# Lecture 13. Gene Finding & RNAseq

Michael Schatz
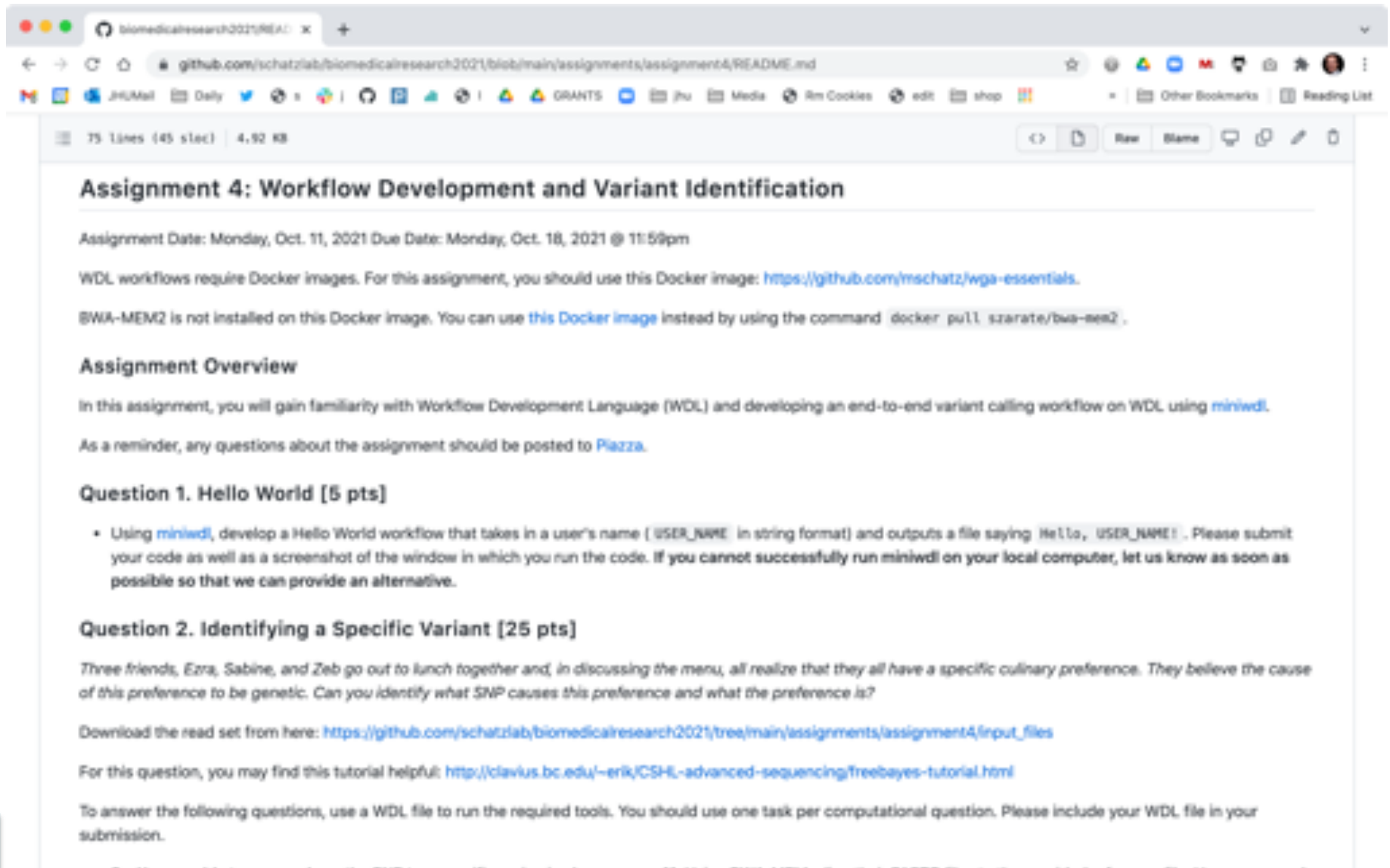
October 13, 2021
Advanced Biomedical Research

# Assignment 4: WDLs
# Due Oct 18 @ 11:59pm



https://github.com/schatzlab/biomedicalresearch2021

# Goal: Genome Annotations

aatgcatgcggctatgctaatgcatgcggctatgctaagctgggatccgatgacaatgcatgcggctatgctaat
gcatgcggctatgcaagctgggatccgatgactatgctaagctgggatccgatgacaatgcatgcggctatgct
aatgaatggtcttgggatttaccttggaatgctaagctgggatccgatgacaatgcatgcggctatgctaatgaa
tggtcttgggatttaccttggaatatgctaatgcatgcggctatgctaagctgggatccgatgacaatgcatgcg
gctatgctaatgcatgcggctatgcaagctgggatccgatgactatgctaagctgcggctatgctaatgcatgcg
gctatgctaagctgggatccgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctgggatcct
gcggctatgctaatgaatggtcttgggatttaccttggaatgctaagctgggatccgatgacaatgcatgcggct
atgctaatgaatggtcttgggatttaccttggaatatgctaatgcatgcggctatgctaagctgggaatgcatgcg
gctatgctaagctgggatccgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctgggatccg
atgactatgctaagctgcggctatgctaatgcatgcggctatgctaagctcatgcggctatgctaagctgggaat
gcatgcggctatgctaagctgggatccgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctg
ggatccgatgactatgctaagctgcggctatgctaatgcatgcggctatgctaagctcggctatgctaatgaatg
gtcttgggatttaccttggaatgctaagctgggatccgatgacaatgcatgcggctatgctaatgaatggtcttgg
gatttaccttggaatatgctaatgcatgcggctatgctaagctgggaatgcatgcggctatgctaagctgggatc
cgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctgggatccgatgactatgctaagctgcg
gctatgctaatgcatgcggctatgctaagctcatgcgg

# Goal: Genome Annotations

aatgcatgcggctatgctaatgcatgcggctatgctaagctgggatccgatgacaatgcatgcggctatgctaat
gcatgcggctatgcaagctgggatccgatgactatgctaagctgggatccgatgacaatgcatgcggctatgct
aatgaatggtcttgggatttaccttggaatgctaagctgggatccgatgacaatgcatgcggctatgctaatgaa
tggtcttgggatttaccttggaatatgctaatgcatgcggctatgctaagctgggatccgatgacaatgcatgcg
gctatgctaatgcatgcggctatgcaagctgggatccgatgactatgctaagctgcggctatgctaatgcatgcg
gctat<span style="color:red">gctaagctgggatccgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctgggatcct</span>
<span style="color:red">gcggctatgctaatgaatggtcttgggatttaccttggaatgctaagctggatccgatgacaatgcatgcggct</span>
<span style="color:red">atgctaatgaatggtcttgggatt</span> **Gene!** <span style="color:red">ctatgctaagctgggaatgcatgcg</span>
<span style="color:red">gctatgctaagctgggatccgat</span> <span style="color:red">atgcggctatgcaagctgggatccg</span>
<span style="color:red">atgactatgctaagctgcggctatgctaatgcatgcggctatgctaagctcatgcggctatgctaagctgggaat</span>
<span style="color:red">gcatgcggctatgctaa</span>gctgggatccgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctg
ggatccgatgactatgctaagctgcggctatgctaatgcatgcggctatgctaagctcggctatgctaatgaatg
gtcttgggatttaccttggaatgctaagctgggatccgatgacaatgcatgcggctatgctaatgaatggtcttgg
gatttaccttggaatatgctaatgcatgcggctatgctaagctgggaatgcatgcggctatgctaagctgggatc
cgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctgggatccgatgactatgctaagctgcg
gctatgctaatgcatgcggctatgctaagctcatgcgg

# BEDTools to the rescue!

# 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 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

(Altschul et al. 1990)

# 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

```
FAKDFLAGGVAAAI SKTAVAPIERVKLLLQVQ HASKQITADKQYKGIIDCVVRIPKEQGV
|     |      |  |||  ||||||||||||||||||  |||  |    |    ||||  |      |  ||||||
FLIDLASGGTAAAV SKTAVAPIERVKLLLQVQ 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 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, $\lambda$ 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--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV 60


Query  56   KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPA 115
            K HGKKV  A ++ +AH+D++     + LS+LH  KL VDP NF+LL + L+   LA H
Sbjct  61   KAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGK 120


Query  116  EFTPAVHASLDKFLASVSTVLTSKY 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    LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF------DLSHGSAQV  55
            LS  +   V   WGKV A     +G E L R+F   P T   F  F      D   S  +
Sbjct  3    LSDGEWQLVLNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKFKHLKSEDEMKASEDL  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  ERMFLSFPTTKTYFPHFDLSHGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAH  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



- 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**

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

Start:
- AUG

Stop:
- UAA
- UAG
- UGA

# Step One

- Find open reading frames (ORFs).

# Step One

- Find open reading frames (ORFs).

Stop codon

Reverse strand

...ATCTACTTACCGAGA**AAT**CTATTTAAAGTACTTTTTATAACT...

...**TAG**ATGAATGGCTCTTTAGA**TAA**ATTTCATGAAAAATAT**TGA**...

Stop codon

Shifted Stop

Stop codon

- 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

# 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

# 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.

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*:
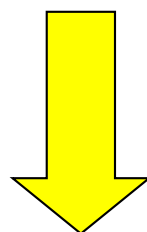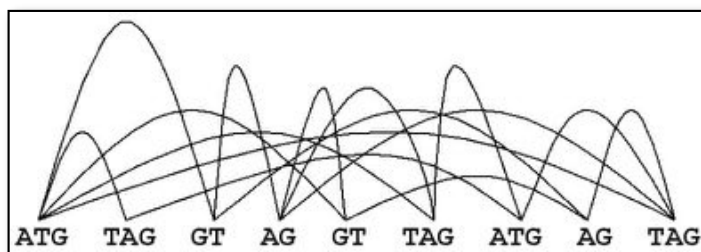


An ORF graph represents all possible *gene parses* (and their scores) for a given set of putative signals. A *path* through the graph represents a single gene parse.

# Conceptual Gene-finding Framework

TATTCCGATCGATCGATCTCTCTAGCGTCTACG
CTATCATCGCTCTCTATTATCGCGCGATCGTCG
ATCGCGCGAGAGTATGCTACGTCGATCGAATTG

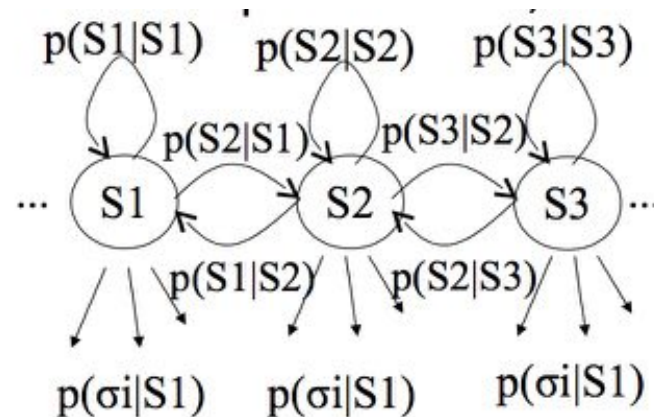identify most promising signals, score signals and content regions between them; induce an ORF graph on the signals



ATG  TAG  GT  AG  GT  TAG  ATG  AG  TAG

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

Duke
UNIVERSITY

# HMMs for Gene Finding

- 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.
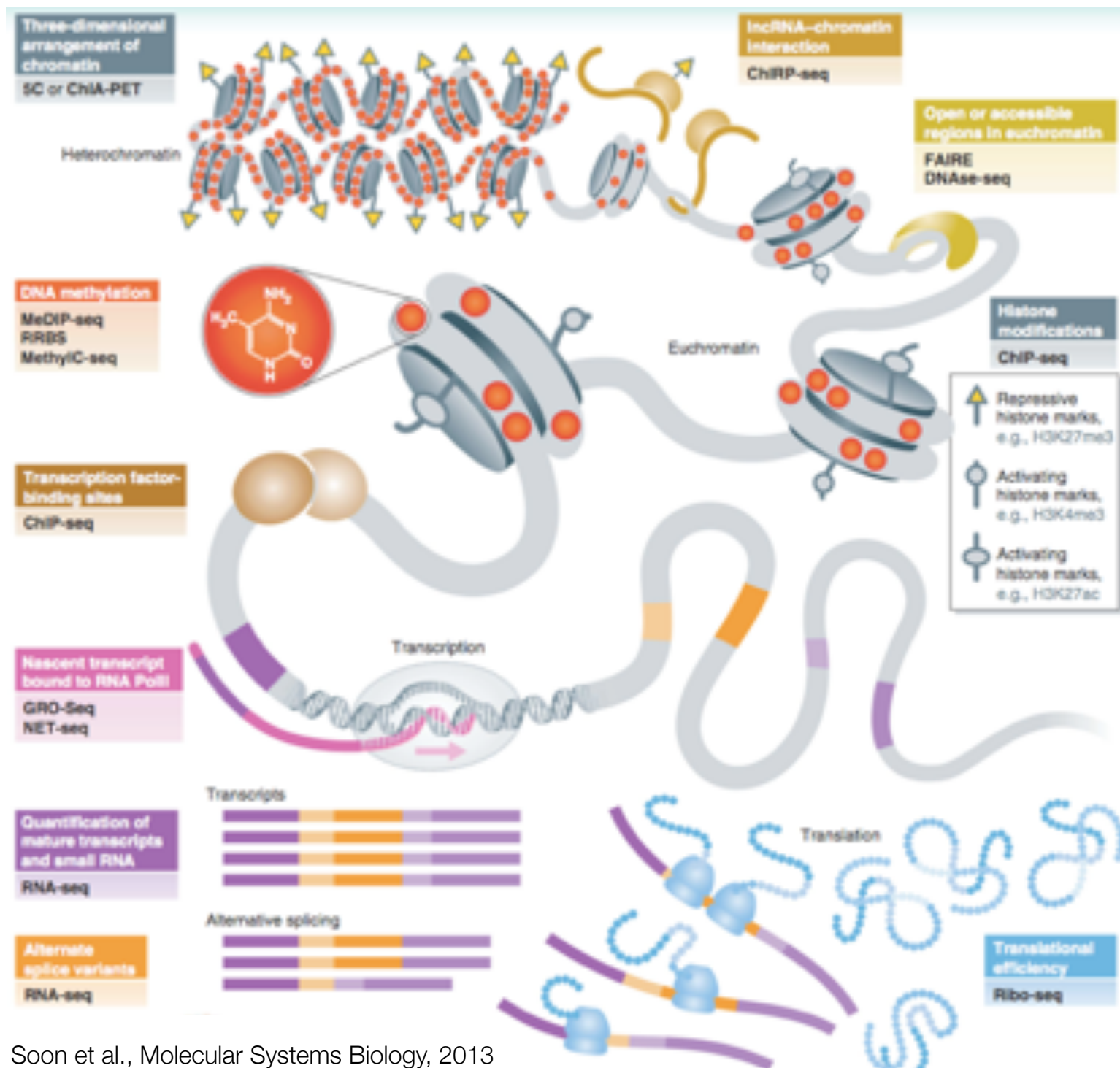


AAAGCATGCATTTAACGTGAGCACAATAGATTACA

# Outline

1. Alignment to other genomes

2. Prediction aka "Gene Finding"

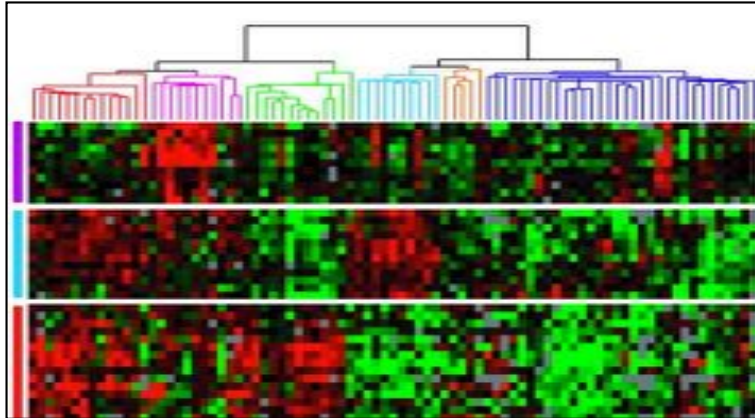3. **Experimental & Functional Assays**

# 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.

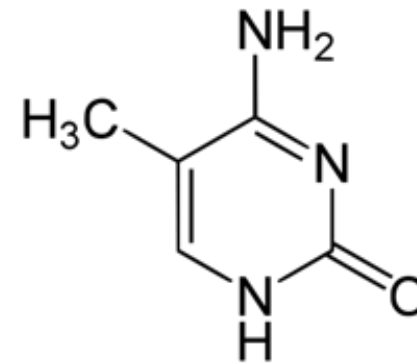14. Andrew M. Smith et al., "Quantitative Phenotyping via Deep Barcode Sequencing," Genome Research (July 21, 2009).
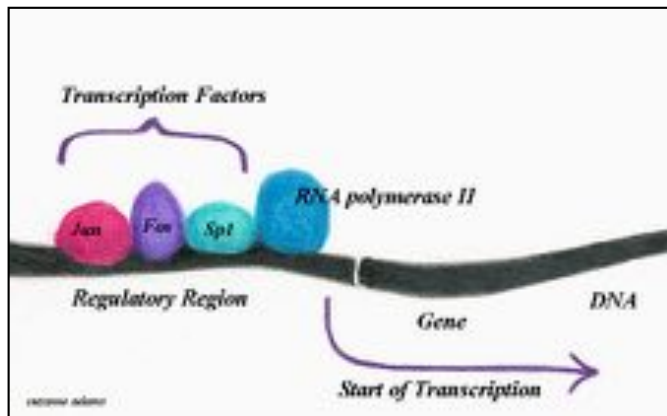
Soon et al., Molecular Systems Biology, 2013
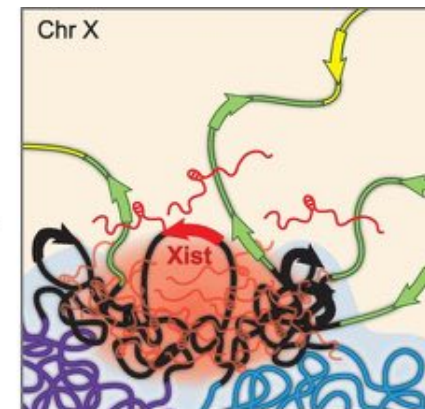
# *-seq in 4 short vignettes

## RNA-seq



## Methyl-seq



## ChIP-seq



## Hi-C

# RNA-seq



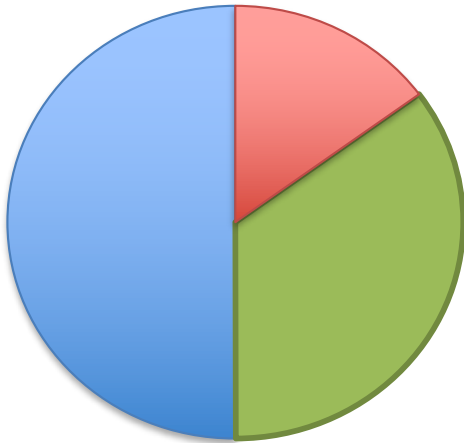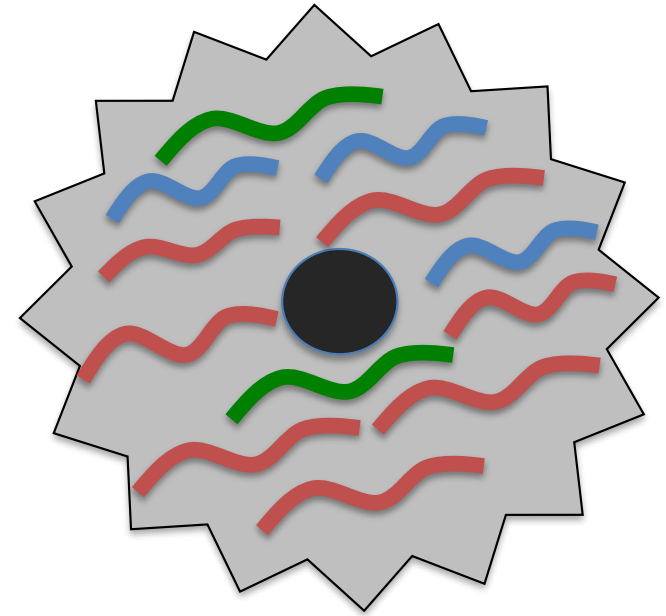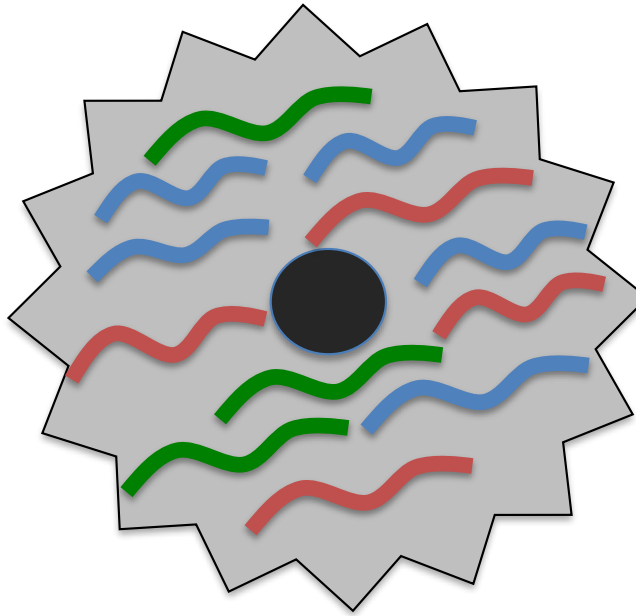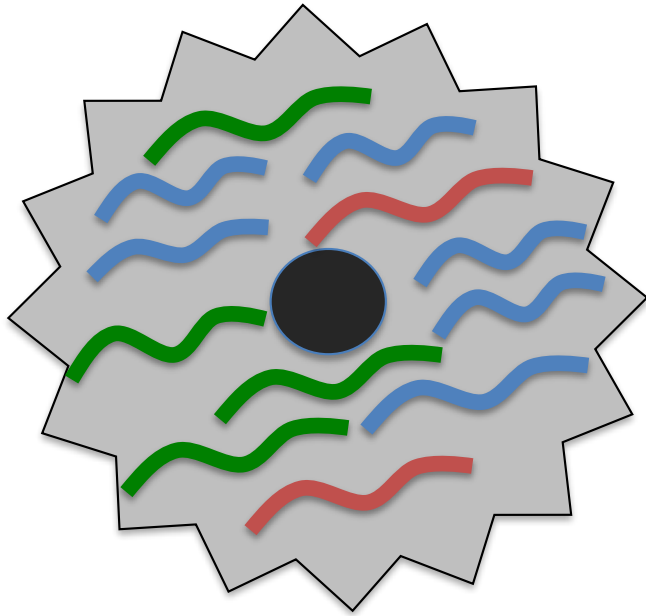**Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications.**
Sørlie et al (2001) *PNAS*. 98(19):10869-74.

# RNA-seq Overview

# RNA-seq Overview

# RNA-seq Overview

**Samples of interest**

Condition 1
(e.g. tumor)
Condition 2
(e.g. normal)

**Isolate RNAs**

Poly(A) tail

**Generate cDNA, fragment, size select, add linkers**

**Sequence ends**

100s of millions of paired reads
10s of billions bases of sequence

**Map to genome, transcriptome, and predicted exon junctions**

Intron   pre-mRNA

Exon

Unsequenced RNA    RNA reads

Transcript

Short reads

Short reads split by intron

Short insert

**Downstream analysis**

# RNA-seq Challenges



***Challenge 1: Eukaryotic genes are spliced***
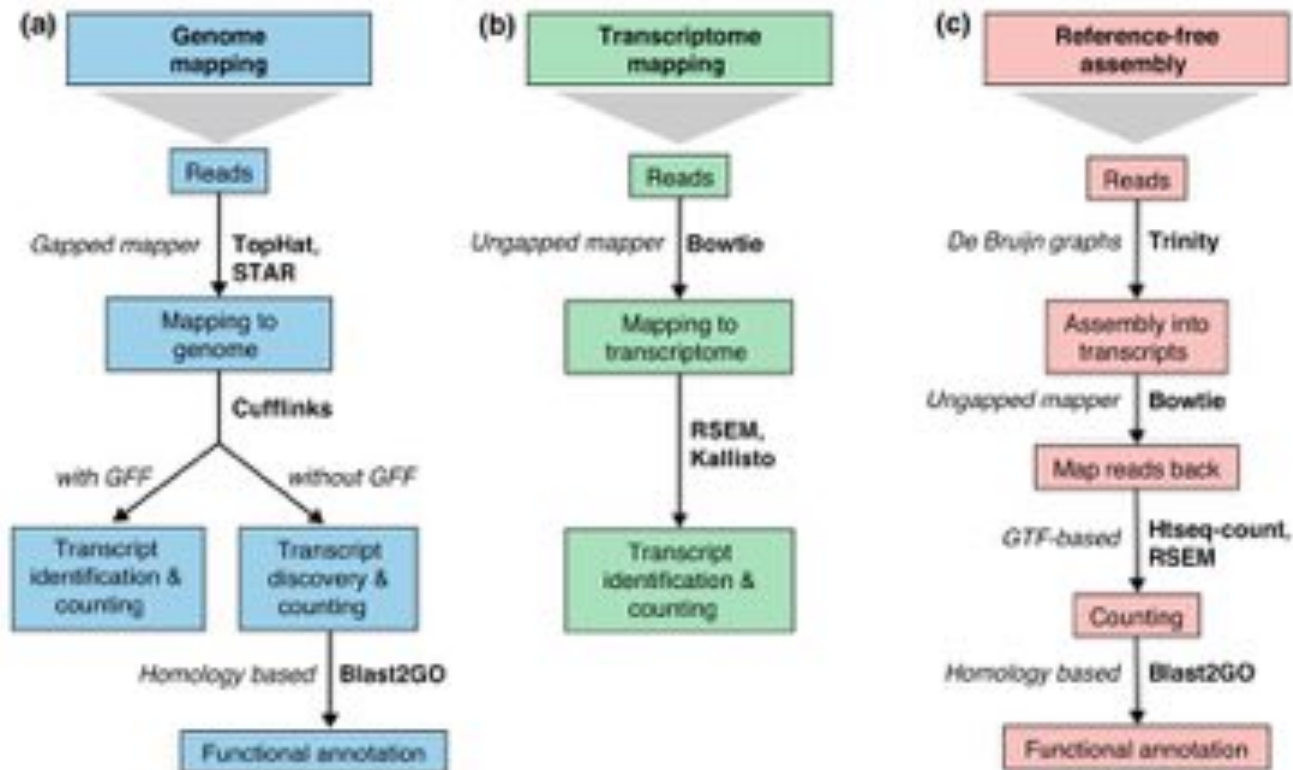
# RNA-Seq Approaches



**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
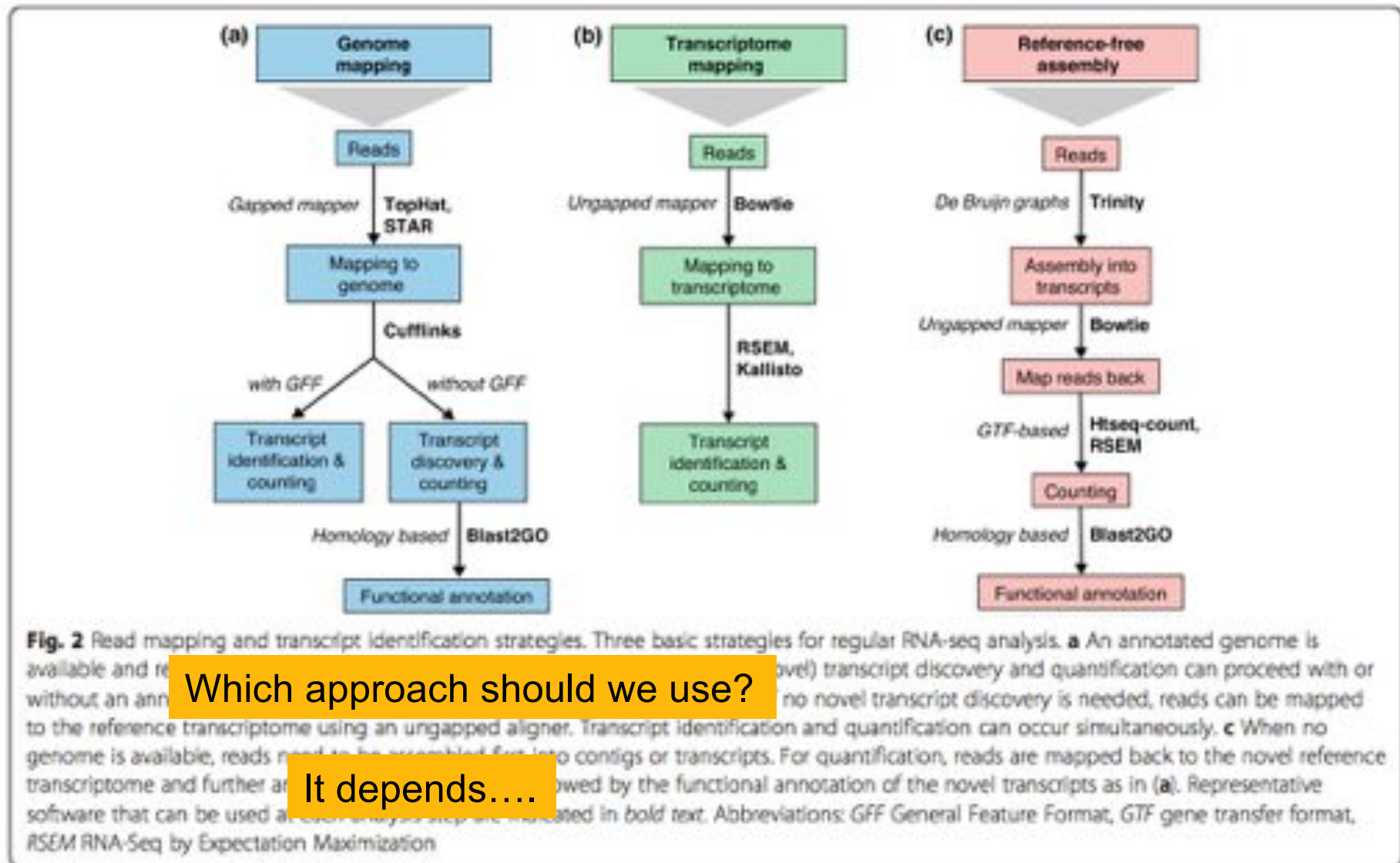
# RNA-Seq Approaches



**Fig. 2** Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and re[...] (novel) transcript discovery and quantification can proceed with or without an ann[...] 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 an[...]owed by the functional annotation of the novel transcripts as in (a). Representative software that can be used a[...] 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*
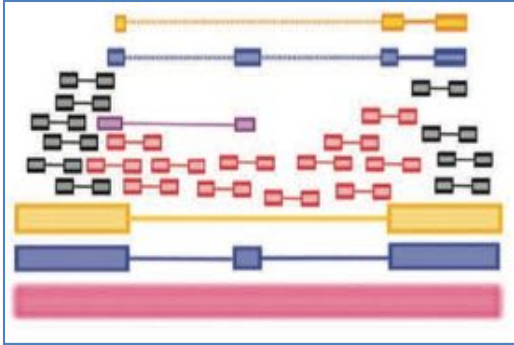Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8
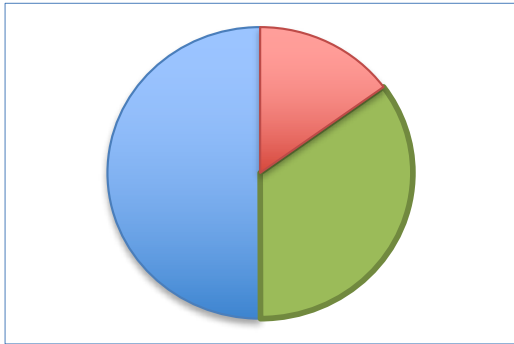
# RNA-seq Challenges



**_Challenge 1: Eukaryotic genes are spliced_**

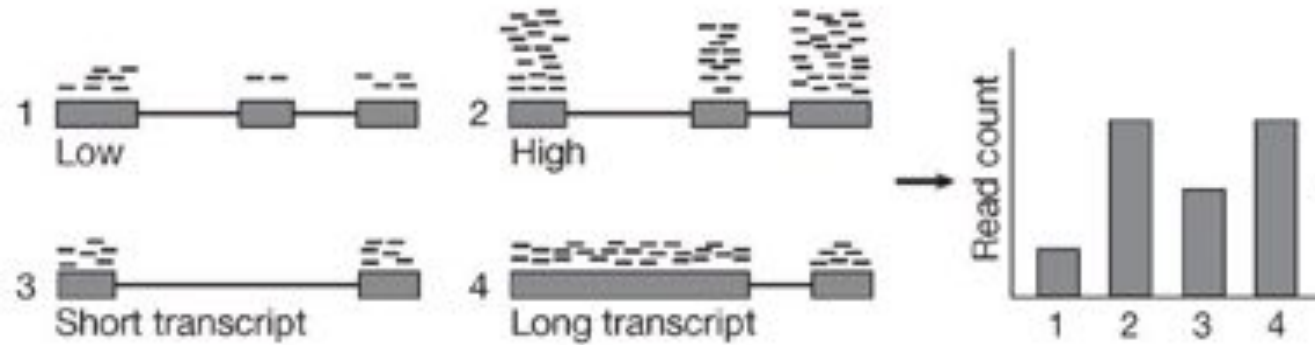Solution: Use a spliced aligner, and assemble isoforms

**_TopHat: discovering spliced junctions with RNA-Seq._**
Trapnell et al (2009) _Bioinformatics_. 25:0 1105-1111



**_Challenge 2: Read Count != Transcript abundance_**
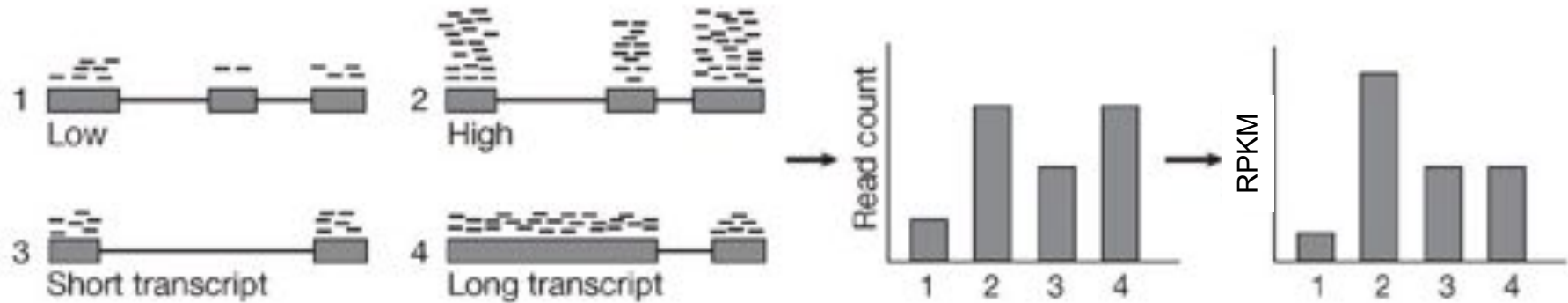
# RPKM, FPKM, TPM



***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



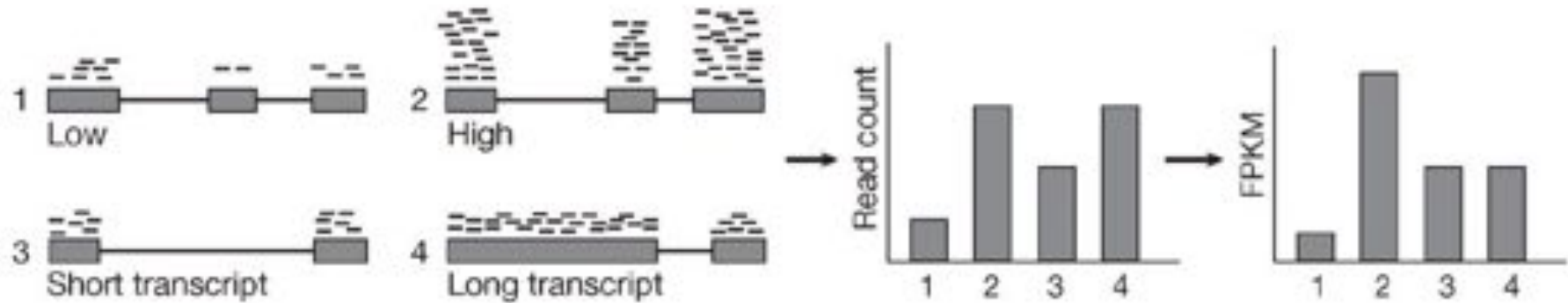***Counting Reads that align to a gene DOESN'T work!***
- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

***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 arent 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

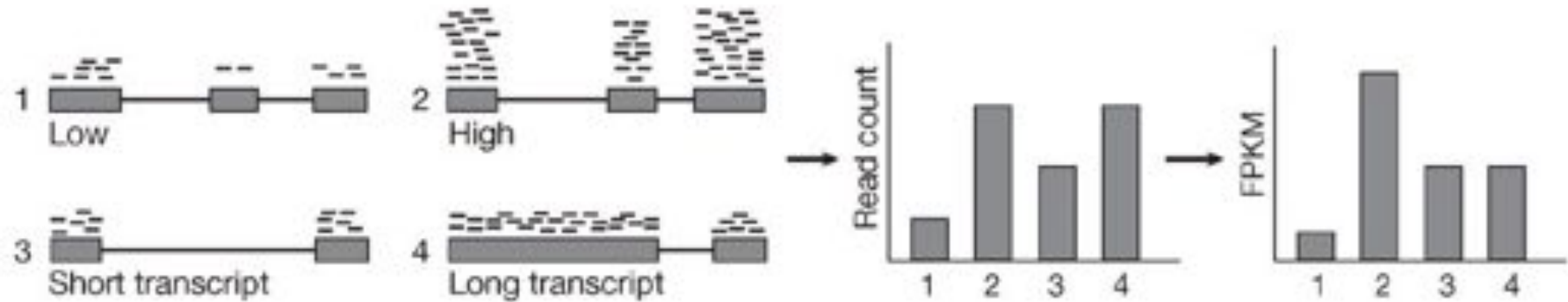**1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)**
=> Wait a second, reads in a pair arent independent!

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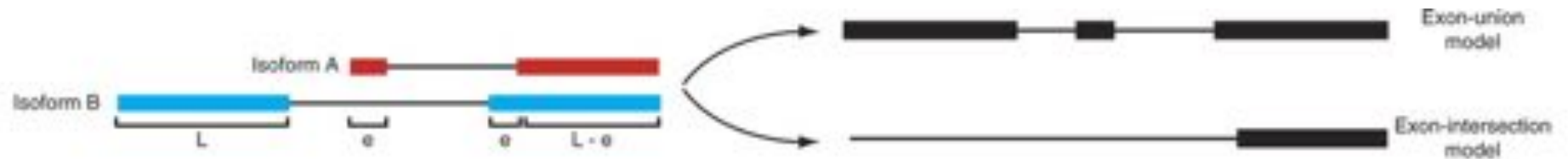⇒ If you were to sequence one million full length transcripts, TPM is the number of transcripts you would have seen of type i, given the abundances of the other transcripts in your sample

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

$$\text{TPM}_i = \left( \frac{\text{FPKM}_i}{\sum_j \text{FPKM}_j} \right) \cdot 10^6$$

# Gene or Isoform Quantification?

# Gene or Isoform Quantification?



**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
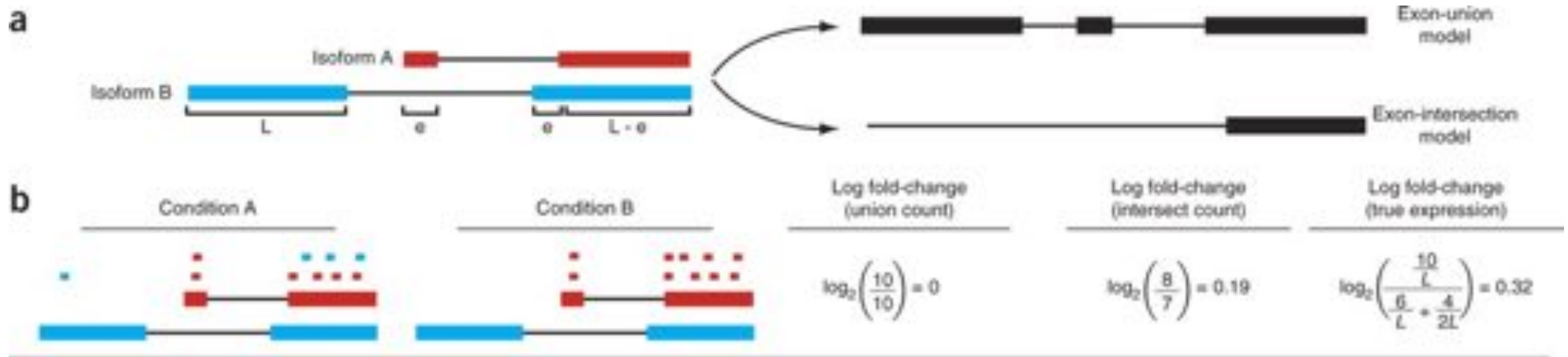
# Gene or Isoform Quantification?



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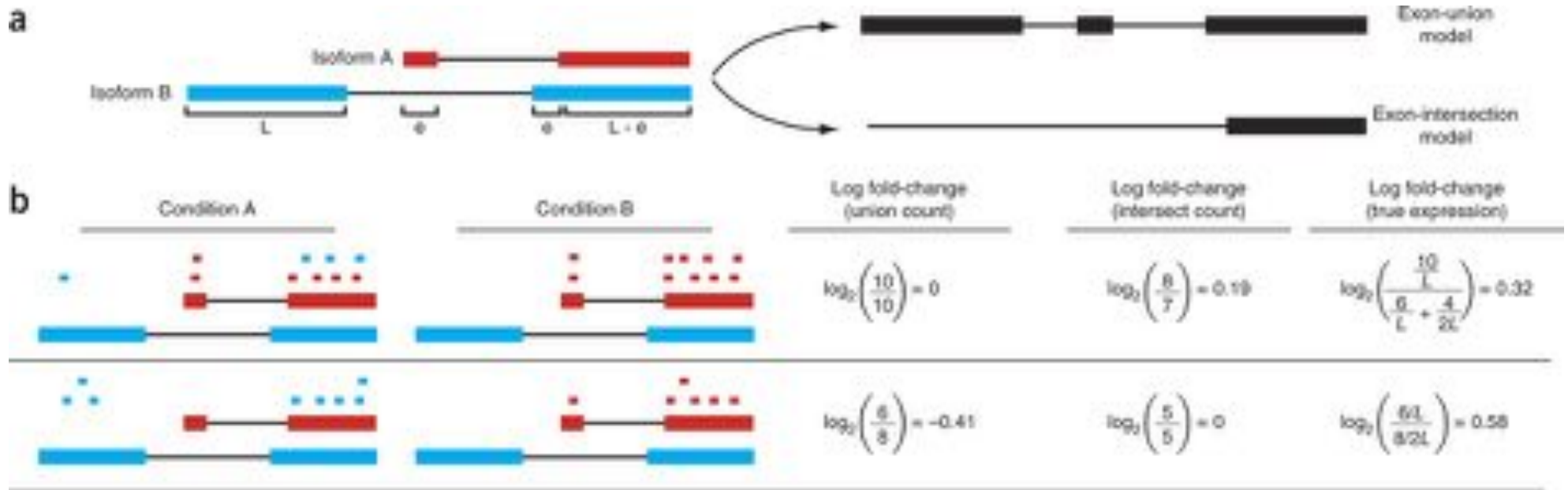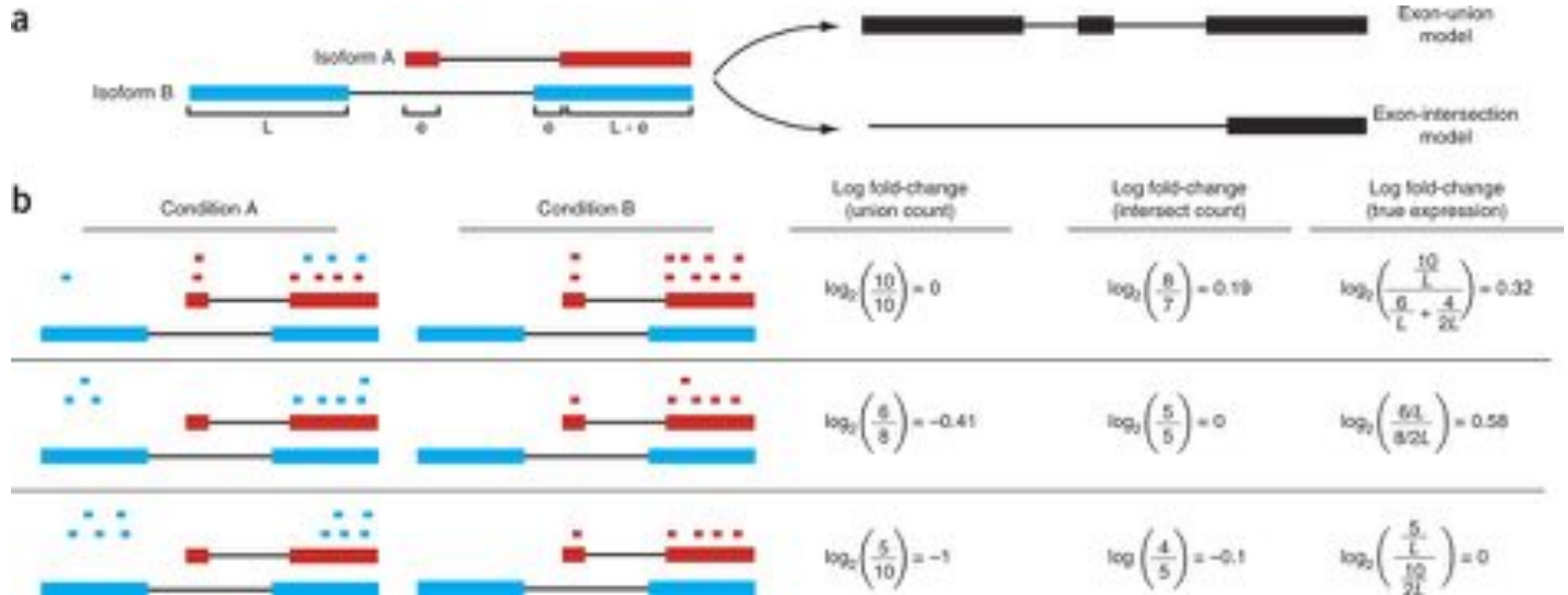Trapnell et al (2013) Nature Biotechnology 31, 46–53. doi:10.1038/nbt.2450
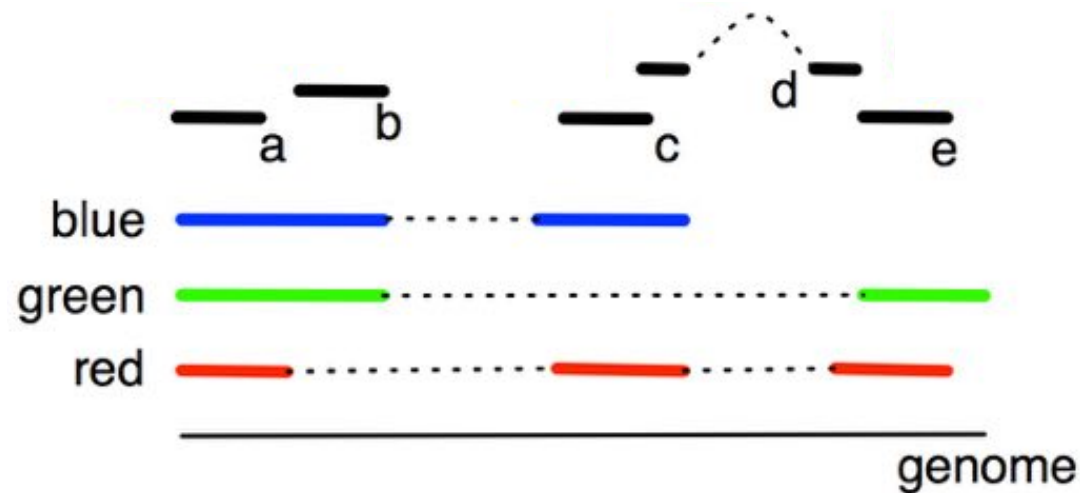
# Gene or Isoform Quantification?



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

# 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
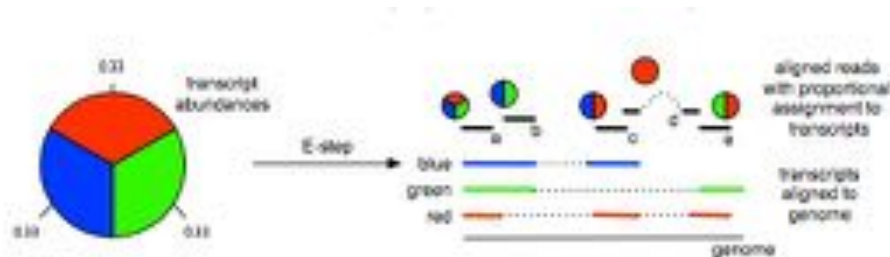- Reads b,c,e map to each of the three pairs of isoforms.

What is the most likely expression level of each isoform?

***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)

***Models for transcript quantification from RNA-seq***
Pachter, L (2011) arXiv. 1104.3889 [q-bio.GN]

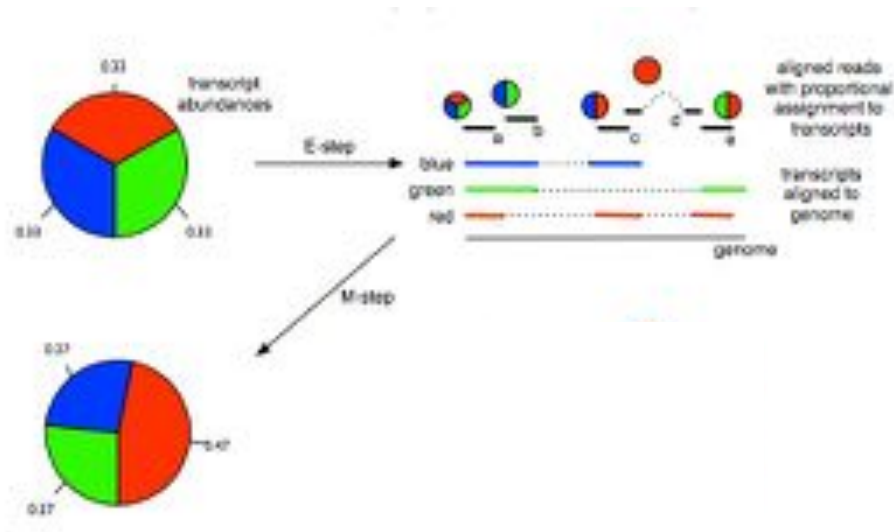# Multi-mapping? Isoform ambiguity? Expectation Maximization to the Rescue



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:
red:     0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)
blue:   0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)
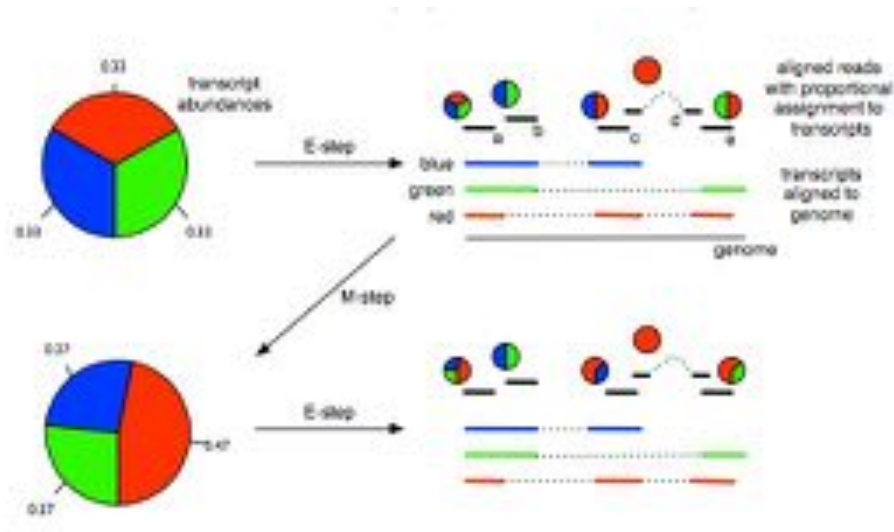green: 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)

***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:
red:    0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)
blue:   0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)
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:

red:    0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)
blue:   0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)
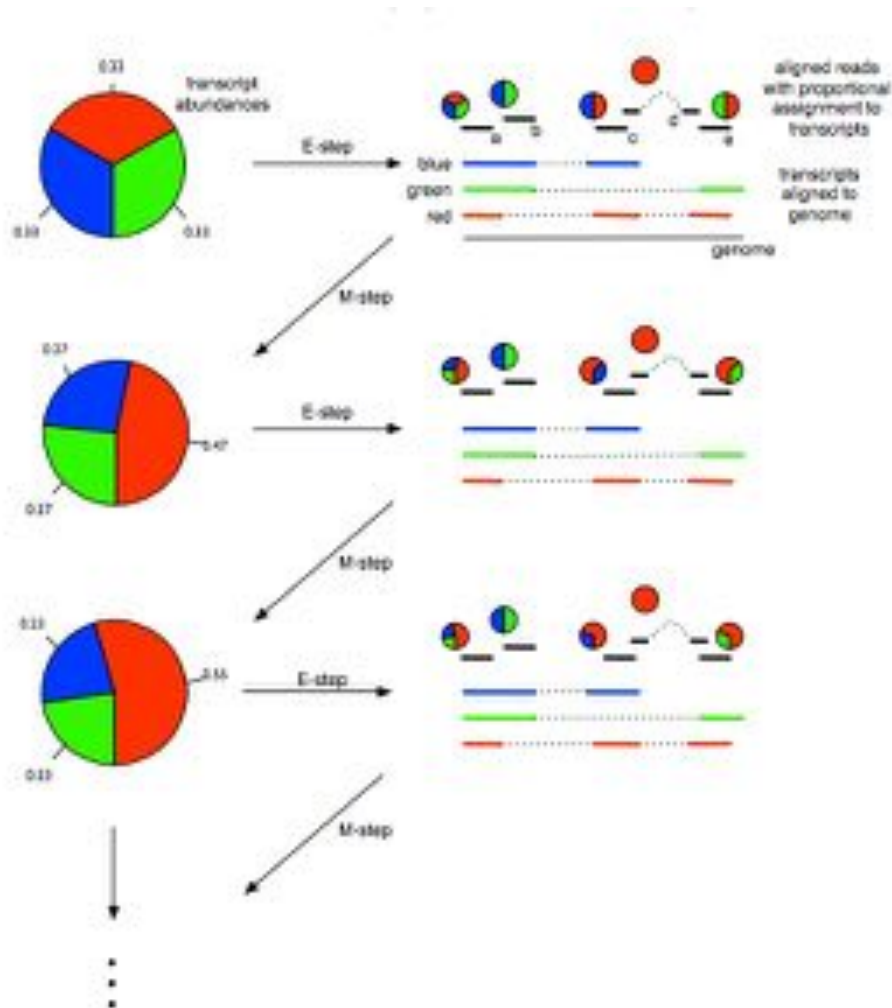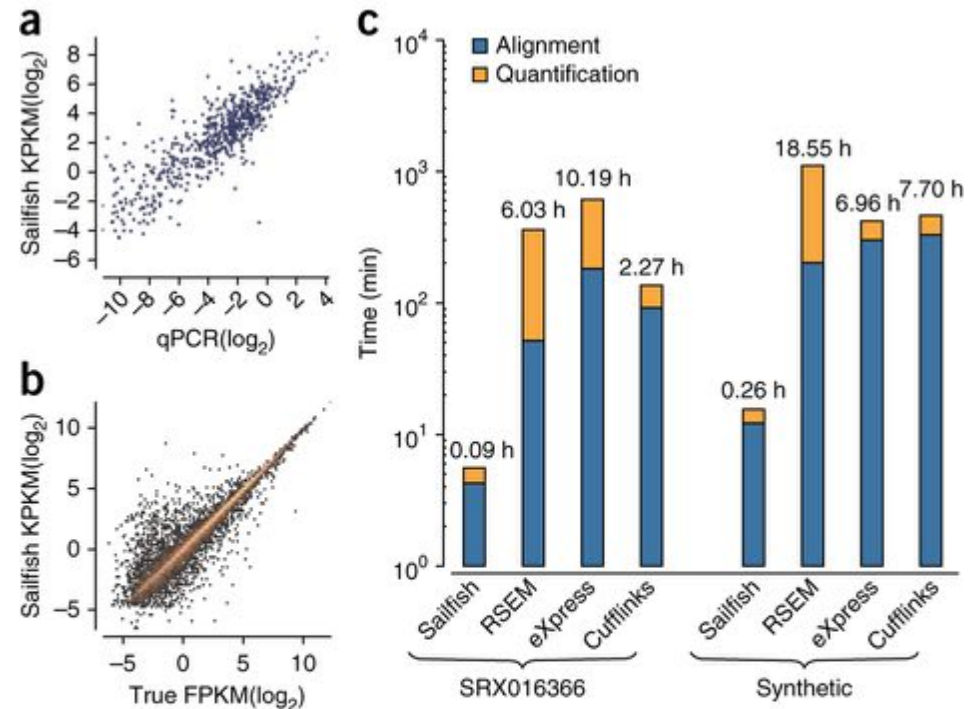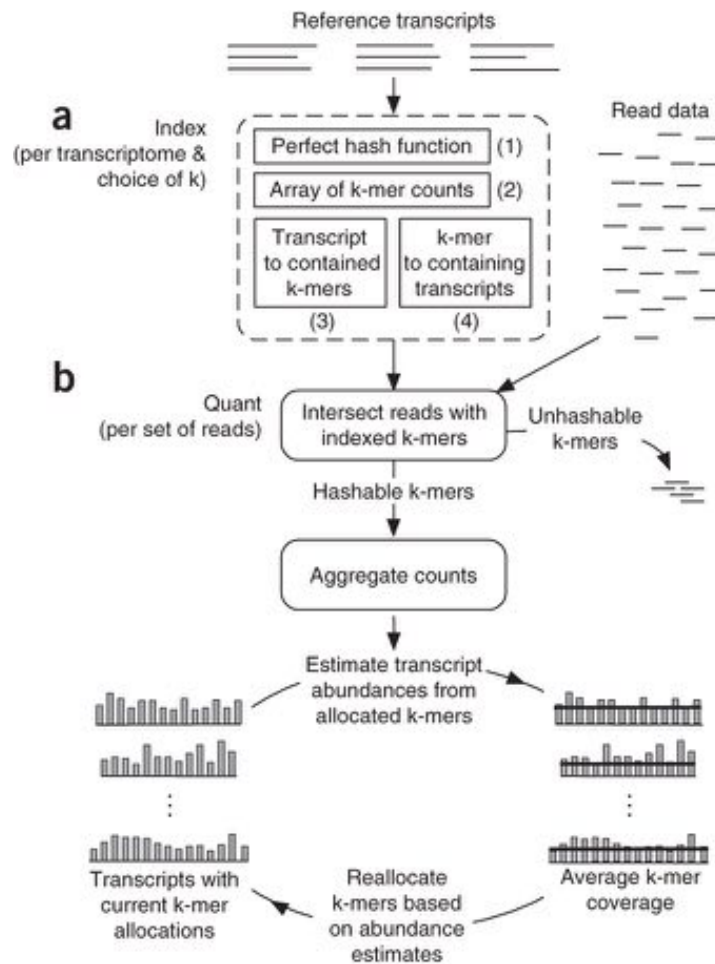green: 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)

Repeat until convergence!

***Models for transcript quantification from RNA-seq***
Pachter, L (2011) arXiv. 1104.3889 [q-bio.GN]

# 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

    1. 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 ☺

# Machine Learning Primer:

# Hidden Markov Models

# 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
    - {Modern, Neanderthal} for Ancestry



  - 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 | Fair) = .5, P(T | Fair) = .5, P(H | Biased) = .9, P(T | Biased) = .1
    - P(A | Gene) = .9, P(A | Not Gene) = .1 …

  - A probability of taking a transition
    - $P(s_i=Fair|s_{i-1}=Fair) = .9$, $P(s_i=Bias|s_{i-1} = Fair)$ .1
    - $P(s_i=Exon | s_{i-1}=Intergenic)$, …

# Why Hidden?

- 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 (exon/intron/intergenic/etc).
  - But we can *infer* the most likely hidden states of an HMM based on the given sequence of emitted symbols.



AAAGCATGCATTTAACGTGAGCACAATAGATTACA

# HMM Example - Casino Coin



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

# Three classic HMM problems

1. **Evaluation**: given a model and an output sequence, what is the probability that the model generated that output?

2. **Decoding**: given a model and an output sequence, what is the most likely state sequence through the model that generated the output?

3. **Learning**: given a model and a set of observed sequences, how do we set the model's parameters so that it has a high probability of generating those sequences?

# Three classic HMM problems

1. **Evaluation**: given a model and an output sequence, what is the probability that the model generated that output?

- To answer this, we consider all possible paths through the model

- Example: we might have a set of HMMs representing protein families -> pick the model with the best score

# Solving the Evaluation problem: The Forward algorithm

- To solve the Evaluation problem (probability that the model generated the sequence), we use the HMM and the data to build a *trellis*

- Filling in the trellis will give tell us the probability that the HMM generated the data by finding all possible paths that could do it
  - Especially useful to evaluate from which models, a given sequence is most likely to have originated

# Our sample HMM



Let $S_1$ be initial state, $S_2$ be final state

# A trellis for the Forward Algorithm

# A trellis for the Forward Algorithm



Time

|  | t=0 | t=1 | t=2 | t=3 |
|---|---|---|---|---|

$S_1$  1.0  $\xrightarrow{(0.6)(0.8)(1.0)}$  0.48  $\xrightarrow{(0.6)(0.2)(0.48)}$  .0576 + .018 = .0756

State

$S_2$  0.0  $\xrightarrow{(0.9)(0.3)(0)}$  0.20  $\xrightarrow{(0.9)(0.7)(0.2)}$  .126 + .096 = .222

$(0.1)(0.1)(0)$

$(0.4)(0.5)(1.0)$

$(0.1)(0.9)(0.2)$

$(0.4)(0.5)(0.48)$

Output:   A   C   C

# A trellis for the Forward Algorithm



Time

| | t=0 | t=1 | t=2 | t=3 |
|---|---|---|---|---|

S₁

S_1 row: 1.0 → (0.6)(0.8)(1.0) → 0.48 → (0.6)(0.2)(0.48) → .009072 + .01998 = .029052

State

S₂

S_2 row: 0.0 → 0.20 → .13986 + .01512 = .15498

Edge labels: (0.1)(0.1)(0) , (0.4)(0.5)(1.0) , (0.9)(0.3)(0) , (0.1)(0.9)(0.2) , (0.4)(0.5)(0.48) , (0.9)(0.7)(0) , (0.1)(0.9)(0.222) , (0.4)(0.5)(0.0756)

Output:  A    C    C

# A trellis for the Forward Algorithm



Time

t=0   t=1   t=2   t=3

S₁

| 1.0 | (0.6)(0.8)(1.0) → | 0.48 | (0.6)(0.2)(0.48) → | .0756 | (0.6)(0.2)(.0756) → | .029 |

State

S₂

| 0.0 | → | 0.20 | → | .222 | → | .155 |

(0.1)(0.1)(0)
(0.4)(0.5)(1.0)
(0.9)(0.3)(0)

(0.1)(0.9)(0.2)
(0.4)(0.5)(0.48)
(0.9)(0.7)(0.2)

(0.1)(0.9)(0.222)
(0.4)(0.5)(0.0756)
(0.9)(0.7)(0.222)

S2 is final state➔ 15.5% probability of this sequence given this model was used

# Probability of the model

- The Forward algorithm computes *P(y|M)*

- If we are comparing two or more models, we want the likelihood that each model generated the data: *P(M|y)*

  – Use Bayes' law:
  $$P(M \mid y) = \frac{P(y \mid M)P(M)}{P(y)}$$

  – Since P(y) is constant for a given input, we just need to maximize *P(y|M)P(M)*

# Three classic HMM problems

2. **Decoding**: given a model and an output sequence, what is the most likely state sequence through the model that generated the output?

- A solution to this problem gives us a way to match up an observed sequence and the states in the model.

AAAGCATGCATTTAACGAGAGCACAAGGGCTCTAATGCCG

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

# Three classic HMM problems

2. **Decoding**: given a model and an output sequence, what is the most likely state sequence through the model that generated the output?

- A solution to this problem gives us a way to match up an observed sequence and the states in the model.

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

# Solving the Decoding Problem: The Viterbi algorithm

- To solve the decoding problem (find the most likely sequence of states), we evaluate the Viterbi algorithm

$$V_i(t) = \begin{cases} 0 & : \quad t = 0 \wedge i \neq S_I \\ 1 & : \quad t = 0 \wedge i = S_I \\ \max V_j(t-1)a_{ji}b_{ji}(y) & : \quad t > 0 \end{cases}$$

Where $V_i(t)$ is the probability that the HMM is in state $i$ after generating the sequence $y_1, y_2, \ldots, y_t$, following the *most probable path* in the HMM

# A trellis for the Viterbi Algorithm



Time

| | t=0 | t=1 | t=2 | t=3 |
|---|---|---|---|---|

$S_1$

1.0   (0.6)(0.8)(1.0)   0.48

max

(0.1)(0.1)(0)

(0.4)(0.5)(1.0)

State

$S_2$

0.0   max   0.20

(0.9)(0.3)(0)

Output:    **A**        **C**        **C**

# A trellis for the Viterbi Algorithm

# A trellis for the Viterbi Algorithm



Time

| | t=0 | | t=1 | | t=2 | | t=3 |
|---|---|---|---|---|---|---|---|

$S_1$  1.0  $\xrightarrow{(0.6)(0.8)(1.0)}$  0.48  $\xrightarrow{(0.6)(0.2)(0.48)}$  max(.006912,.01134) = .01134

max

$(0.1)(0.1)(0)$   $(0.1)(0.9)(0.2)$   $(0.1)(0.9)(0.126)$

State

$(0.4)(0.5)(1.0)$   $(0.4)(0.5)(0.48)$   $(0.4)(0.5)(0.0576)$

$S_2$  0.0  $\xrightarrow{(0.9)(0.3)(0)}$  0.20  $\xrightarrow{(0.9)(0.7)(0.2)}$  max(.01152,.07938) = .07938

Output:   **A**   **C**   **C**

# A trellis for the Viterbi Algorithm



Time

| | t=0 | t=1 | t=2 | t=3 |
|---|---|---|---|---|

S₁ → $S_1$: 1.0 → (0.6)(0.8)(1.0) → 0.48 → (0.6)(0.2)(0.48) → .0576 → (0.6)(0.2)(0.0576) → .01134

State

S₂ → $S_2$: 0.0 → 0.20 → .126 → .07938

Transition labels:
(0.1)(0.1)(0)
(0.4)(0.5)(1.0)
(0.9)(0.3)(0)
max
(0.1)(0.9)(0.2)
(0.4)(0.5)(0.48)
(0.9)(0.7)(0.2)
max
(0.1)(0.9)(0.126)
(0.4)(0.5)(0.0576)
(0.9)(0.7)(0.126)
max

S2 is final state➜ the most probable sequence of states has a 7.9% probability

# A trellis for the Viterbi Algorithm



Time

t=0                    t=1                    t=2                    t=3

$S_1$   1.0   (0.6)(0.8)(1.0)   0.48   (0.6)(0.2)(0.48)   .0576   (0.6)(0.2)(0.0576)   .01134

max                    max                    max

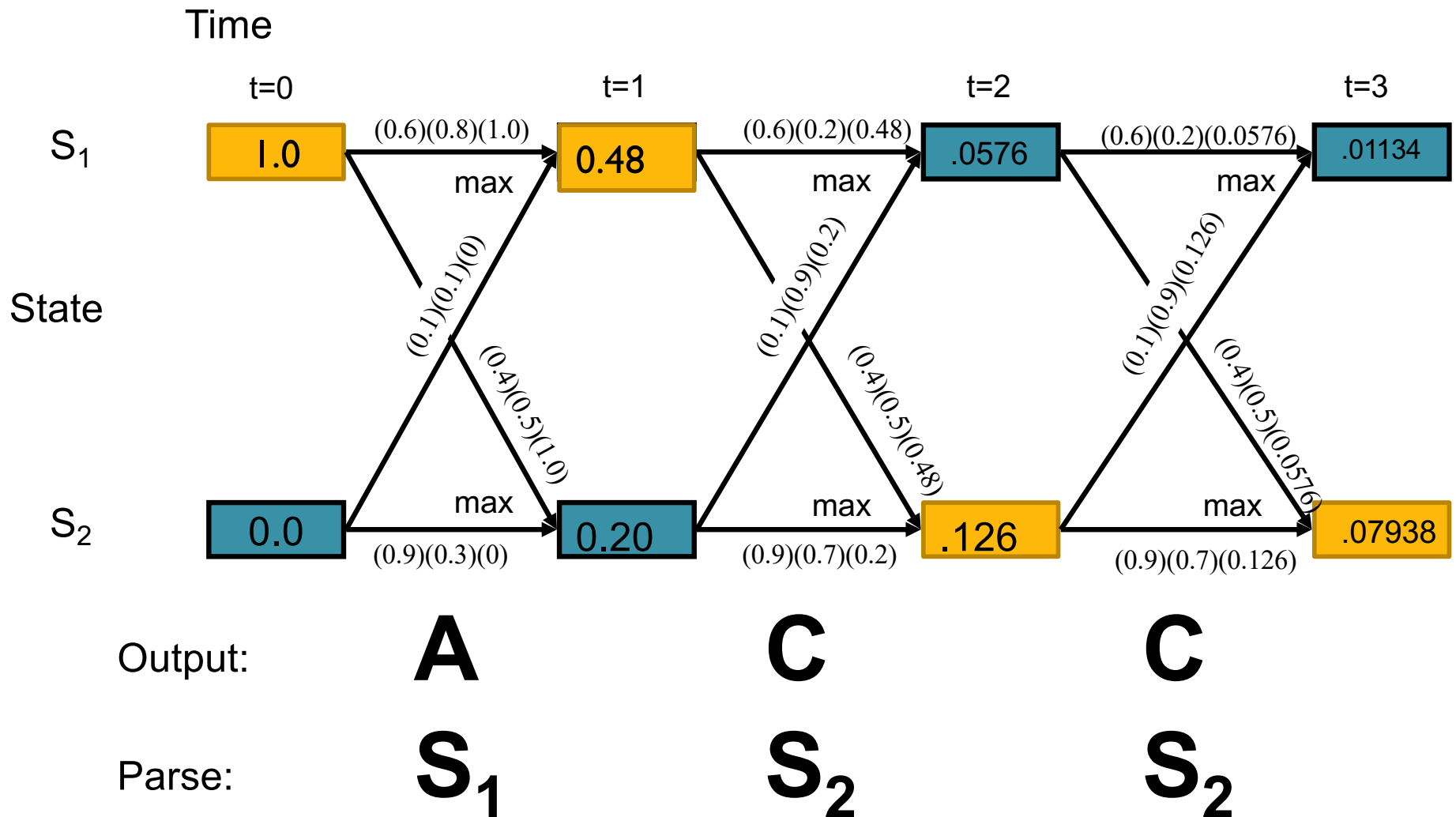(0.1)(1.0)(1.0)        (0.1)(0.9)(0.2)        (0.1)(0.9)(0.126)

(0.4)(0.5)(1.0)        (0.4)(0.5)(0.48)       (0.4)(0.5)(0.0576)

State

$S_2$   0.0   max   0.20   max   .126   max   .07938

(0.9)(0.3)(0)          (0.9)(0.7)(0.2)        (0.9)(0.7)(0.126)

Output:        **A**                    **C**                    **C**

Parse:         **$S_1$**                **$S_2$**                **$S_2$**

# Three classic HMM problems

3. **Learning**: given a model and a set of observed sequences, how do we set the model's parameters so that it has a high probability of generating those sequences?

- This is perhaps the most important, and most difficult problem.

- A solution to this problem allows us to determine all the probabilities in an HMMs by using an ensemble of training data

# Learning in HMMs:

- The learning algorithm uses Expectation-Maximization (E-M)
  - Also called the Forward-Backward algorithm
  - Also called the Baum-Welch algorithm

- In order to learn the parameters in an "empty" HMM, we need:
  - The topology of the HMM
  - Data - the more the better
  - Start with a random (or naïve) probability, repeat until converges

# 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