFLASVSTVLTSKYR 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 ERMFLSFPTTKTYFPHFDSLHGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAAH 89
++M ++P P+F+ +H + + +A AL N ++DD+ +LSA D

Sbjct 59 QKMLGNYPEV---LPYFNKAHQISL--SQPRILAFALLNYAKNIDDL-TSLSAFMDQIVV 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

The screenshot shows the NCBI BLAST homepage. At the top, there's a navigation bar with links for Home, Recent Results, Saved Strategies, Help, and My NCBI (Sign In/Register). Below the navigation, a banner says "BLAST finds regions of similarity between biological sequences." and includes a link to "Primer-BLAST". A news sidebar on the right discusses the transition from the two-sequence aligner to the standard BLAST submission form, dated Tuesday, 03 Feb 2009, 16:00:00 EST. It also links to "More BLAST News...". A "Tip of the Day" section provides instructions for doing batch BLAST jobs. The main content area is divided into sections: "BLAST Assembled Genomes" (listing Human, Mouse, Rat, Arabidopsis thaliana, Oryza sativa, Bos taurus, Danio rerio, Gallus gallus, Pan troglodytes, Microbes, and Apis mellifera melanogaster), "Basic BLAST" (listing nucleotide blast, protein blast, blastx, tblastn, and tblastx), and "Specialized BLAST" (listing Primer-BLAST, trace archives, conserved domains, conserved domain architecture, gene expression profiles, IgBLAST, and SNPs). The bottom of the page shows the URL <http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Nucleotides&PROGRAM=blastn&MECABLAST=on&BLA...>, the FoxyProxy status, and a page number of 53.

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



Outline

1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. Experimental & Functional Assays



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

Genetic Code

	Second letter				
	U	C	A	G	
First letter	UUU } Phe UUC UUA } Leu UUG }	UCU } UCC UCA } Ser UCG }	UAU } Tyr UAC UAA Stop UAG Stop	UGU } Cys UGC UGA Stop UGG Trp	U C A G
C	CUU } CUC CUA } Leu CUG }	CCU } CCC CCA } Pro CCG }	CAU } His CAC CAA } Gln CAG }	CGU } CGC CGA } Arg CGG }	U C A G
A	AUU } AUC AUA } Ile AUG Met	ACU } ACC ACA } Thr ACG }	AAU } Asn AAC AAA } Lys AAG }	AGU } Ser AGC AGA } Arg AGG }	U C A G
G	GUU } GUC GUA } Val GUG }	GCU } GCC GCA } Ala GCG }	GAU } Asp GAC GAA } Glu GAG }	GGU } GGC GGA } Gly GGG }	U C A G

- Start:
- AUG
- Stop:
- UAA
 - UAG
 - UGA

Step One

- Find open reading frames (ORFs).



...TAGATGAATGGCTCTTTAGATAAATTTCATGAAAAAATTGA...

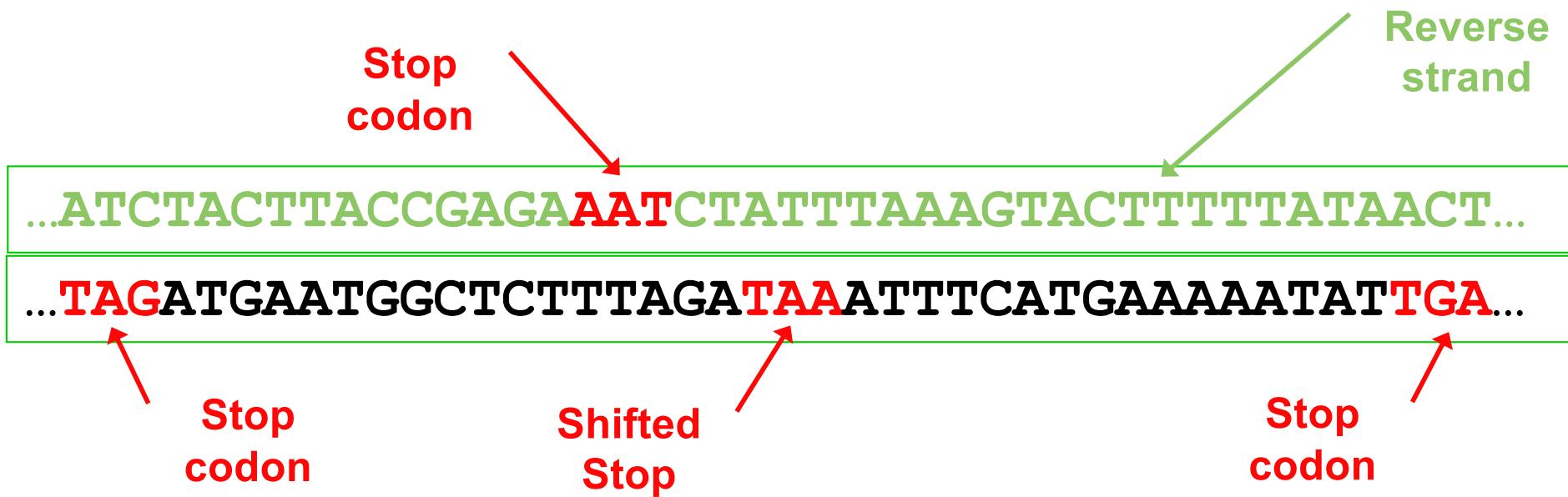
Start codon

Stop codon

Stop codon

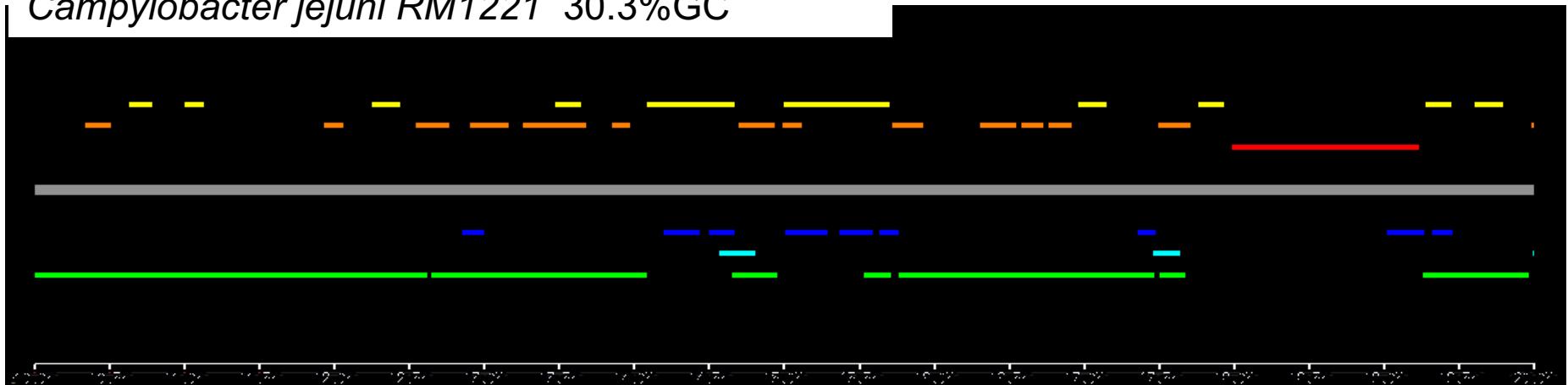
Step One

- Find open reading frames (ORFs).



- But ORFs generally overlap ...

Campylobacter jejuni RM1221 30.3%GC

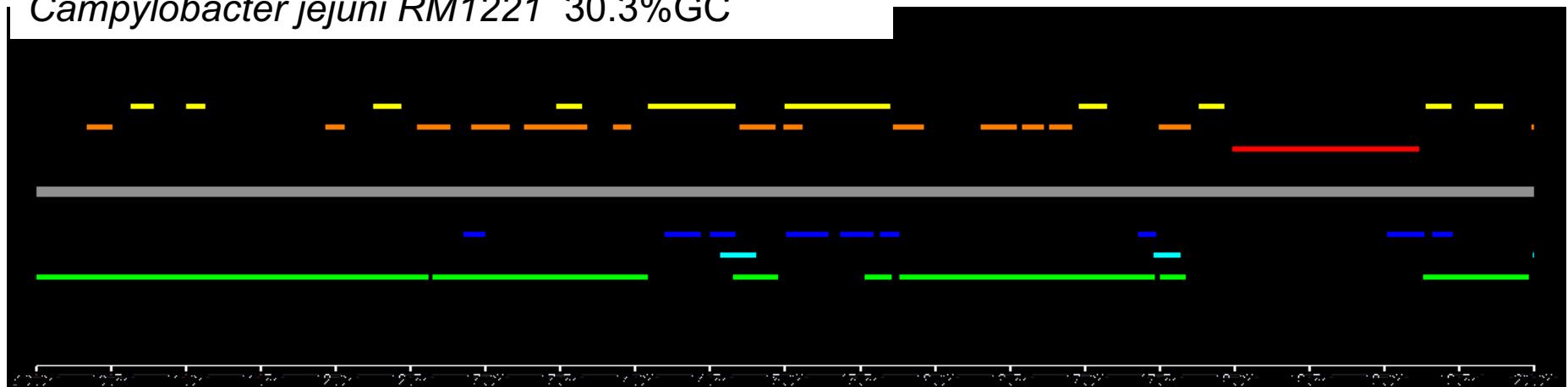


All ORFs longer than 100bp on both strands shown
- color indicates reading frame
Longest ORFs likely to be protein-coding genes

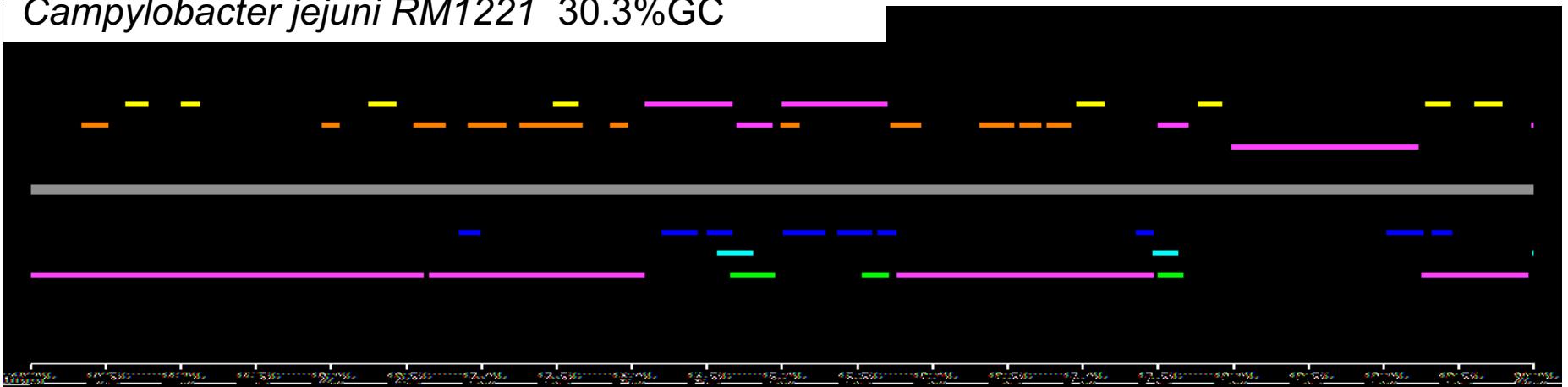
Note the low GC content

All genes are ORFs but not all ORFs are genes

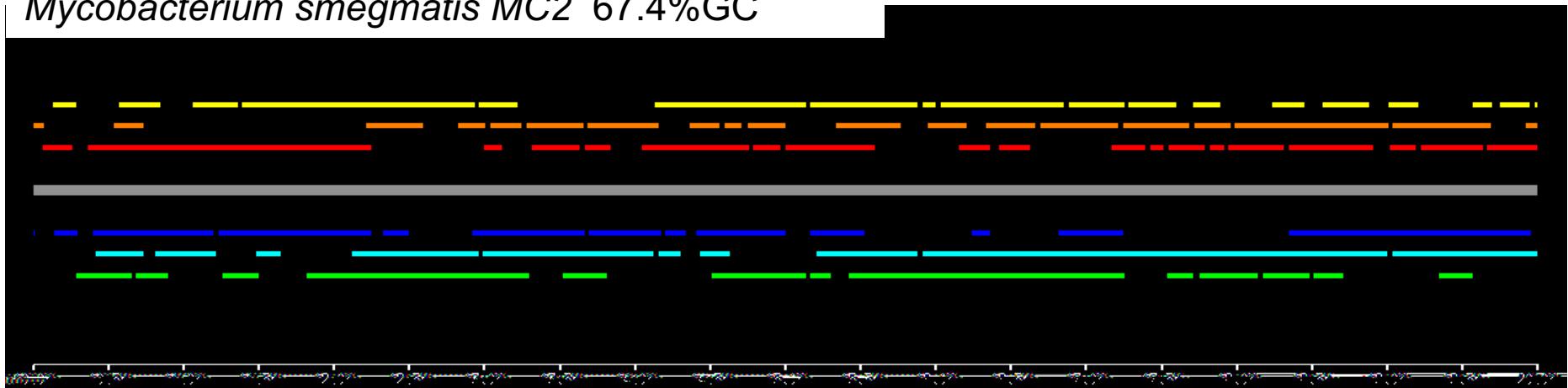
Campylobacter jejuni RM1221 30.3%GC



Campylobacter jejuni RM1221 30.3%GC

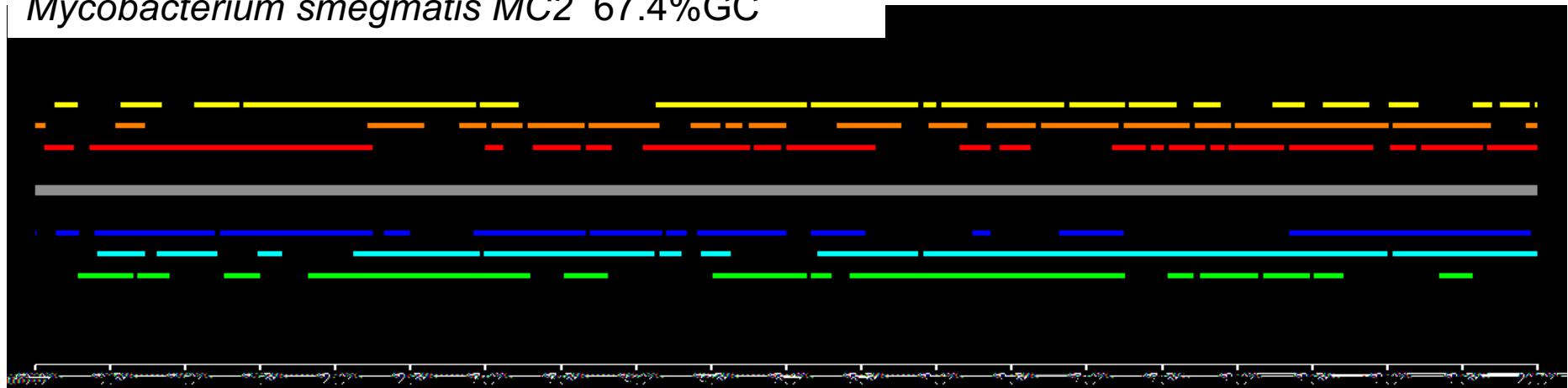


Mycobacterium smegmatis MC2 67.4%GC

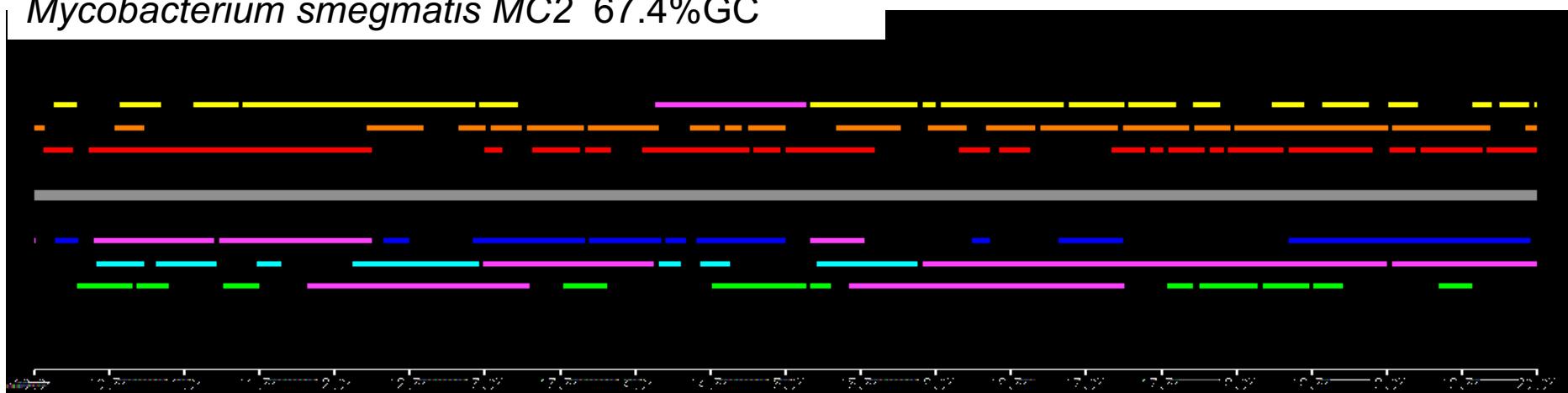


Note what happens in a high-GC genome

Mycobacterium smegmatis MC2 67.4%GC



Mycobacterium smegmatis MC2 67.4%GC



Flipping a Biased Coin

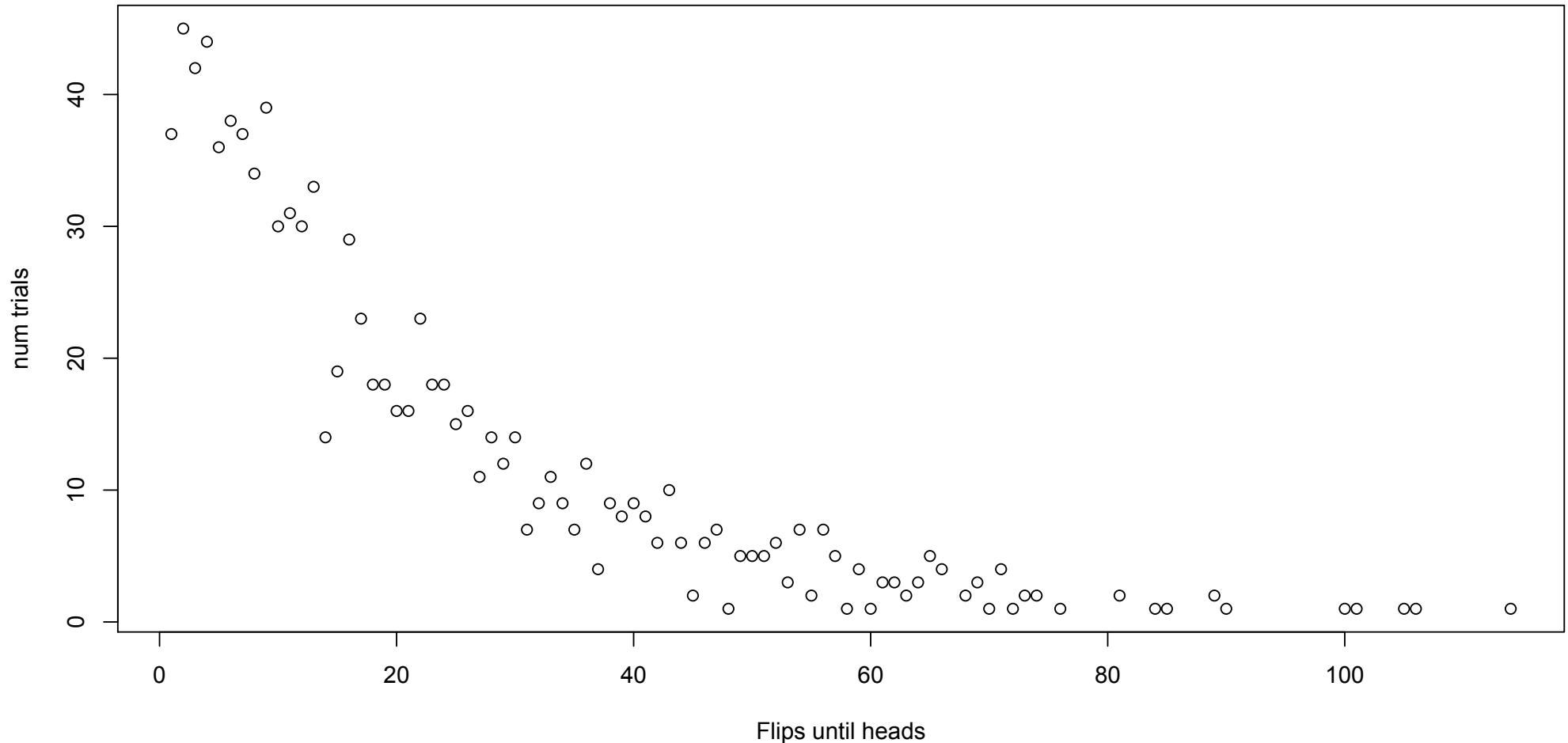
$P(\text{heads}) = 61/64 (95.4\%)$ $P(\text{tails}) = 3/64 (4.6\%)$

How many flips until my first tail?

Flipping a Biased Coin

$$P(\text{heads}) = 61/64 (95.4\%) \quad P(\text{tails}) = 3/64 (4.6\%)$$

How many flips until my first tail?

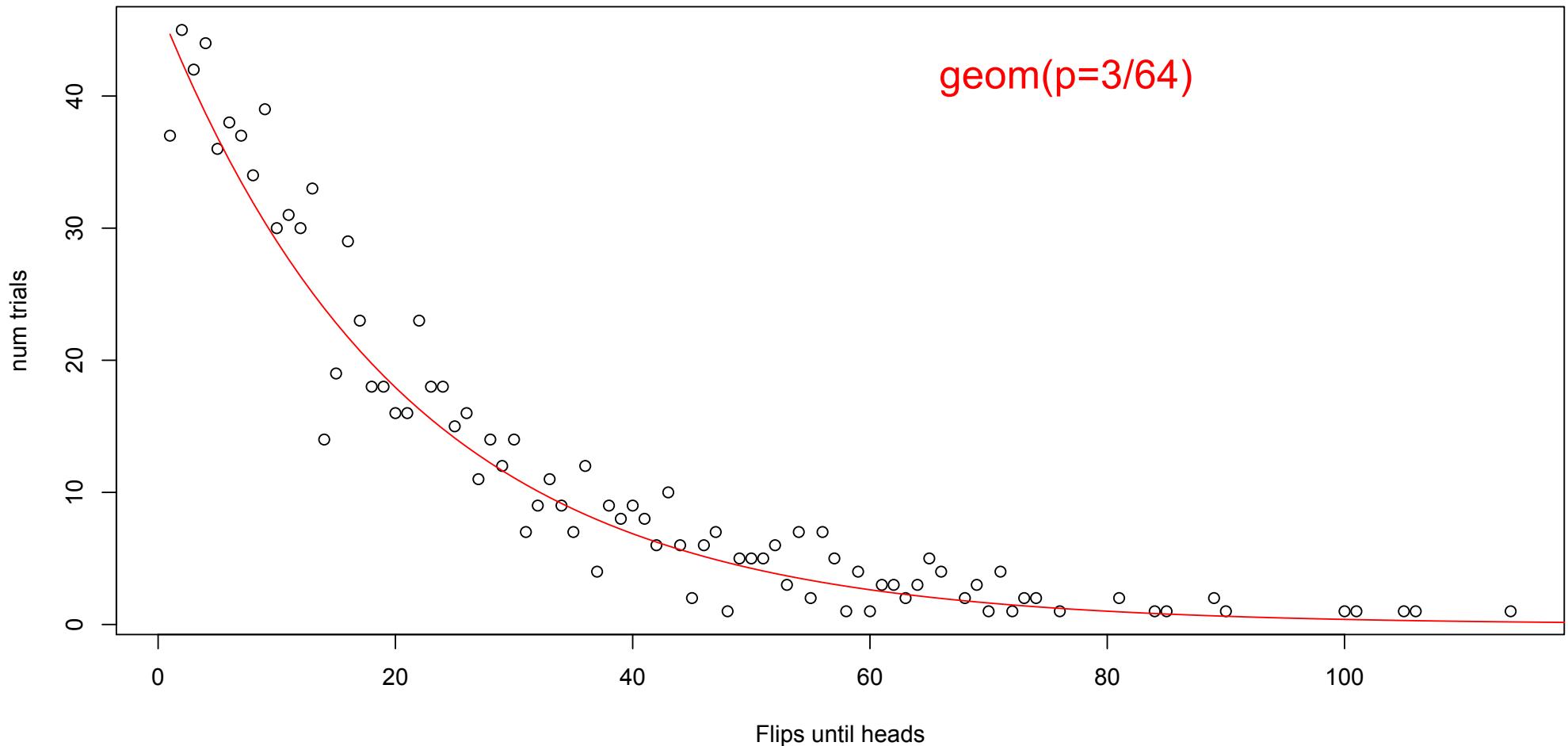


Flipping a Biased Coin

$$P(\text{heads}) = 61/64 (95.4\%) \quad P(\text{tails}) = 3/64 (4.6\%)$$

How many flips until my first tail?

Geometric Distribution: $P(X=x) = p_{\text{heads}}^{x-1} p_{\text{tails}}$

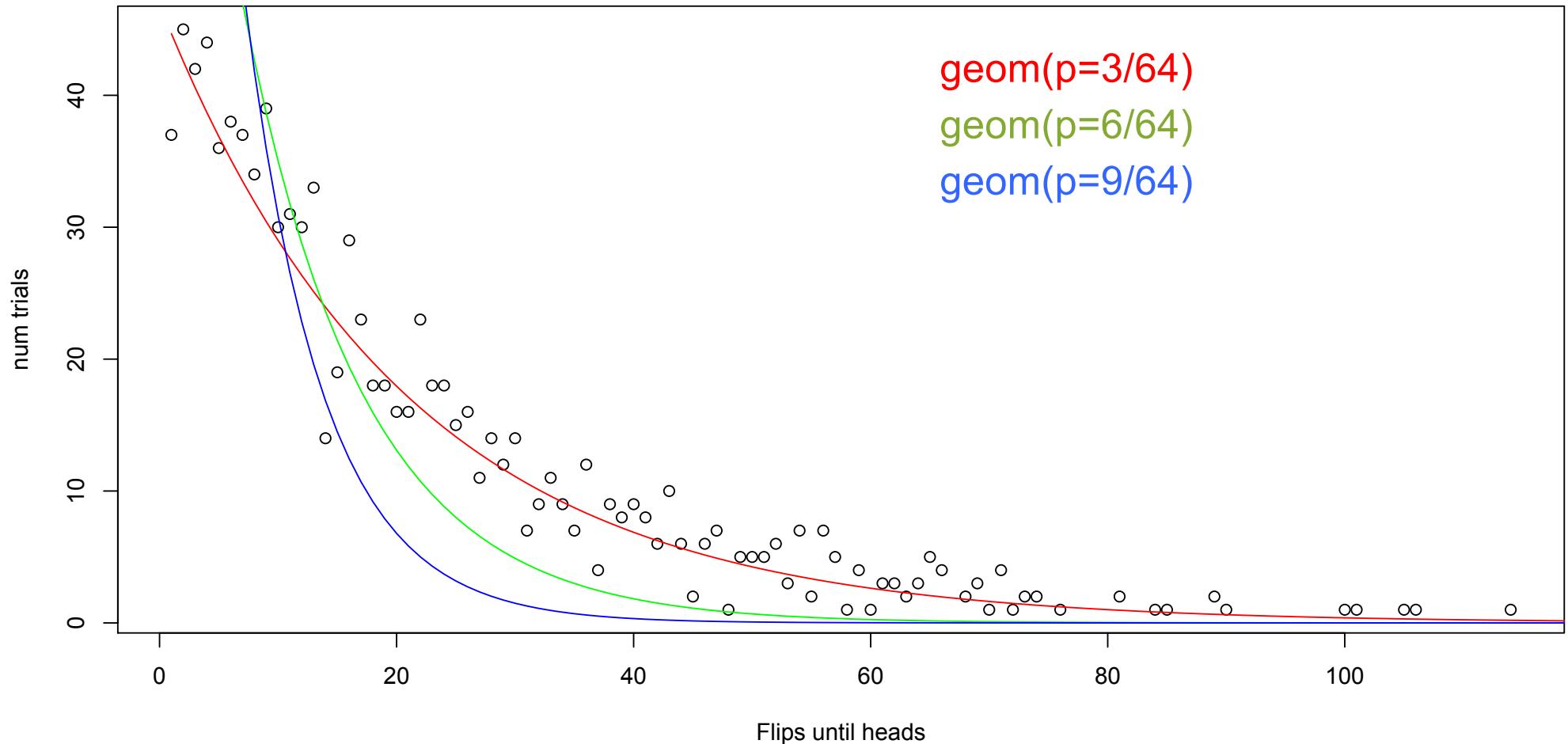


Flipping a Biased Coin

$$P(\text{heads}) = 61/64 (95.4\%) \quad P(\text{tails}) = 3/64 (4.6\%)$$

How many flips until my first tail?

Geometric Distribution: $P(X=x) = p_{\text{heads}}^{x-1} p_{\text{tails}}$

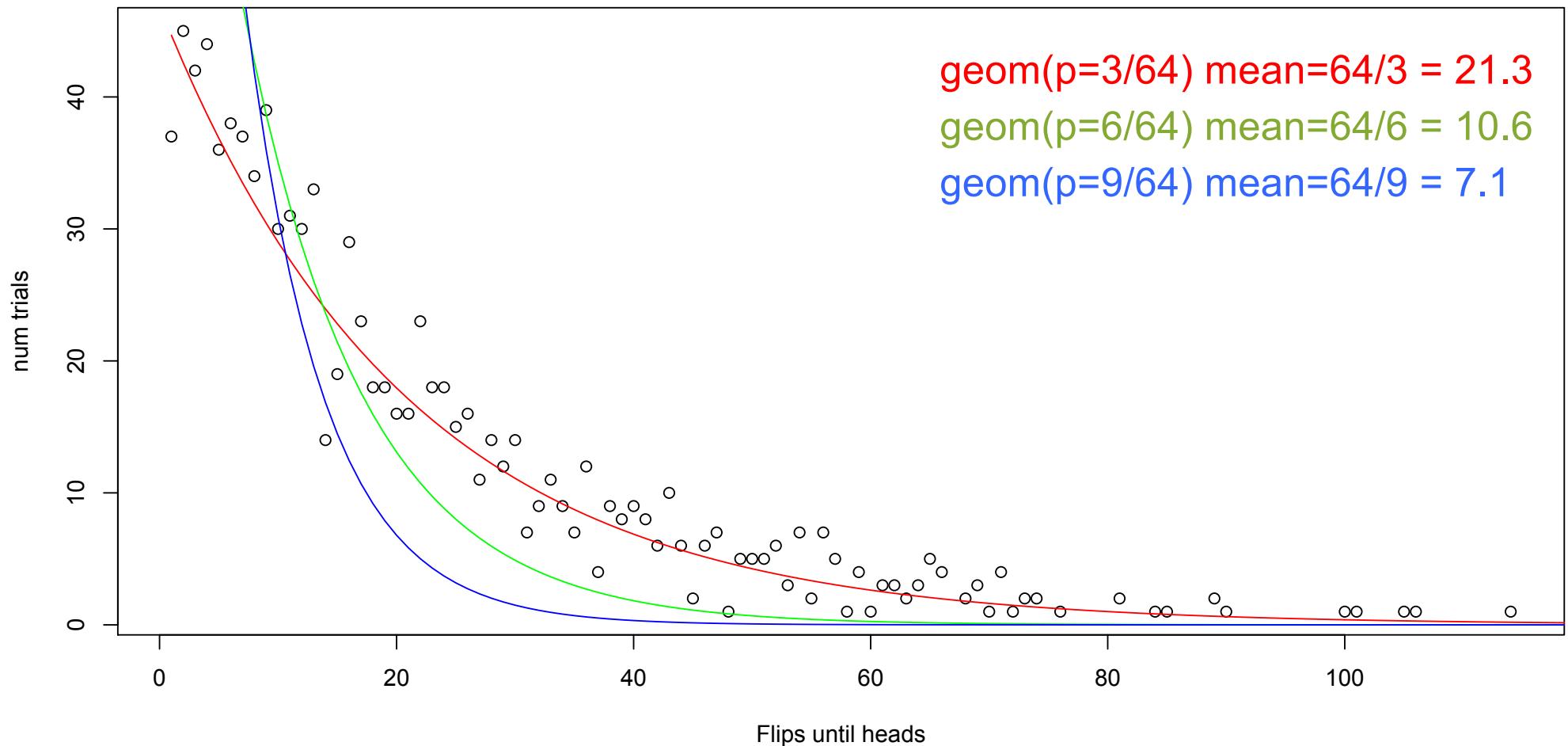


Flipping a Biased Coin

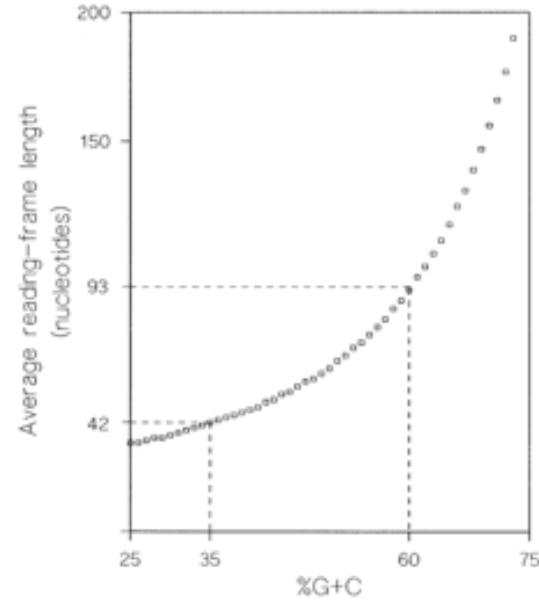
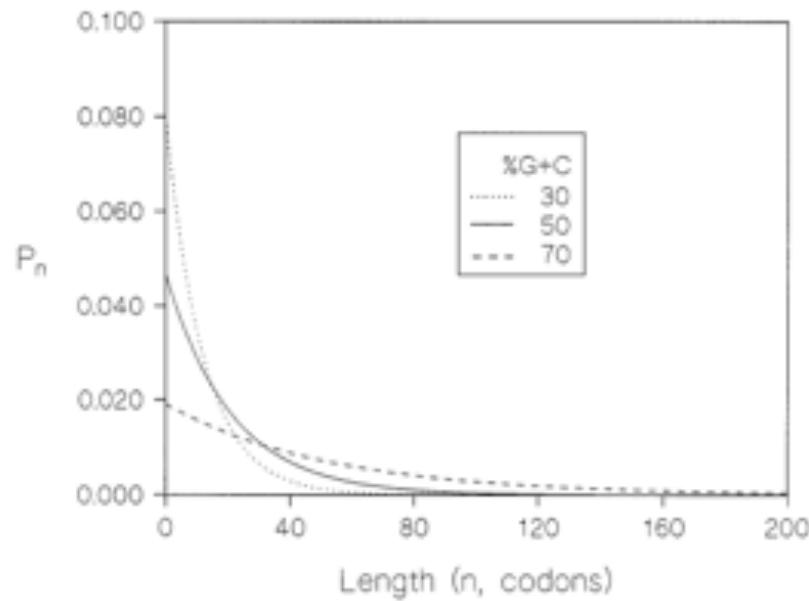
$$P(\text{heads}) = 61/64 (95.4\%) \quad P(\text{tails}) = 3/64 (4.6\%)$$

How many flips until my first tail?

Geometric Distribution: $P(X=x) = p_{\text{heads}}^{x-1} p_{\text{tails}}$



Stop Codon Frequencies



If the sequence is mostly A+T, then likely to form stop codons by chance!

In High A+T (Low G+C):

Frequent stop codons; Short Random ORFs; long ORFs likely to be true genes

In High G+C (Low A+T):

Rare stop codons; Long Random ORFs; harder to identify true genes

A relationship between GC content and coding-sequence length.

Oliver & Marín (1996) J Mol Evol. 43(3):216-23.

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

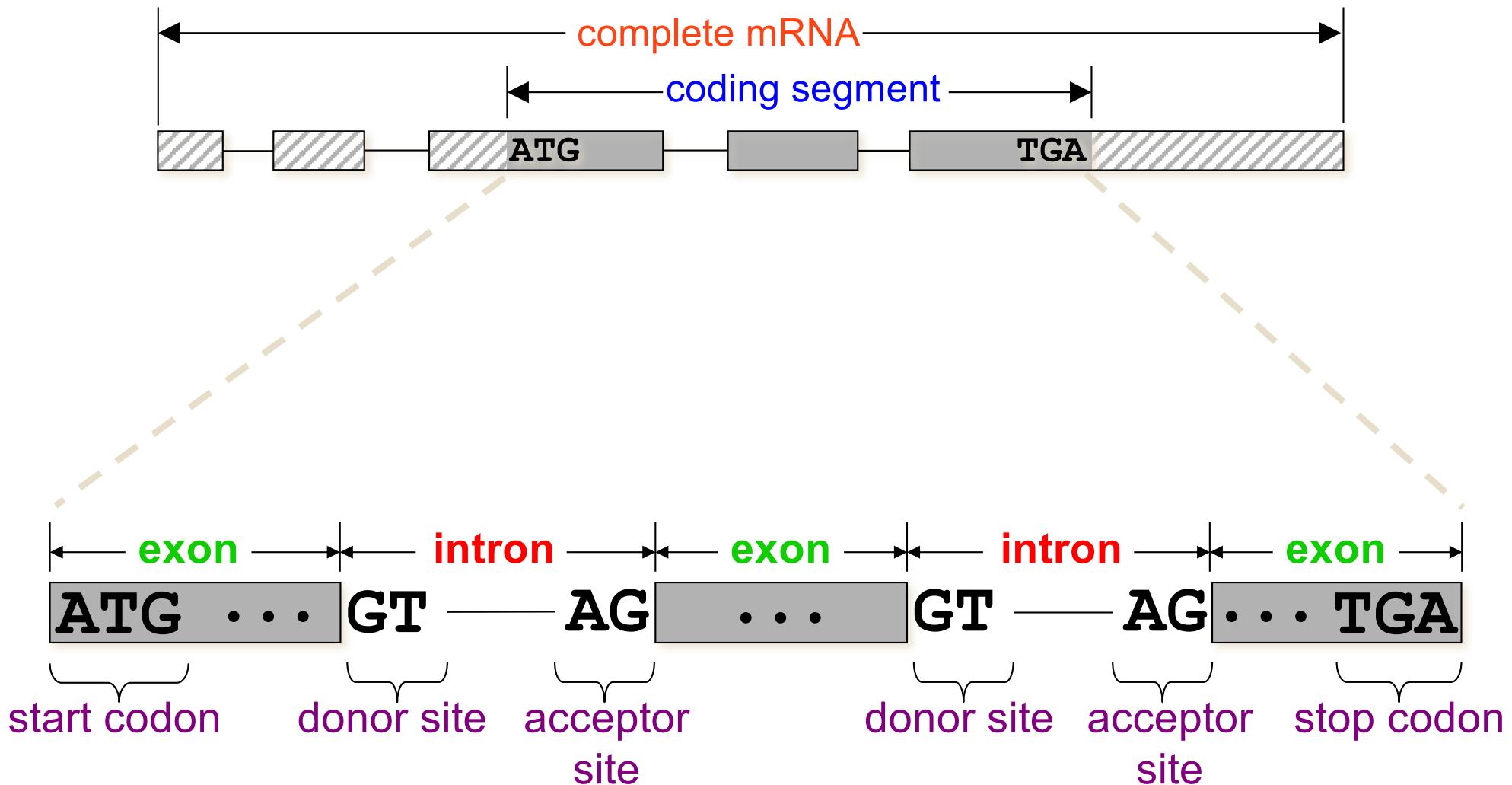


Overview of Eukaryotic Gene Prediction

CBB 231 / COMPSCI 261

W.H. Majoros

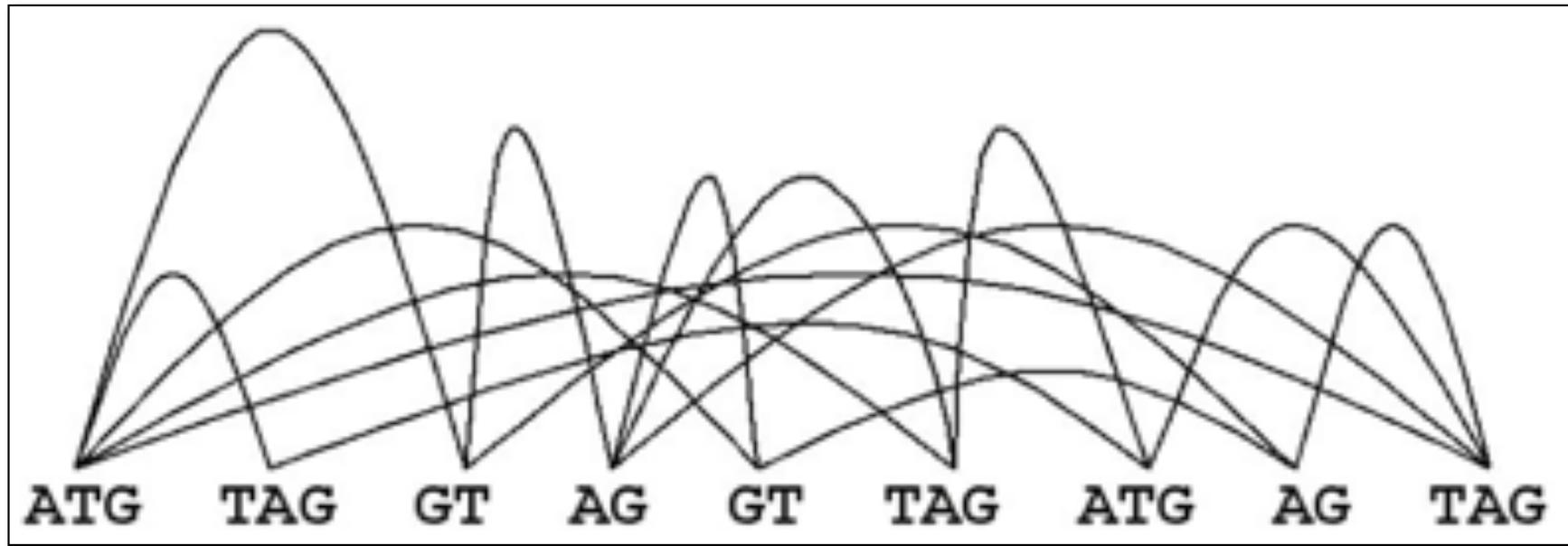
Eukaryotic Gene Syntax



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

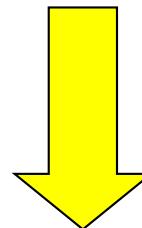
After identifying the most promising (i.e., highest-scoring) signals in an input sequence, we can apply the gene syntax rules to connect these into an *ORF graph*:



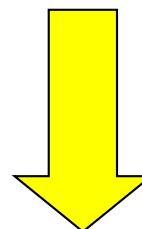
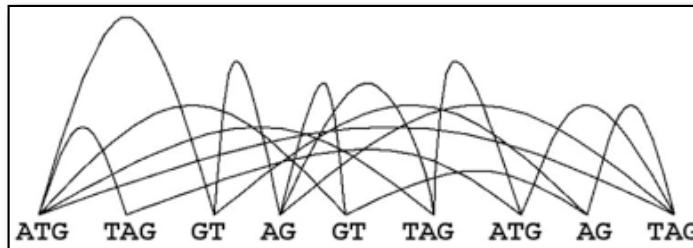
An ORF graph represents all possible *gene parses* (and their scores) for a given set of putative signals. A *path* through the graph represents a single gene parse.

Conceptual Gene-finding Framework

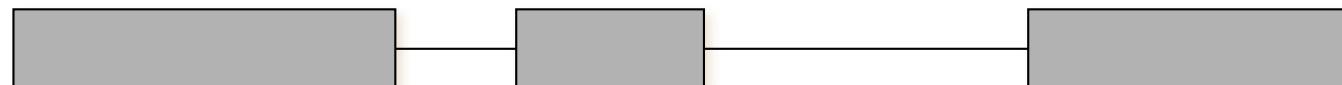
TATTCCGATCGATCGATCTCTCTAGCGTCTACG
CTATCATCGCTCTCTATTATCGCGCGATCGTCG
ATCGCGCGAGAGTATGCTACGTGATCGAATTG



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

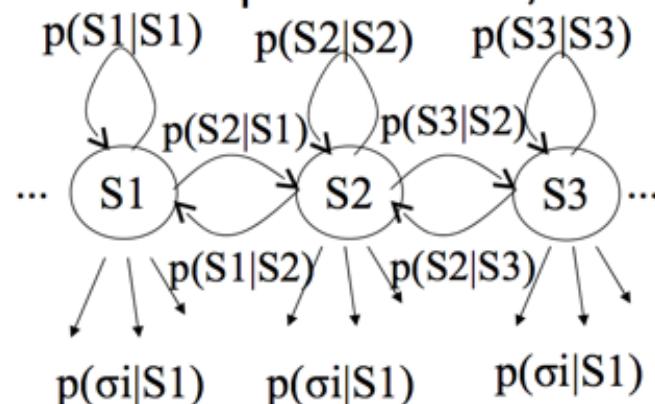


find highest-scoring path through ORF graph;
interpret path as a gene parse = gene structure



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.



AAAGCATGCATTAACGTGAGCACAAATAGATTACA



Eukaryotic Gene Finding with **GlimmerHMM**

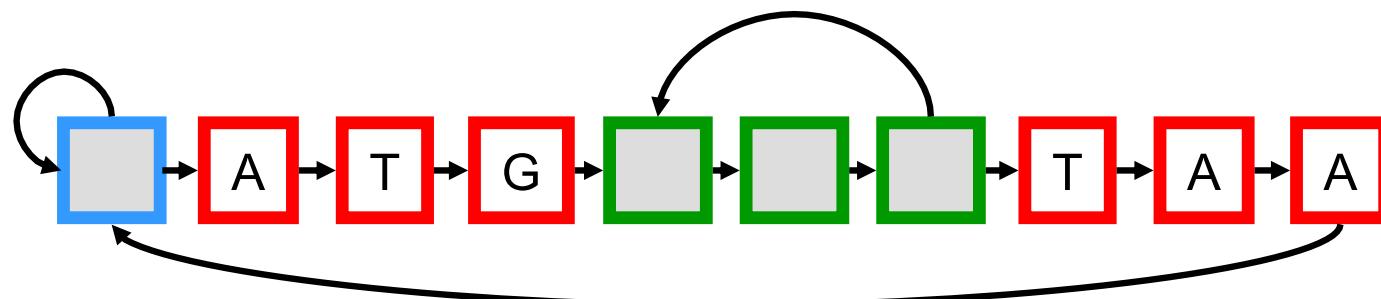
Mihaela Pertea

JHU

HMMs and Gene Structure

- Nucleotides $\{A,C,G,T\}$ are the observables
- Different states generate nucleotides at different frequencies

A simple HMM for unspliced genes:



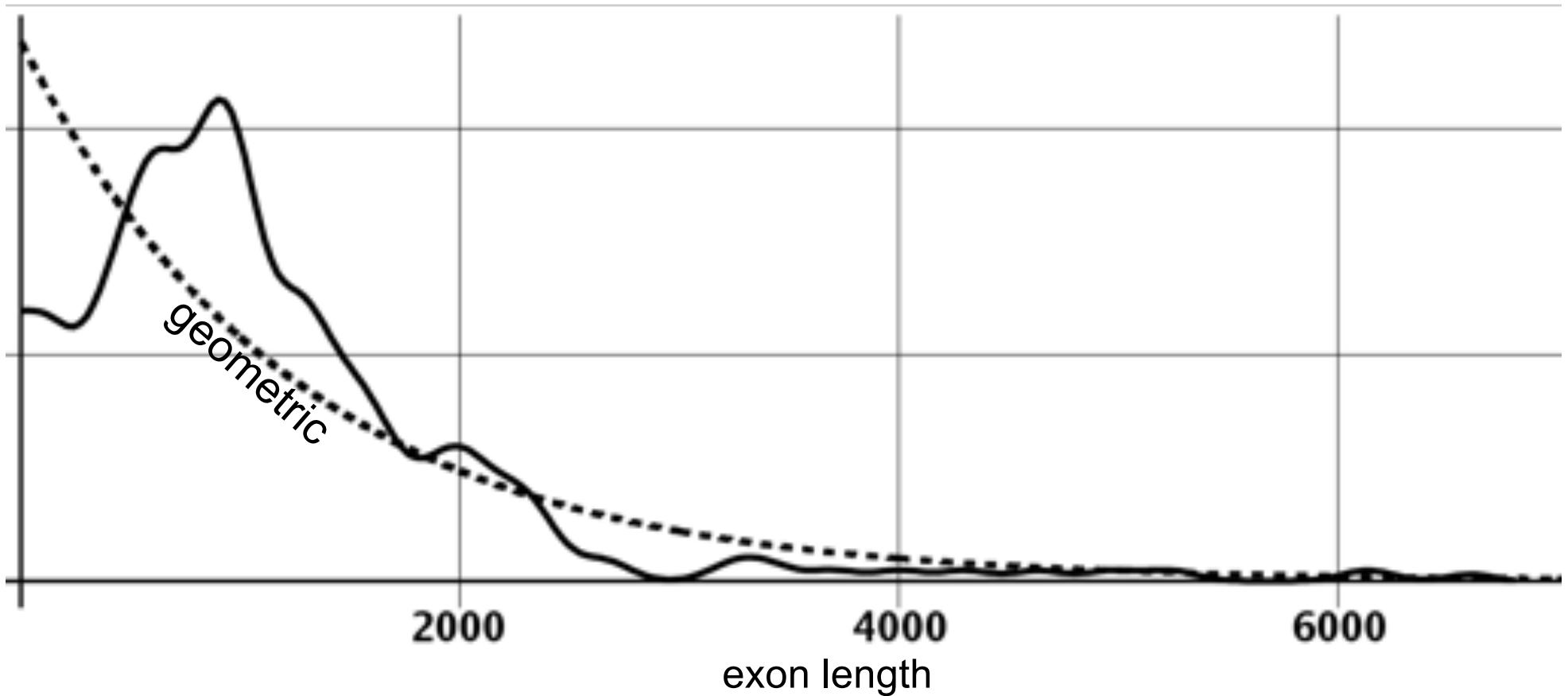
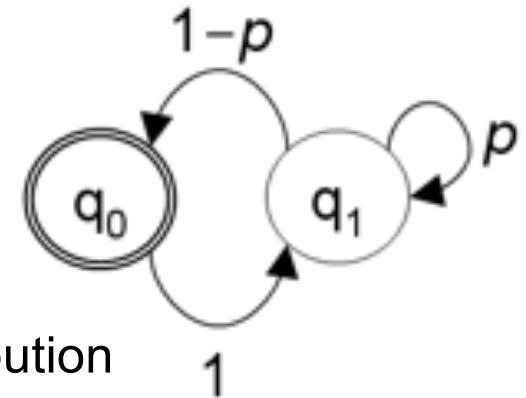
AAAGC ATG CAT TTA ACG AGA GCA CAA GGG CTC TAA TGCCG

- The sequence of states is an annotation of the generated string – each nucleotide is generated in **intergenic**, **start/stop**, **coding** state

HMMs & Geometric Feature Lengths

$$P(x_0 \dots x_{d-1} | \theta) = \left(\prod_{i=0}^{d-1} P_e(x_i | \theta) \right) p^{d-1} (1-p)$$

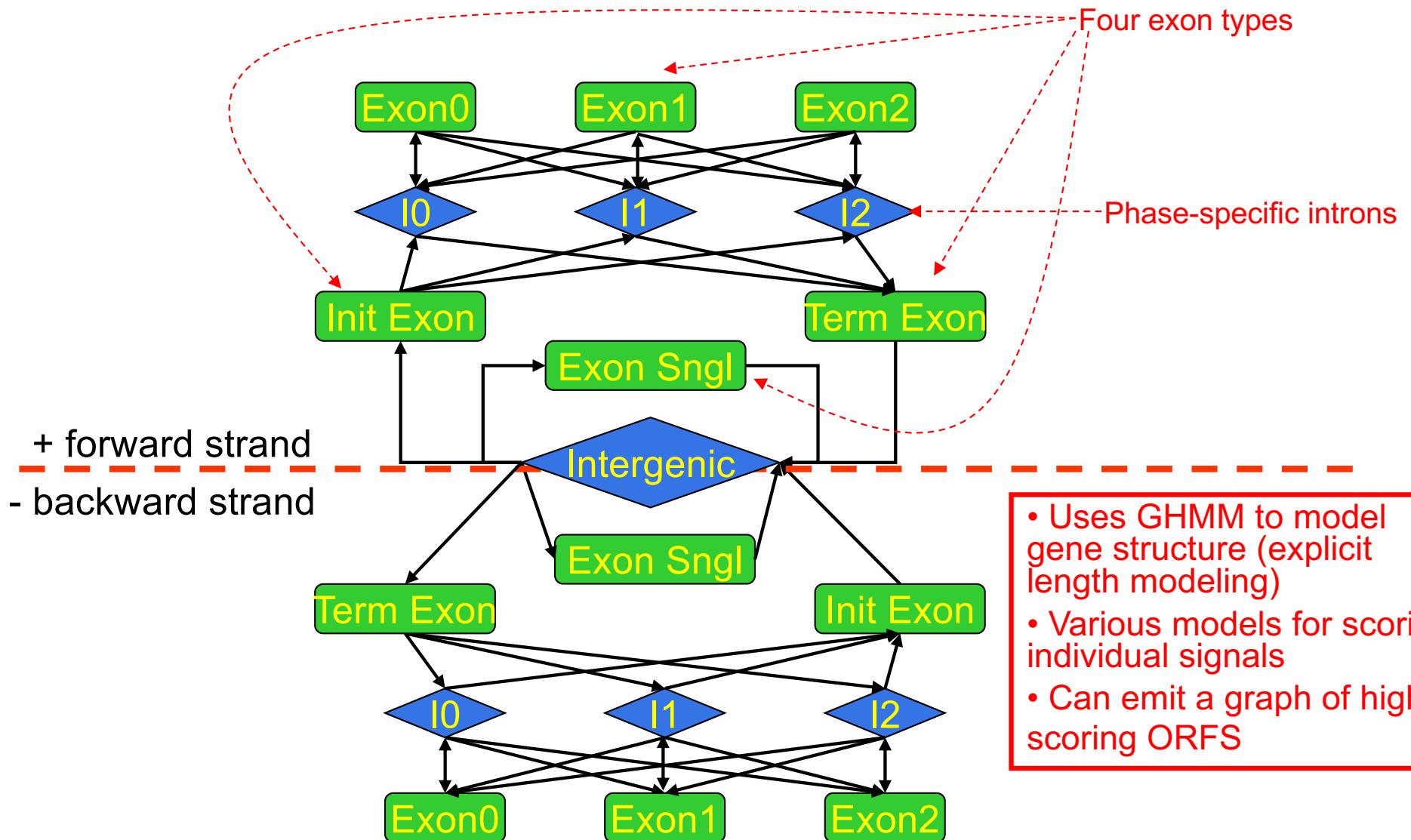
geometric distribution



Generalized HMMs Summary

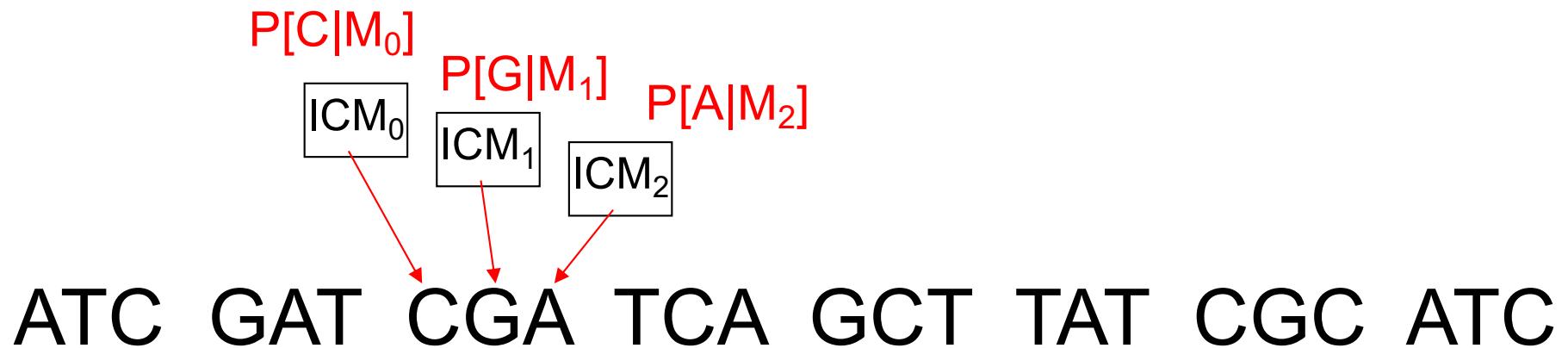
- GHMMs generalize HMMs by allowing each state to emit a **subsequence** rather than just a single symbol
- Whereas HMMs model all feature lengths using a **geometric distribution**, coding features can be modeled using an arbitrary **length distribution** in a GHMM
- Emission models within a GHMM can be any arbitrary probabilistic model (“**submodel abstraction**”), such as a neural network or decision tree
- GHMMs tend to have many **fewer states** => simplicity & modularity

GlimmerHMM architecture



Coding vs Non-coding

A three-periodic ICM uses three ICMs in succession to evaluate the different codon positions, which have different statistics:



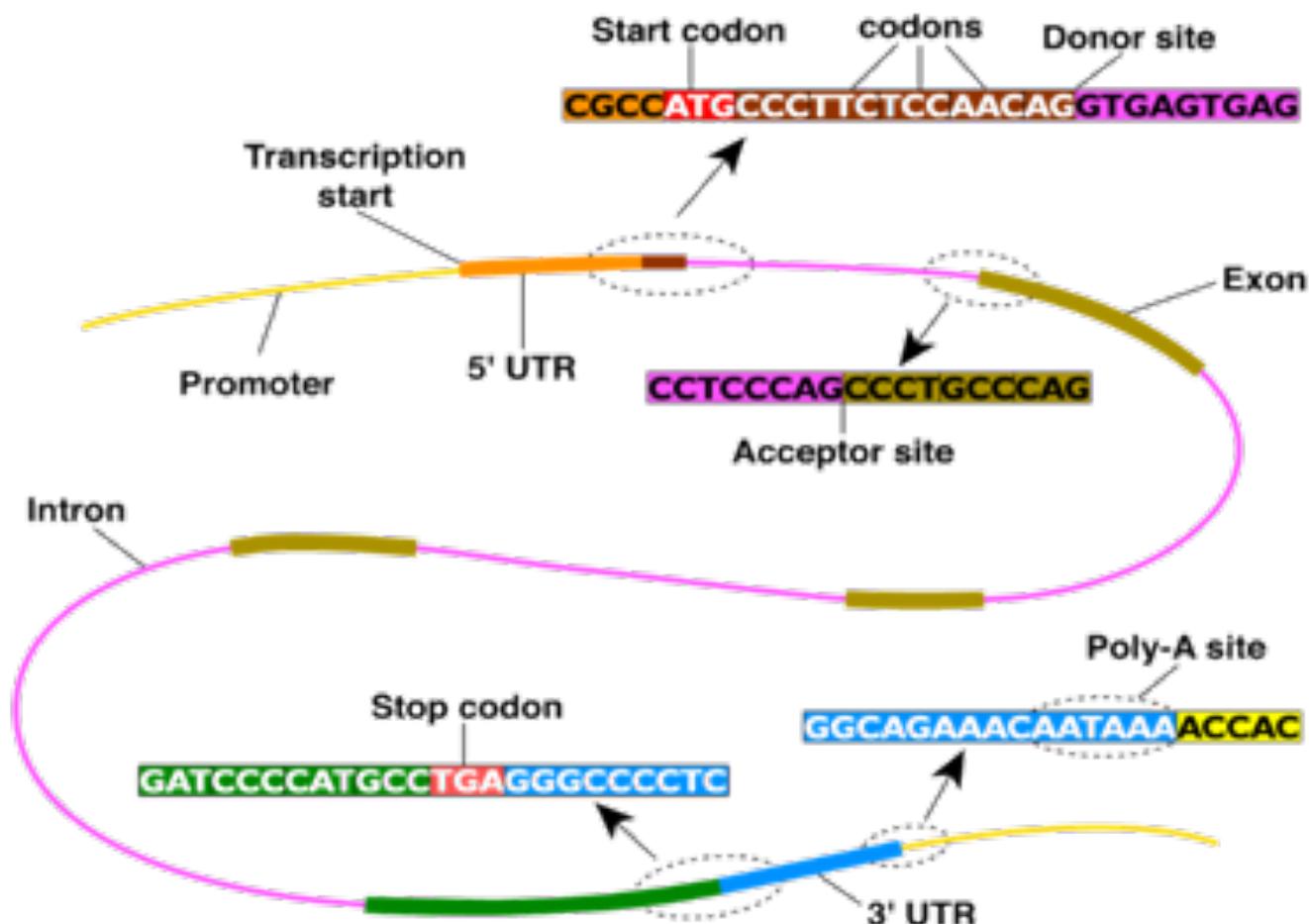
The three ICMs correspond to the three phases. Every base is evaluated in every phase, and the score for a given stretch of (putative) coding DNA is obtained by multiplying the phase-specific probabilities in a mod 3 fashion:

$$\prod_{i=0}^{L-1} P_{(f+i)(\text{mod } 3)}(x_i)$$

GlimmerHMM uses 3-periodic ICMs for coding and homogeneous (non-periodic) ICMs for noncoding DNA.

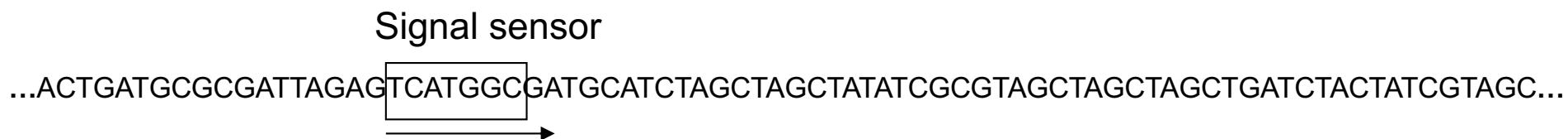
Signal Sensors

Signals – short sequence patterns in the genomic DNA that are recognized by the cellular machinery.



Identifying Signals In DNA

We slide a fixed-length model or “window” along the DNA and evaluate score (signal) at each point:

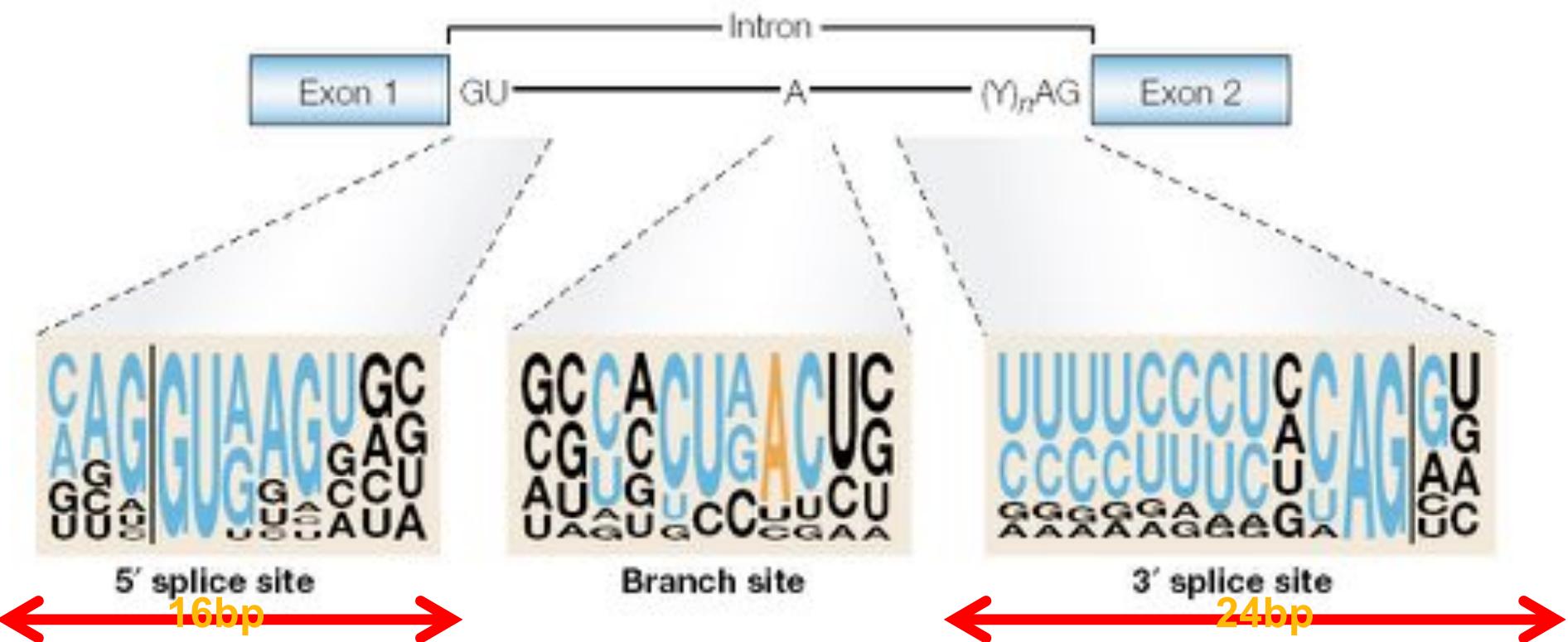


When the score is greater than some threshold (determined empirically to result in a desired sensitivity), we remember this position as being the potential site of a signal.

The most common signal sensor is the Position Weight Matrix:

A = 31% T = 28% C = 21% G = 20%	A = 18% T = 32% C = 24% G = 26%	A 100%	T 100%	G 100%	A = 19% T = 20% C = 29% G = 32%	A = 24% T = 18% C = 26% G = 32%
--	--	------------------	------------------	------------------	--	--

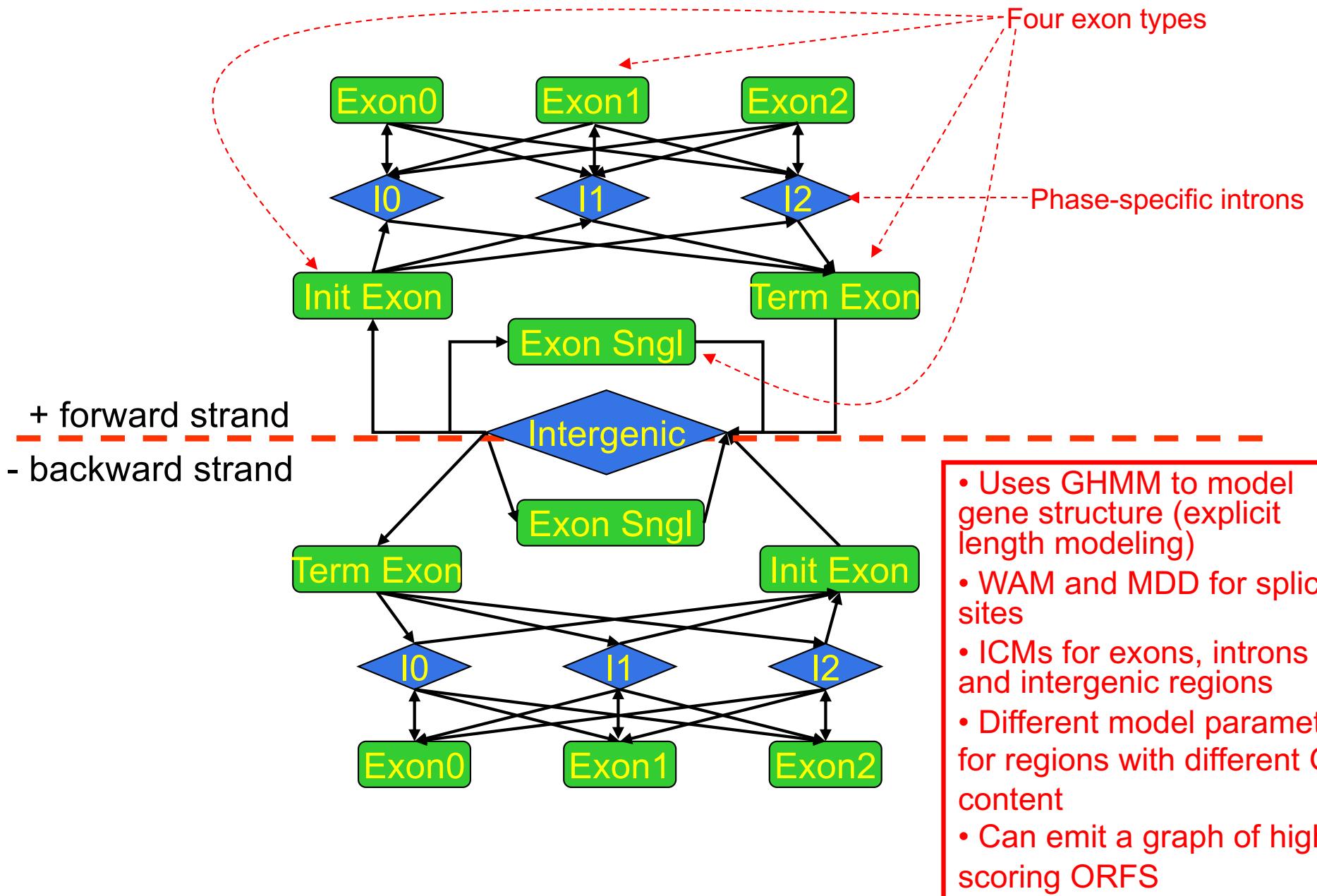
Splice site prediction



The splice site score is a combination of:

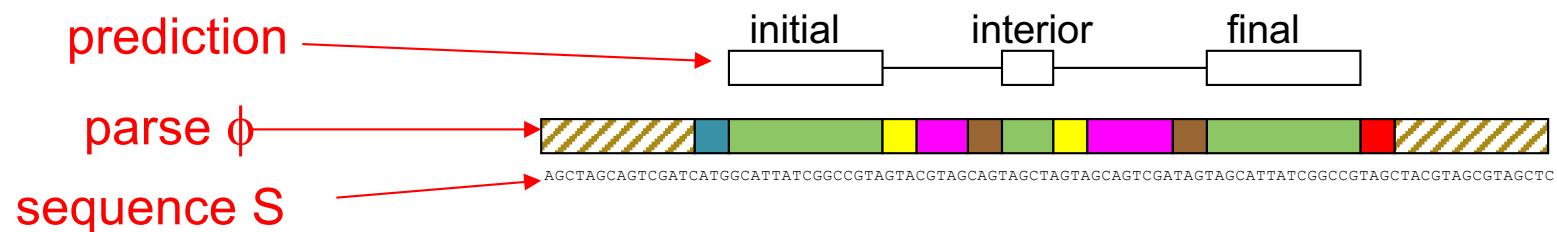
- first or second order inhomogeneous Markov models on windows around the acceptor and donor sites
- Maximal dependence decomposition (MDD) decision trees
- longer Markov models to capture difference between coding and non-coding on opposite sides of site (optional)
- maximal splice site score within 60 bp (optional)

GlimmerHMM architecture



Gene Prediction with a GHMM

Given a sequence S , we would like to determine the parse ϕ of that sequence which segments the DNA into the most likely exon/intron structure:

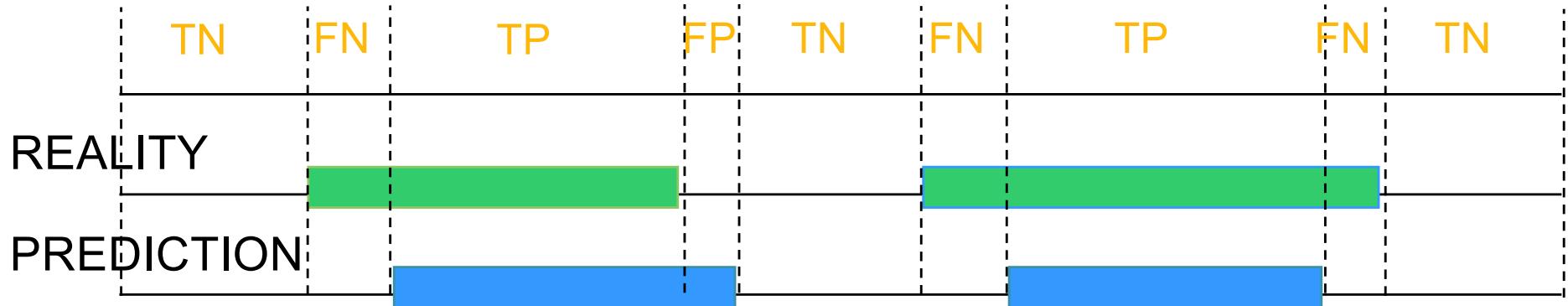


The parse ϕ consists of the coordinates of the predicted exons, and corresponds to the precise sequence of states during the operation of the GHMM (and their duration, which equals the number of symbols each state emits).

This is the same as in an HMM except that in the HMM each state emits bases with fixed probability, whereas in the GHMM each state emits an entire feature such as an exon or intron.

Evaluation of Gene Finding Programs

Nucleotide level accuracy



Sensitivity:

$$Sn = \frac{TP}{TP + FN}$$

What fraction of reality did you predict?

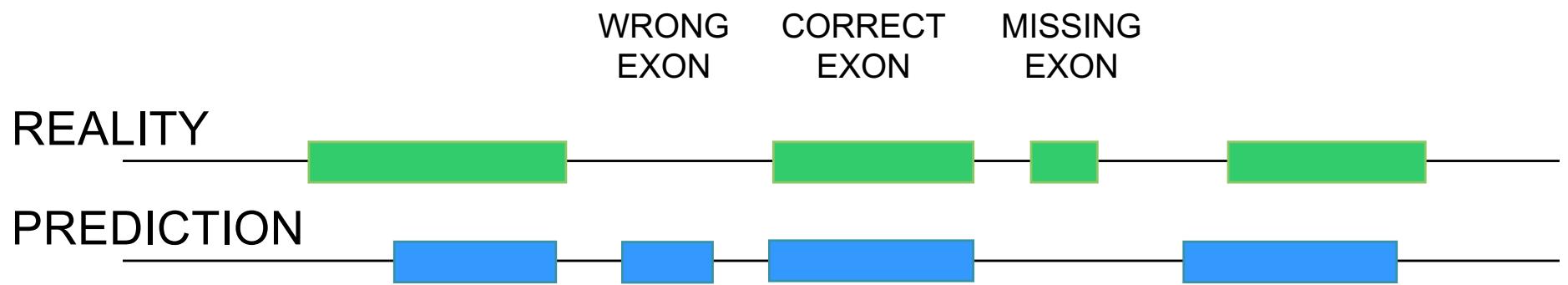
Specificity:

$$Sp = \frac{TN}{TP + FP}$$

What fraction of your predictions are real?

More Measures of Prediction Accuracy

Exon level accuracy



$$ExonSn = \frac{TE}{AE} = \frac{\text{number of correct exons}}{\text{number of actual exons}}$$

$$ExonSp = \frac{TE}{PE} = \frac{\text{number of correct exons}}{\text{number of predicted exons}}$$

GlimmerHMM is a high-performance ab initio gene finder

Arabidopsis thaliana test results

	Nucleotide			Exon			Gene		
	Sn	Sp	Acc	Sn	Sp	Acc	Sn	Sp	Acc
GlimmerHMM	97	99	98	84	89	86.5	60	61	60.5
SNAP	96	99	97.5	83	85	84	60	57	58.5
Genscan+	93	99	96	74	81	77.5	35	35	35

- All three programs were tested on a test data set of 809 genes, which did not overlap with the training data set of GlimmerHMM.
- All genes were confirmed by full-length Arabidopsis cDNAs and carefully inspected to remove homologues.

GlimmerHMM on human data

	<i>Nuc Sens</i>	<i>Nuc Spec</i>	<i>Nuc Acc</i>	<i>Exon Sens</i>	<i>Exon Spec</i>	<i>Exon Acc</i>	<i>Exact Genes</i>
<i>GlimmerHMM</i>	86%	72%	79%	72%	62%	67%	17%
<i>Genscan</i>	86%	68%	77%	69%	60%	65%	13%

GlimmerHMM's performance compared to Genscan on 963 human RefSeq genes selected randomly from all 24 chromosomes, non-overlapping with the training set. The test set contains 1000 bp of untranslated sequence on either side (5' or 3') of the coding portion of each gene.

Gene Finding Overview

- Prokaryotic gene finding distinguishes real genes and random ORFs
 - Prokaryotic genes have simple structure and are largely homogenous, making it relatively easy to recognize their sequence composition
- Eukaryotic gene finding identifies the genome-wide most probable gene models (set of exons)
 - “Probabilistic Graphical Model” to enforce overall gene structure, separate models to score splicing/transcription signals
 - Accuracy depends to a large extent on the quality of the training data



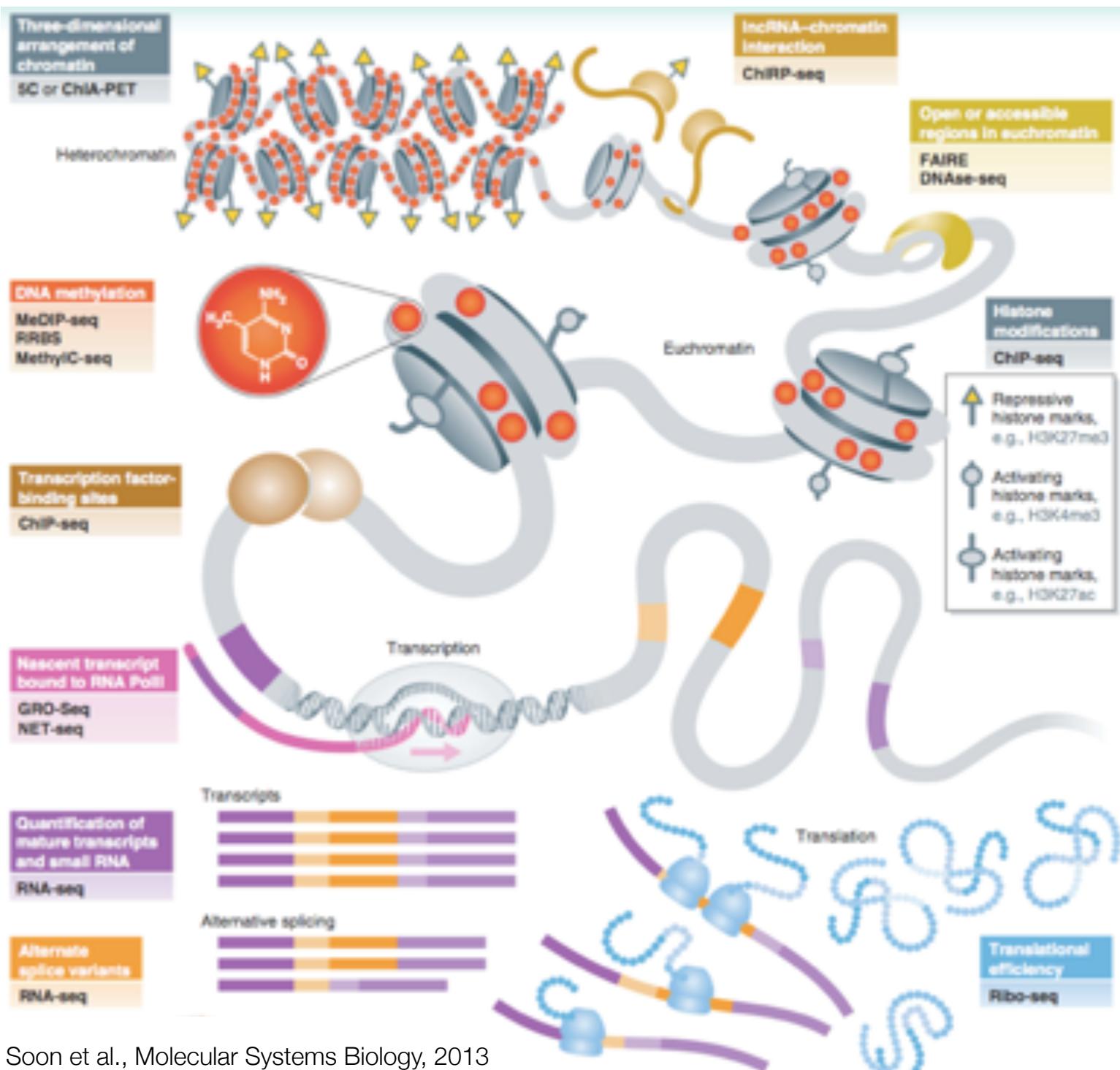
Outline

1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. **Experimental & Functional Assays**

Sequencing Assays

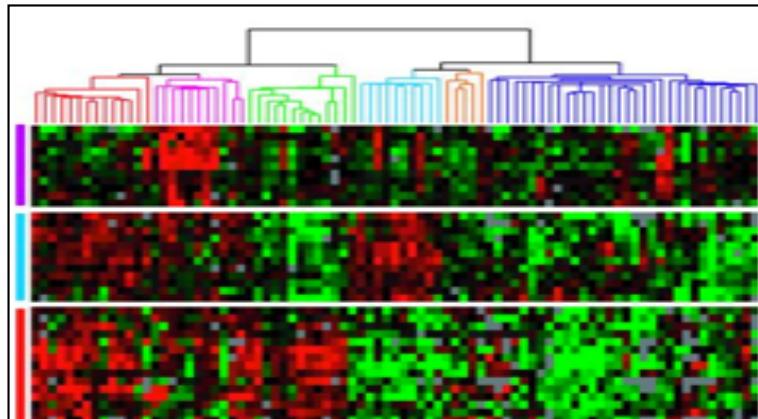
The *Seq List (in chronological order)

1. Gregory E. Crawford et al., “Genome-wide Mapping of DNase Hypersensitive Sites Using Massively Parallel Signature Sequencing (MPSS),” *Genome Research* 16, no. 1 (January 1, 2006): 123–131, doi:10.1101/gr.4074106.
2. David S. Johnson et al., “Genome-Wide Mapping of in Vivo Protein-DNA Interactions,” *Science* 316, no. 5830 (June 8, 2007): 1497–1502, doi:10.1126/science.1141319.
3. Tarjei S. Mikkelsen et al., “Genome-wide Maps of Chromatin State in Pluripotent and Lineage-committed Cells,” *Nature* 448, no. 7153 (August 2, 2007): 553–560, doi:10.1038/nature06008.
4. Thomas A. Down et al., “A Bayesian Deconvolution Strategy for Immunoprecipitation-based DNA Methylome Analysis,” *Nature Biotechnology* 26, no. 7 (July 2008): 779–785, doi:10.1038/nbt1414.
5. Ali Mortazavi et al., “Mapping and Quantifying Mammalian Transcriptomes by RNA-Seq,” *Nature Methods* 5, no. 7 (July 2008): 621–628, doi:10.1038/nmeth.1226.
6. Nathan A. Baird et al., “Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers,” *PLoS ONE* 3, no. 10 (October 13, 2008): e3376, doi:10.1371/journal.pone.0003376.
7. Leighton J. Core, Joshua J. Waterfall, and John T. Lis, “Nascent RNA Sequencing Reveals Widespread Pausing and Divergent Initiation at Human Promoters,” *Science* 322, no. 5909 (December 19, 2008): 1845–1848, doi:10.1126/science.1162228.
8. Chao Xie and Martti T. Tammi, “CNV-seq, a New Method to Detect Copy Number Variation Using High-throughput Sequencing,” *BMC Bioinformatics* 10, no. 1 (March 6, 2009): 80, doi:10.1186/1471-2105-10-80.
9. Jay R. Hesselberth et al., “Global Mapping of protein-DNA Interactions in Vivo by Digital Genomic Footprinting,” *Nature Methods* 6, no. 4 (April 2009): 283–289, doi:10.1038/nmeth.1313.
10. Nicholas T. Ingolia et al., “Genome-Wide Analysis in Vivo of Translation with Nucleotide Resolution Using Ribosome Profiling,” *Science* 324, no. 5924 (April 10, 2009): 218–223, doi:10.1126/science.1168978.
11. Alayne L. Brunner et al., “Distinct DNA Methylation Patterns Characterize Differentiated Human Embryonic Stem Cells and Developing Human Fetal Liver,” *Genome Research* 19, no. 6 (June 1, 2009): 1044–1056, doi:10.1101/gr.088773.108.
12. Mayumi Oda et al., “High-resolution Genome-wide Cytosine Methylation Profiling with Simultaneous Copy Number Analysis and Optimization for Limited Cell Numbers,” *Nucleic Acids Research* 37, no. 12 (July 1, 2009): 3829–3839, doi:10.1093/nar/gkp260.
13. Zachary D. Smith et al., “High-throughput Bisulfite Sequencing in Mammalian Genomes,” *Methods* 48, no. 3 (July 2009): 226–232, doi:10.1016/j.ymeth.2009.05.003.
14. Andrew M. Smith et al., “Quantitative Phenotyping via Deep Barcode Sequencing,” *Genome Research* (July 21, 2009).

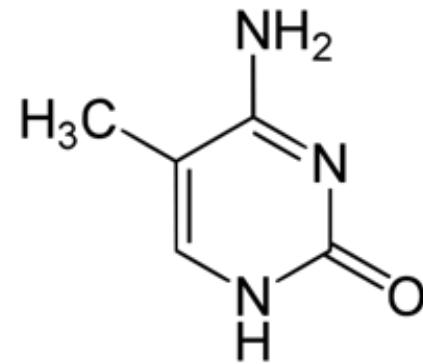


*-seq in 4 short vignettes

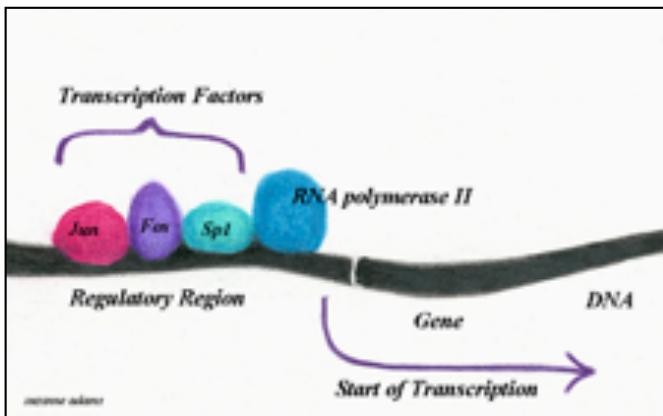
RNA-seq



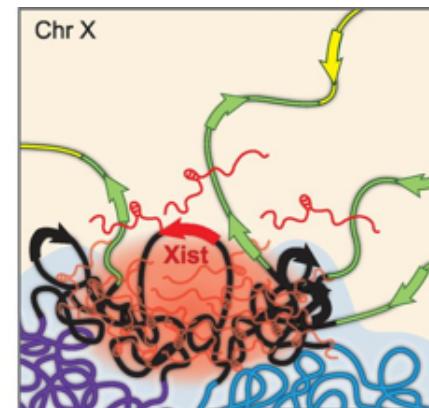
Methyl-seq



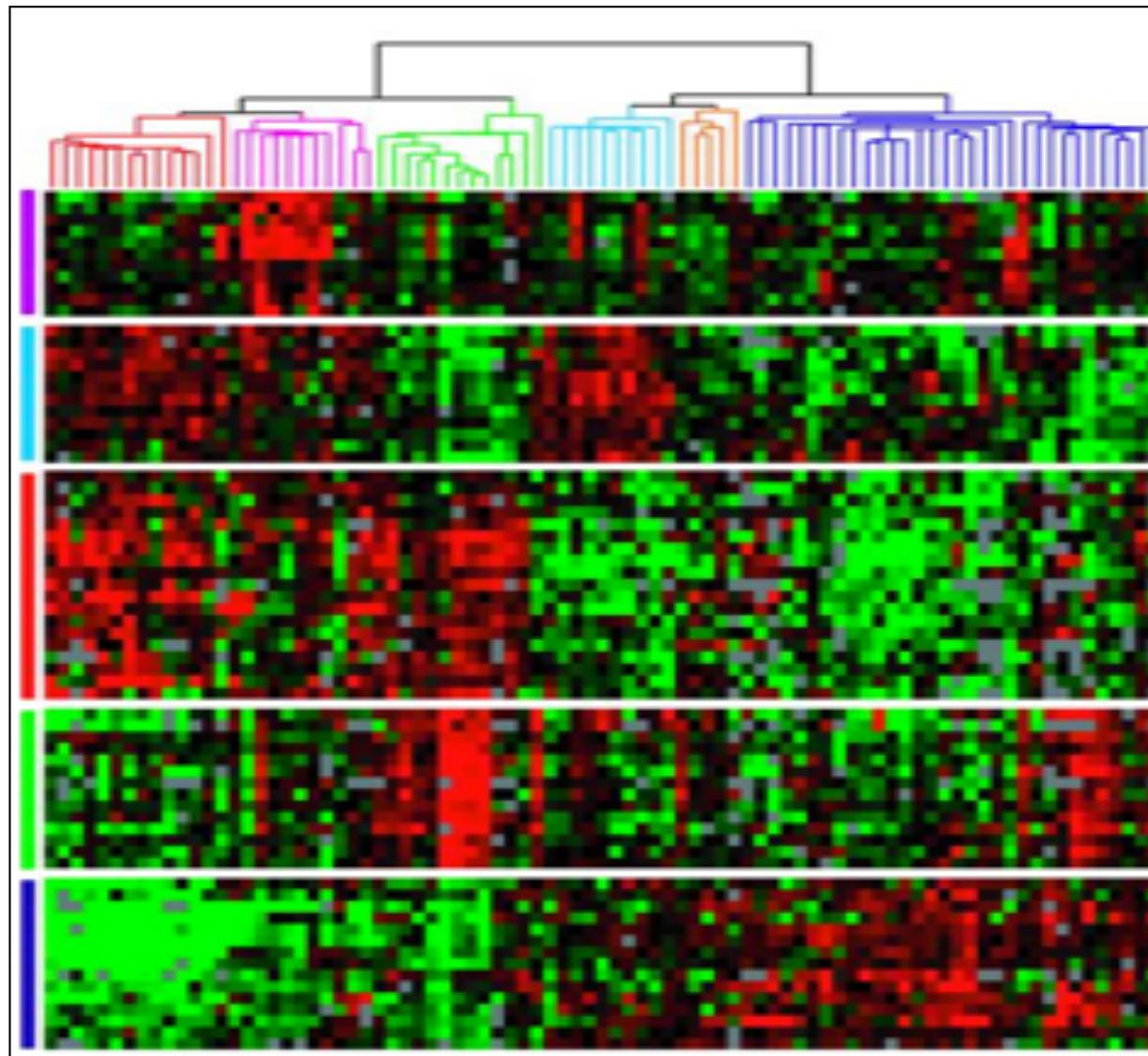
ChIP-seq



Hi-C

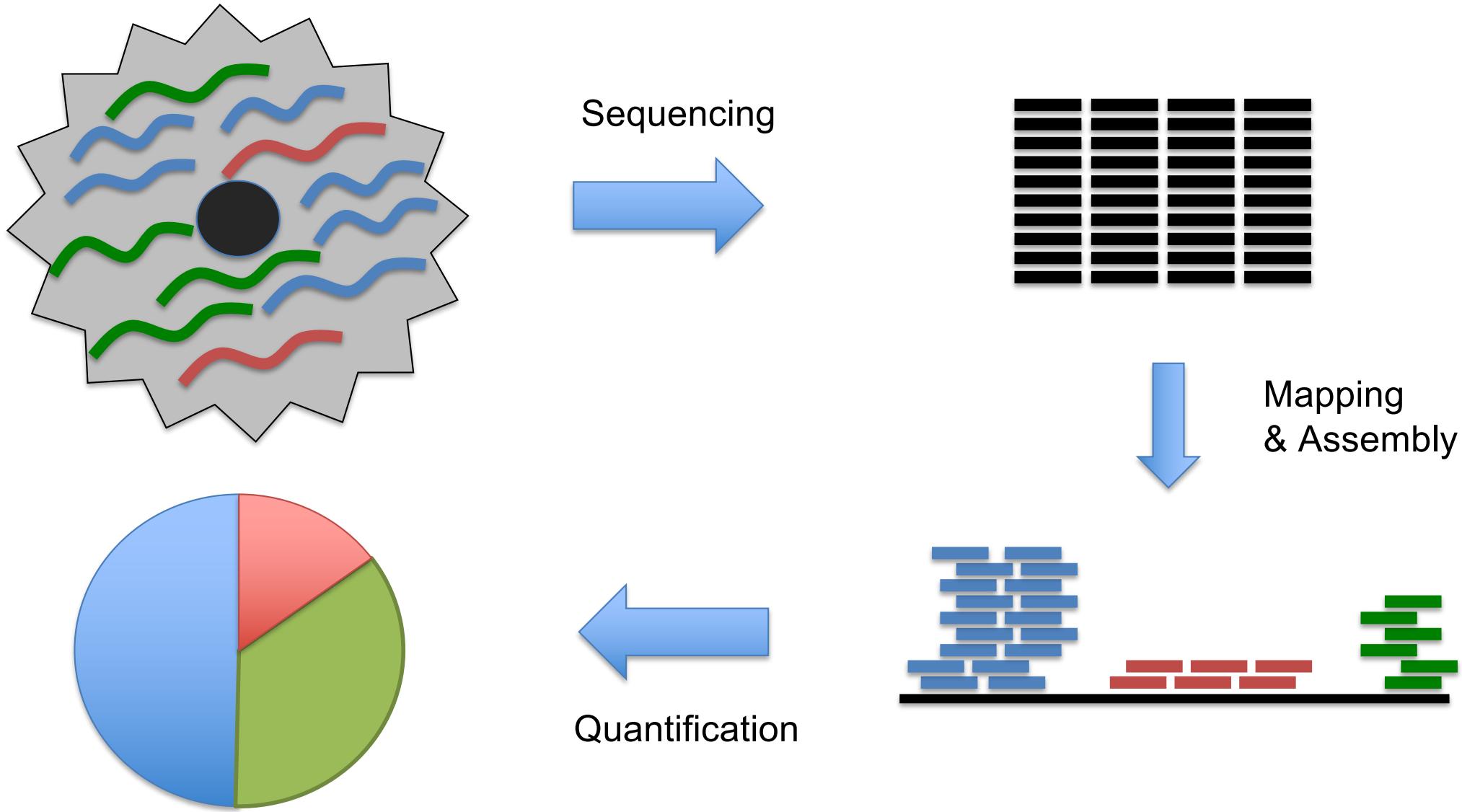


RNA-seq

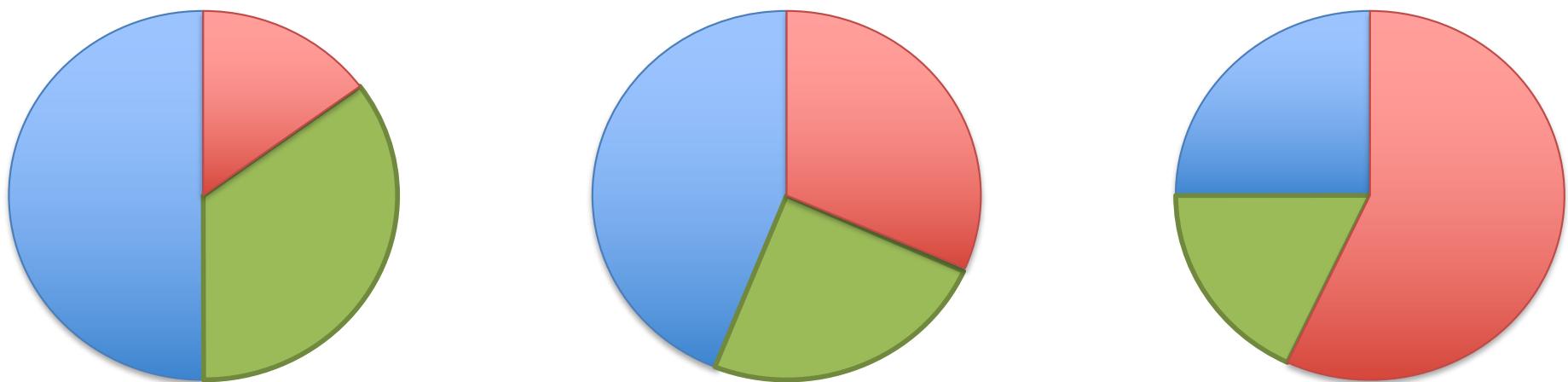
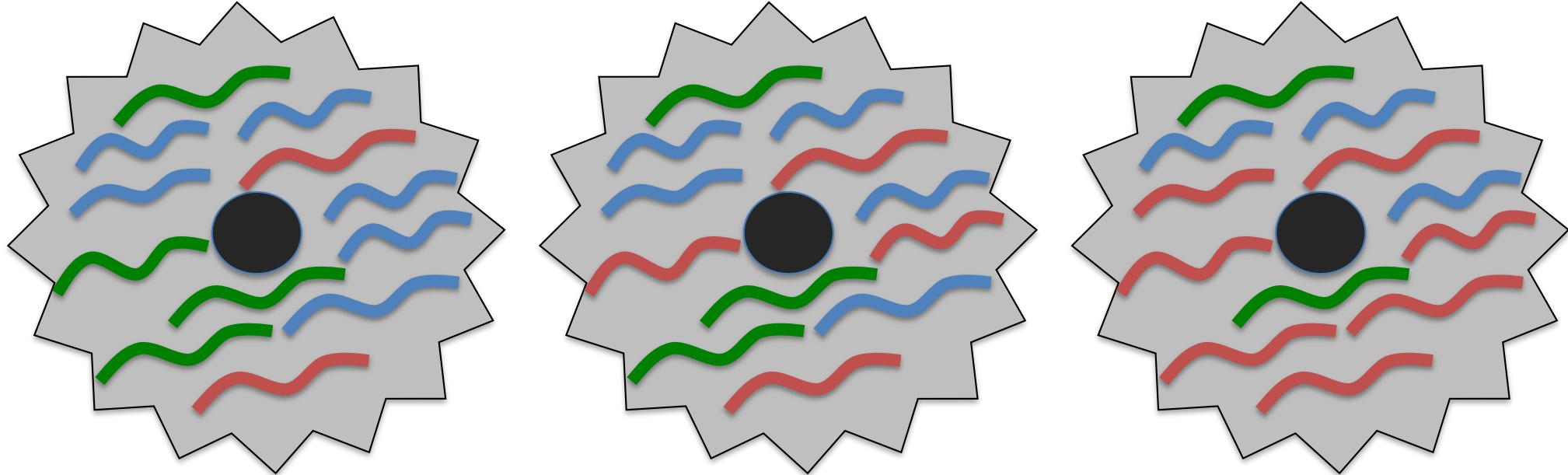


Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications.
Sørlie et al (2001) PNAS. 98(19):10869-74.

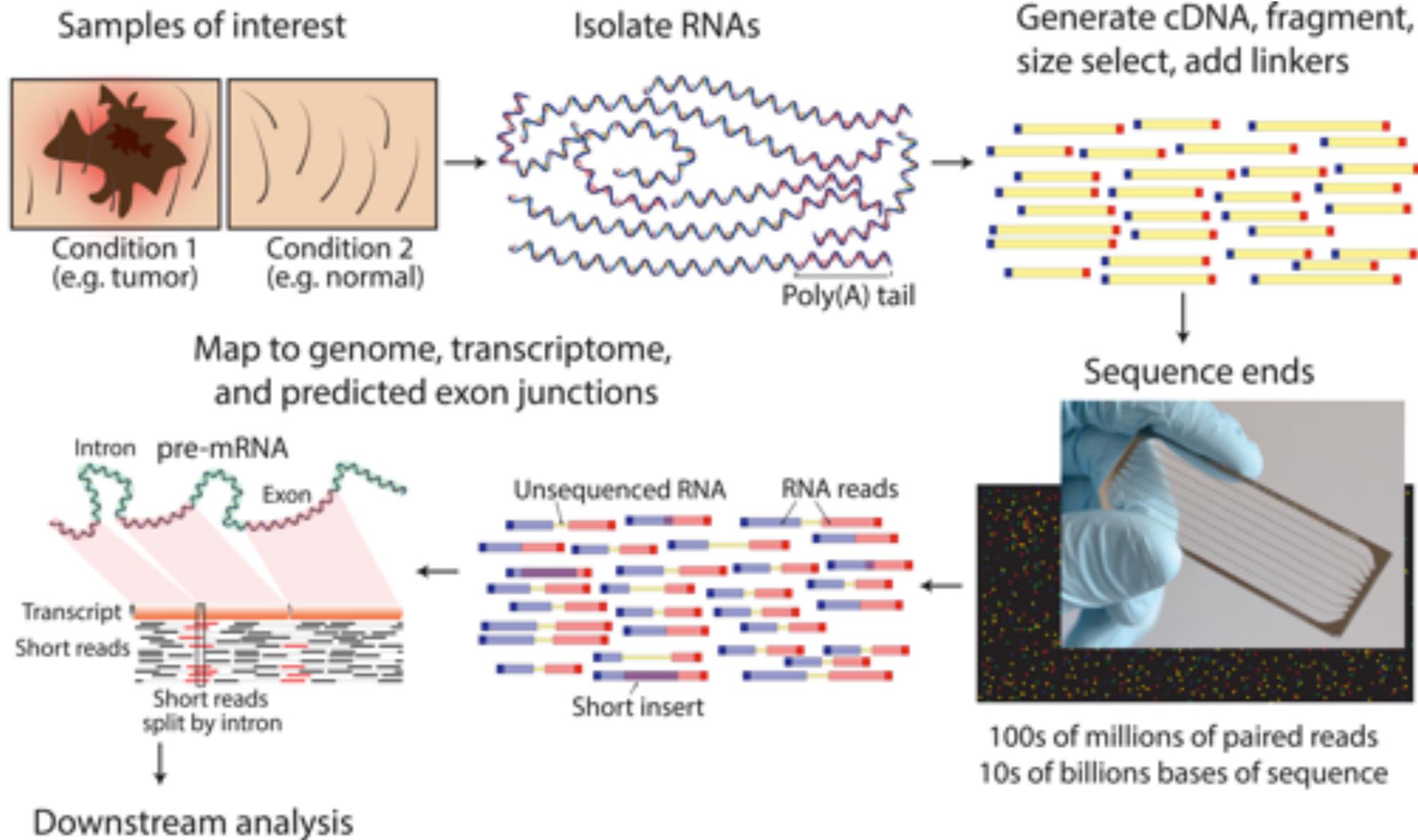
RNA-seq Overview



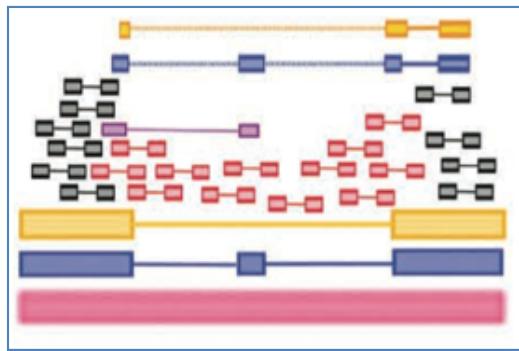
RNA-seq Overview



RNA-seq Overview



RNA-seq Challenges



Challenge I: Eukaryotic genes are spliced

RNA-Seq Approaches

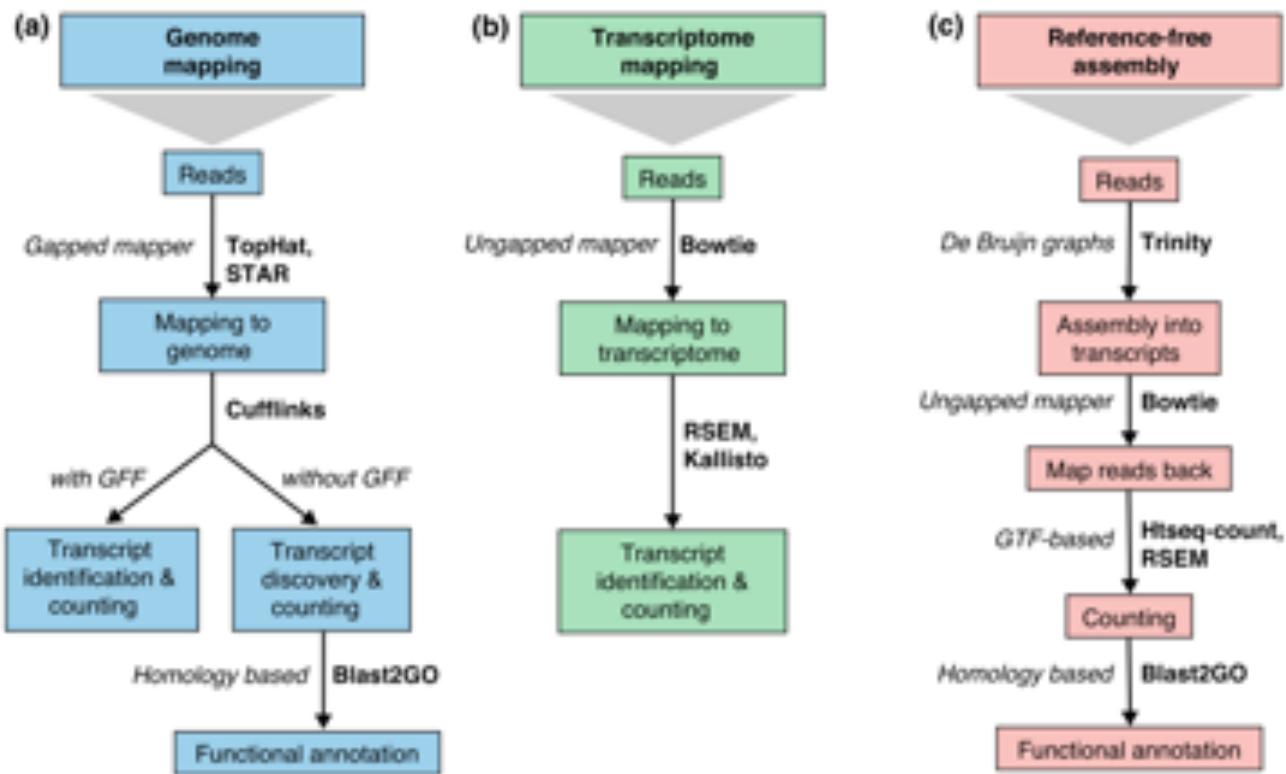


Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome with a gapped mapper. Next (novel) transcript discovery and quantification can proceed with or without an annotation file. Novel transcripts are then functionally annotated. **b** If no novel transcript discovery is needed, reads can be mapped to the reference transcriptome using an ungapped aligner. Transcript identification and quantification can occur simultaneously. **c** When no genome is available, reads need to be assembled first into contigs or transcripts. For quantification, reads are mapped back to the novel reference transcriptome and further analysis proceeds as in (b) followed by the functional annotation of the novel transcripts as in (a). Representative software that can be used at each analysis step are indicated in bold text. Abbreviations: GFF General Feature Format, GTF gene transfer format, RSEM RNA-Seq by Expectation Maximization

A survey of best practices for RNA-seq data analysis

Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

RNA-Seq Approaches

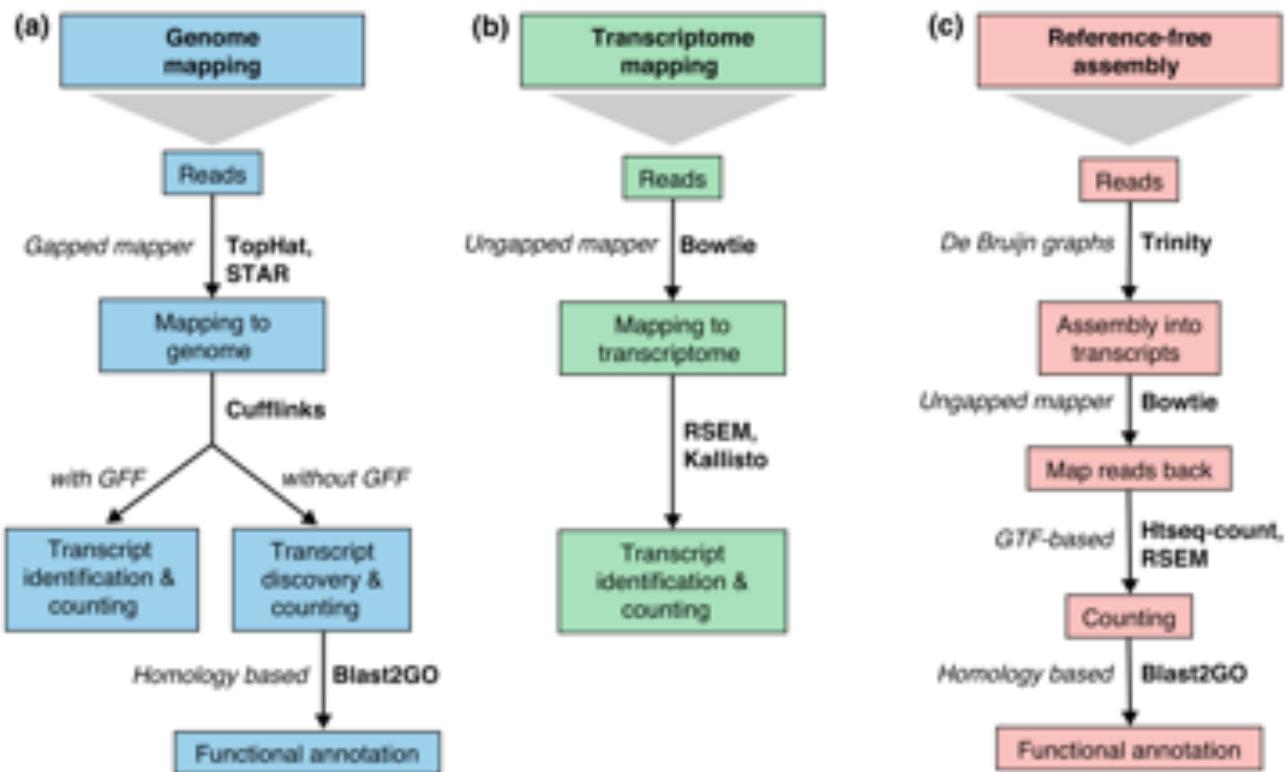


Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome using a gapped aligner (TopHat, STAR). Transcript identification and quantification can proceed with or without an annotation file (GFF). **b** No genome is available and reads are mapped to the reference transcriptome using an ungapped aligner (Bowtie). Transcript identification and quantification can occur simultaneously. **c** When no genome is available, reads need to be assembled first into contigs or transcripts. For quantification, reads are mapped back to the novel reference transcriptome and further analyzed. Functional annotation follows. Abbreviations: GFF General Feature Format, GTF gene transfer format, RSEM RNA-Seq by Expectation Maximization

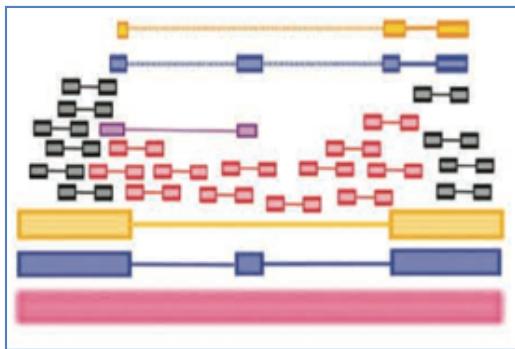
Which approach should we use?

It depends....

A survey of best practices for RNA-seq data analysis

Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

RNA-seq Challenges



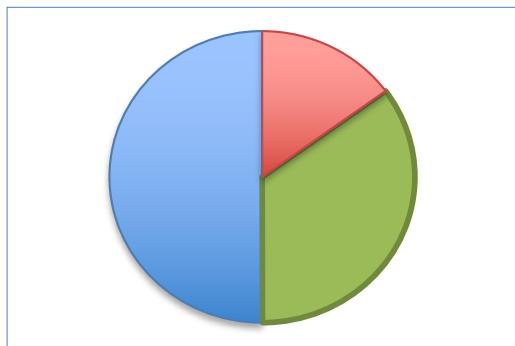
Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

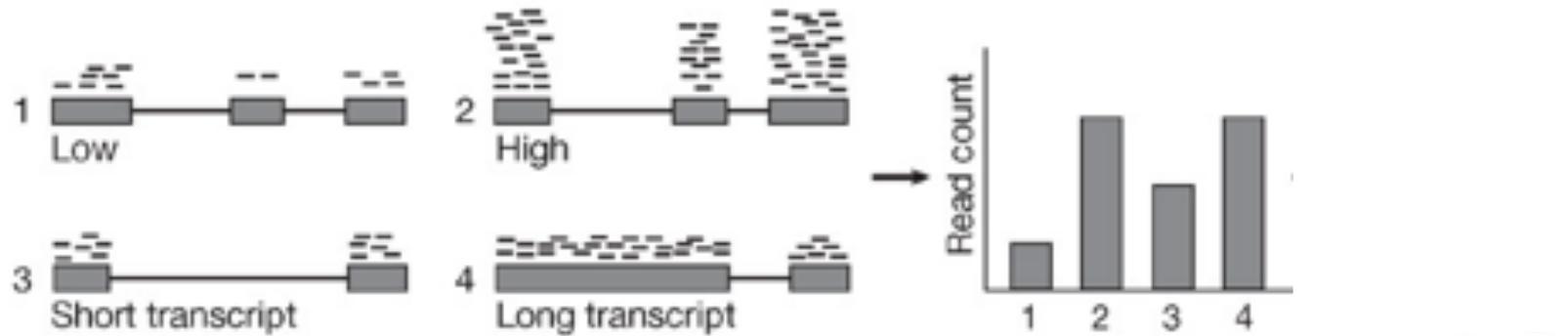
TopHat: discovering spliced junctions with RNA-Seq.

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

Challenge 2: Read Count != Transcript abundance



RPKM, FPKM, TPM

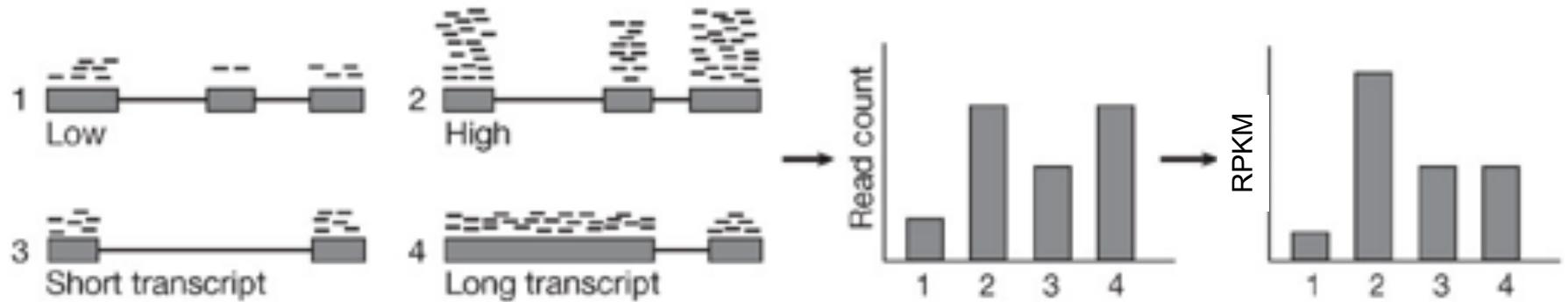


Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

RPKM, FPKM, TPM



Counting Reads that align to a gene DOESN'T work!

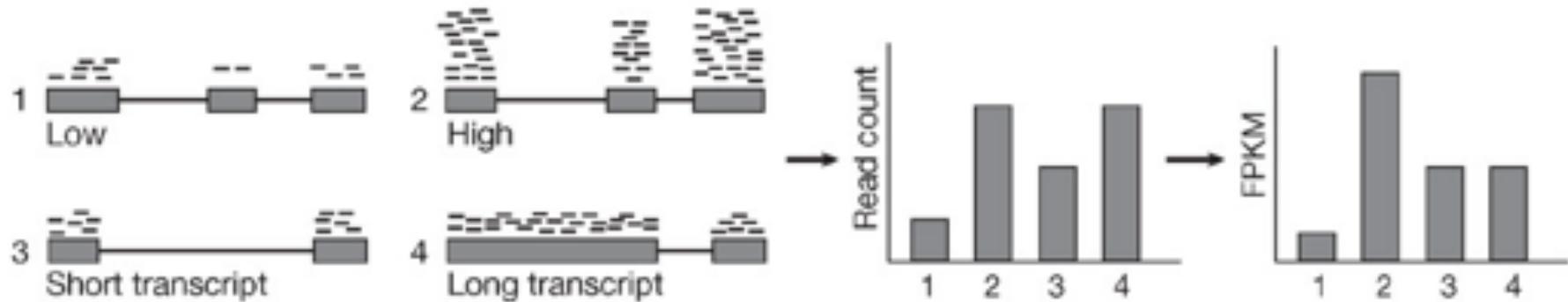
- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

(Count reads aligned to gene) / (length of gene in kilobases) / (# millions of read mapped)

=> Wait a second, reads in a pair aren't independent!

RPKM, FPKM, TPM



Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

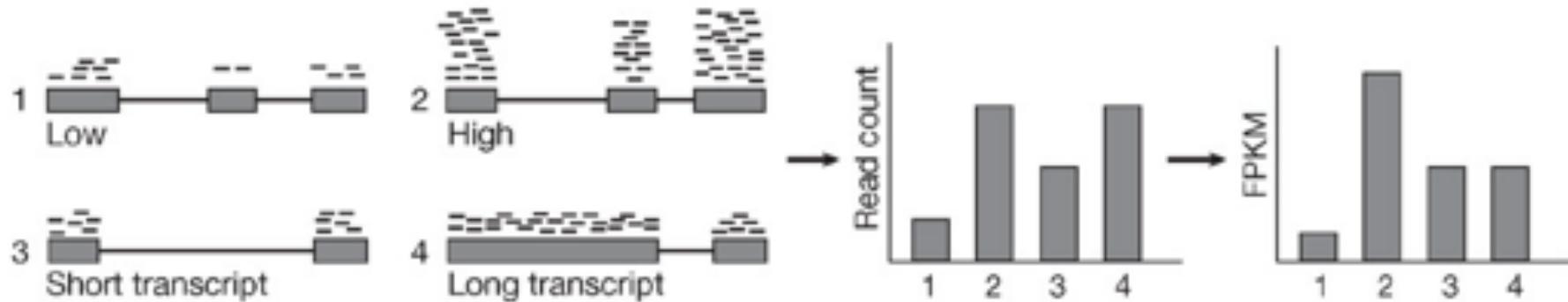
=> Wait a second, reads in a pair aren't independent!

2. FPKM: Fragments Per Kilobase of Exon Per Million Reads Mapped (Trapnell et al, 2010)

⇒ Does a much better job with short exons & short genes by boosting coverage

⇒ Wait a second, FPKM depends on the average transcript length!

RPKM, FPKM, TPM



Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

=> Wait a second, reads in a pair aren't independent!

2. FPKM: Fragments Per Kilobase of Exon Per Million Reads Mapped (Trapnell et al, 2010)

=> Wait a second, FPKM depends on the average transcript length!

3. TPM: Transcripts Per Million (Li et al, 2011)

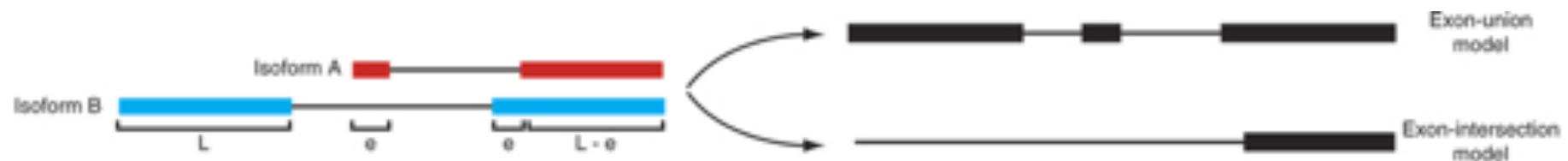
⇒ If you were to sequence one million full length transcripts, TPM is the number of transcripts you would have seen of type i, given the abundances of the other transcripts in your sample

=> Recommend you use TPM for all analysis, easy to compute given FPKM

$$TPM_i = \left(\frac{FPKM_i}{\sum_j FPKM_j} \right) \cdot 10^6$$

Gene or Isoform Quantification?

a



Differential analysis of gene regulation at transcript resolution with RNA-seq
Trapnell et al (2013) Nature Biotechnology 31, 46–53. doi:10.1038/nbt.2450

Gene or Isoform Quantification?

a



b

Condition A	Condition B	Log fold-change (union count)	Log fold-change (intersect count)	Log fold-change (true expression)
		$\log_2\left(\frac{10}{10}\right) = 0$	$\log_2\left(\frac{8}{7}\right) = 0.19$	$\log_2\left(\frac{\frac{10}{L}}{\frac{6}{L} + \frac{4}{2L}}\right) = 0.32$

Differential analysis of gene regulation at transcript resolution with RNA-seq
Trapnell et al (2013) Nature Biotechnology 31, 46–53. doi:10.1038/nbt.2450

Gene or Isoform Quantification?

a

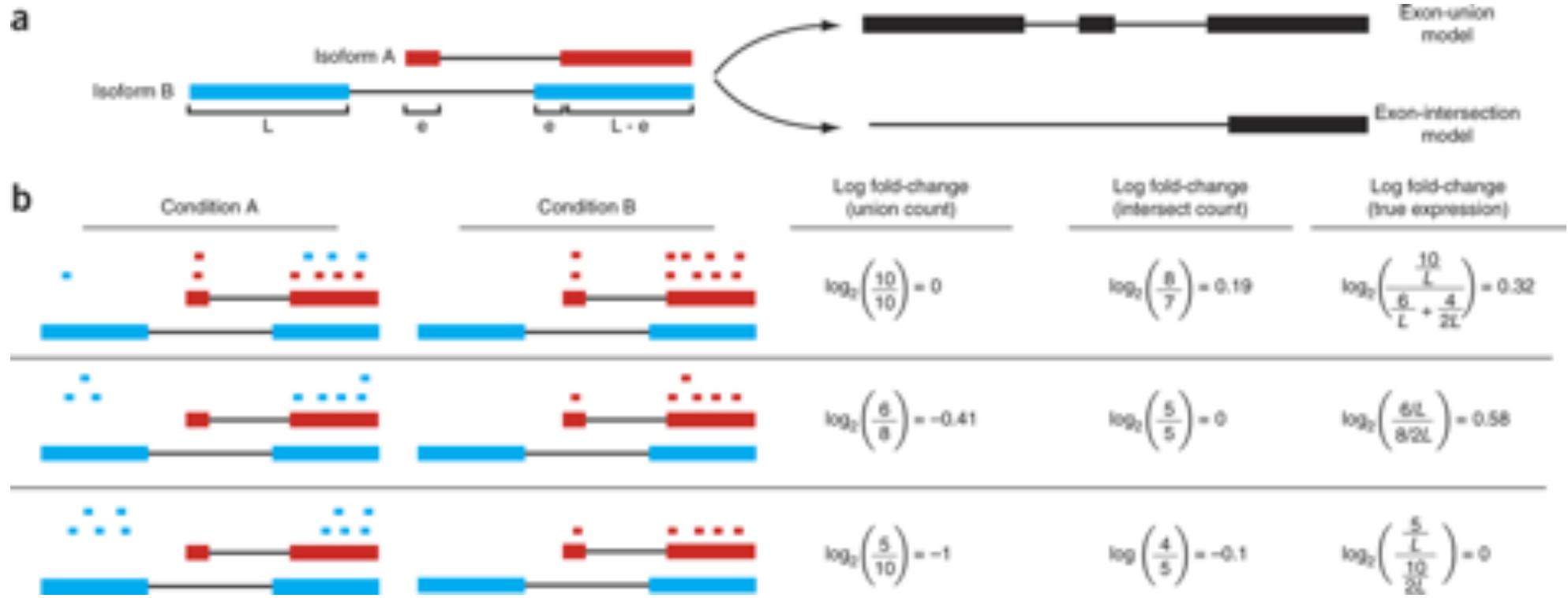


b

Condition A	Condition B	Log fold-change (union count)	Log fold-change (intersect count)	Log fold-change (true expression)
		$\log_2\left(\frac{10}{7}\right) = 0$	$\log_2\left(\frac{7}{7}\right) = 0.19$	$\log_2\left(\frac{10}{\frac{1}{L} + \frac{4}{2L}}\right) = 0.32$
		$\log_2\left(\frac{8}{5}\right) = -0.41$	$\log_2\left(\frac{5}{5}\right) = 0$	$\log_2\left(\frac{8/L}{8/2L}\right) = 0.58$

Differential analysis of gene regulation at transcript resolution with RNA-seq
Trapnell et al (2013) Nature Biotechnology 31, 46–53. doi:10.1038/nbt.2450

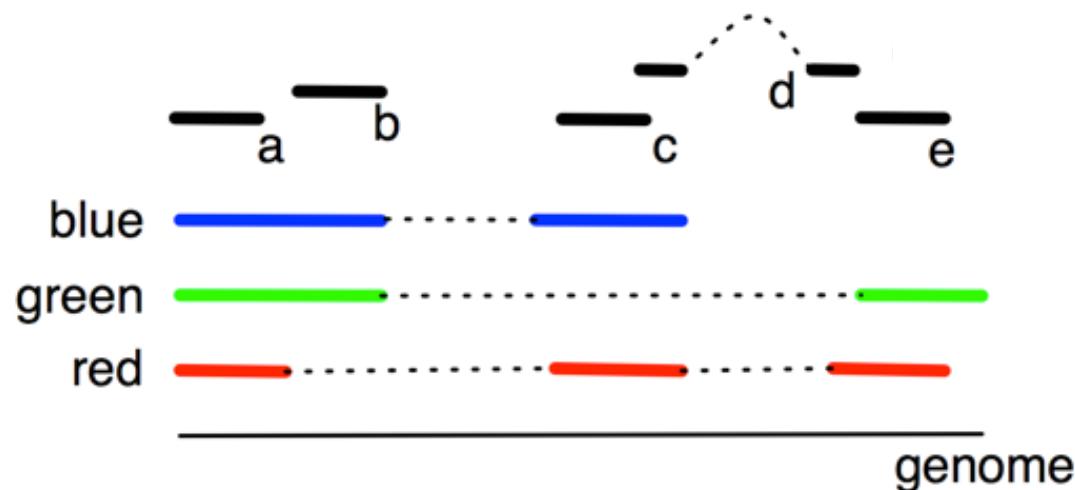
Gene or Isoform Quantification?



Key point : The length of the actual molecule from which the fragments derive is crucially important to obtaining accurate abundance estimates.

Differential analysis of gene regulation at transcript resolution with RNA-seq
Trapnell et al (2013) Nature Biotechnology 31, 46–53. doi:10.1038/nbt.2450

Multi-mapping? Isoform ambiguity? Expectation Maximization to the Rescue



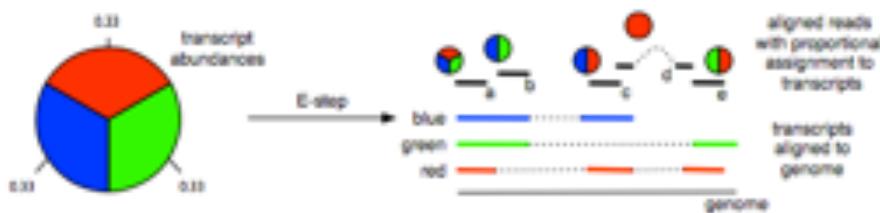
The gene has three isoforms (red, green, blue) of the same length.
Our initial expectation is all 3 isoforms are equally expressed

There are five reads (a,b,c,d,e) mapping to the gene.

- Read a maps to all three isoforms
- Read d only to red
- Reads b,c,e map to each of the three pairs of isoforms.

What is the most likely expression level of each isoform?

Multi-mapping? Isoform ambiguity? Expectation Maximization to the Rescue

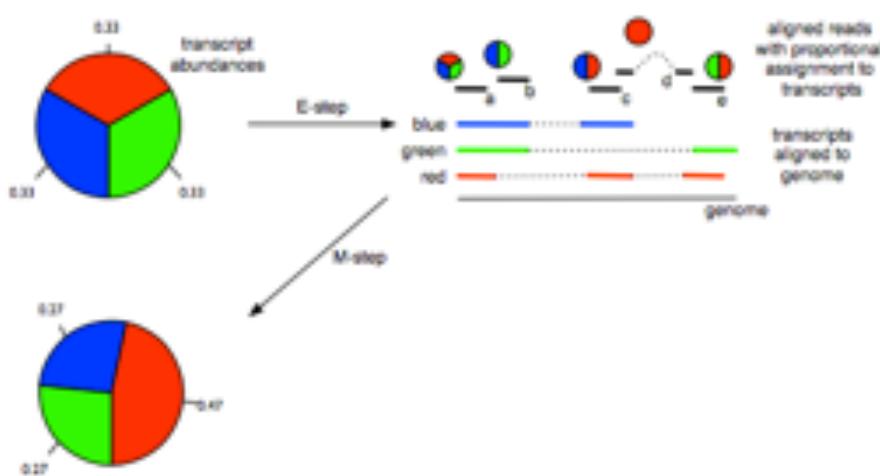


The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)

Multi-mapping? Isoform ambiguity? Expectation Maximization to the Rescue



The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)

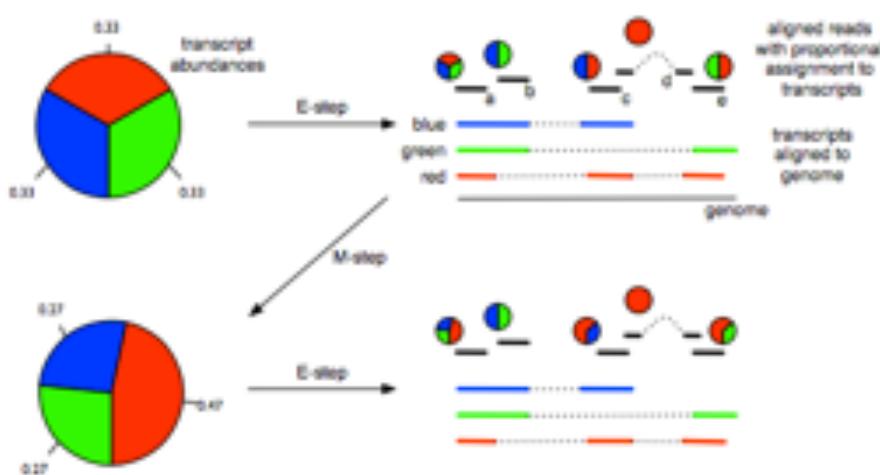
Next, during the maximization (M) step isoform abundances are recalculated from the proportionately assigned read counts:

$$\text{red: } 0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)$$

$$\text{blue: } 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$$

$$\text{green: } 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$$

Multi-mapping? Isoform ambiguity? Expectation Maximization to the Rescue



The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)

Next, during the maximization (M) step isoform abundances are recalculated from the proportionately assigned read counts:

$$\text{red: } 0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)$$

$$\text{blue: } 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$$

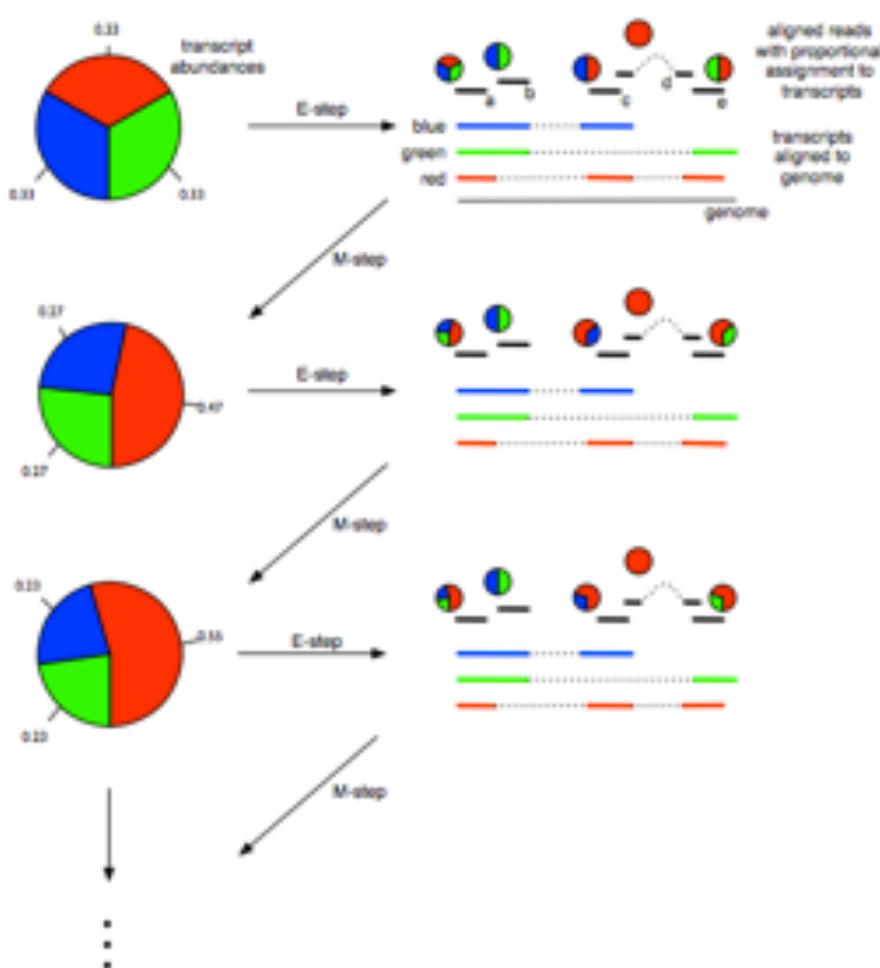
$$\text{green: } 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$$

Repeat until convergence!

Models for transcript quantification from RNA-seq

Pachter, L (2011) arXiv. 1104.3889 [q-bio.GN]

Multi-mapping? Isoform ambiguity? Expectation Maximization to the Rescue



The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)

Next, during the maximization (M) step isoform abundances are recalculated from the proportionately assigned read counts:

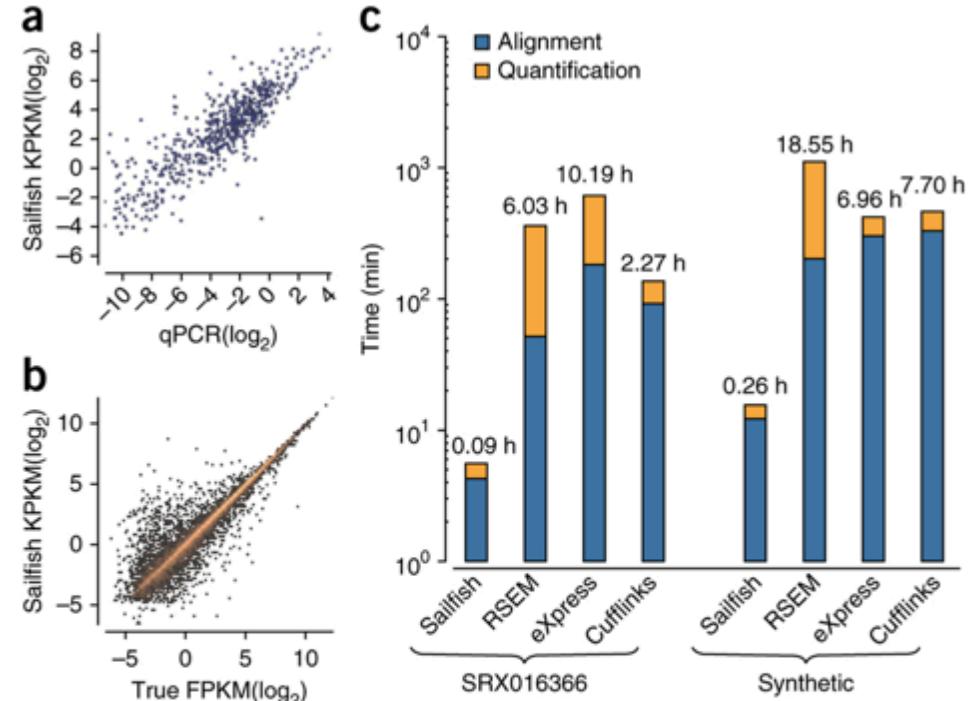
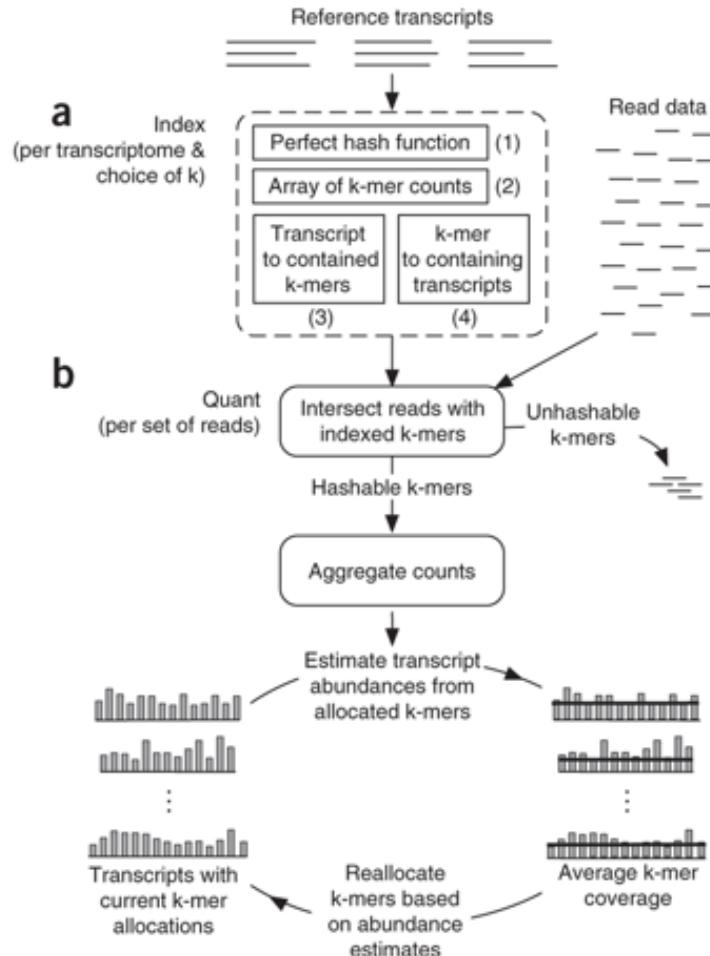
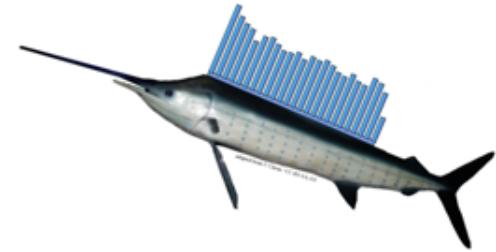
$$\text{red: } 0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)$$

$$\text{blue: } 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$$

$$\text{green: } 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$$

Repeat until convergence!

Sailfish: Fast & Accurate RNA-seq Quantification



Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms
 Patro et al (2014) Nature Biotechnology 32, 462–464 doi:10.1038/nbt.2862

Annotation Summary

- Three major approaches to annotate a genome
 - I. Alignment:
 - Does this sequence align to any other sequences of known function?
 - Great for projecting knowledge from one species to another
 - 2. Prediction:
 - Does this sequence statistically resemble other known sequences?
 - Potentially most flexible but dependent on good training data
 - 3. Experimental:
 - Lets test to see if it is transcribed/methylated/bound/etc
 - Strongest but expensive and context dependent
- Many great resources available
 - Learn to love the literature and the databases
 - Standard formats let you rapidly query and cross reference
 - Google is your number one resource ☺

