

Functional Genomics

Michael Schatz

Oct 14, 2019

Lecture 13: Computational Biomedical Research



Project Pitches

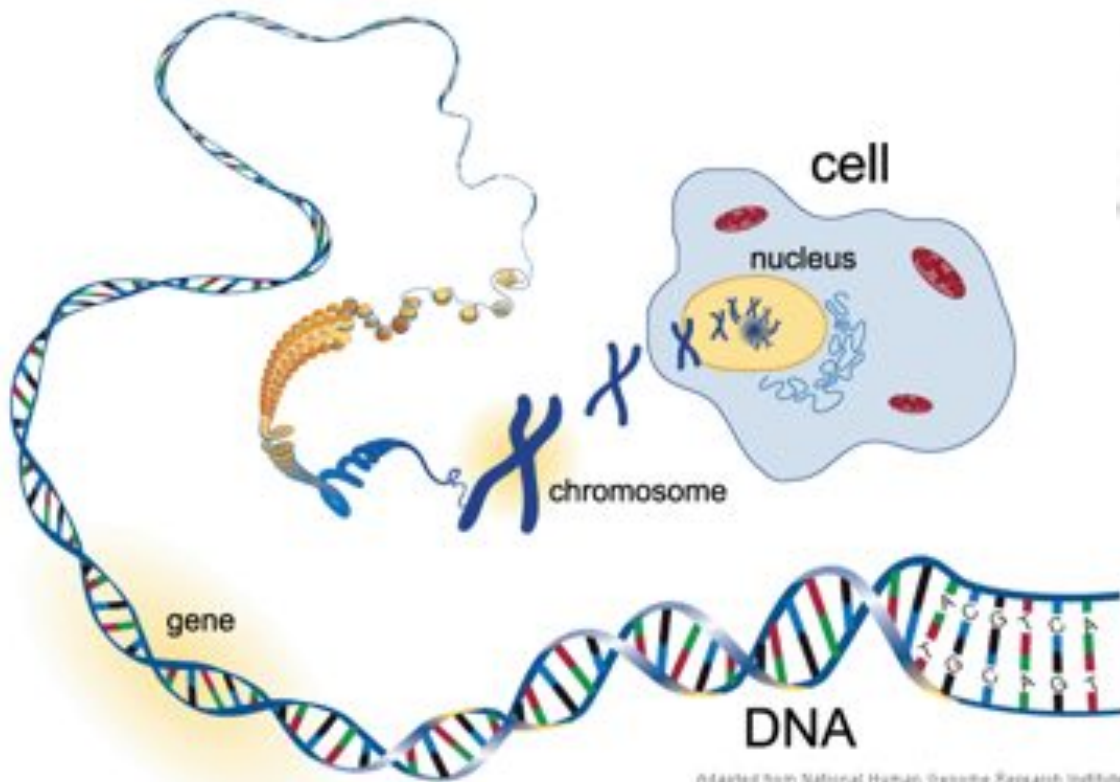
Student	Topic
Mary Joseph	Ethnic origins
Christian Seremetis	Ethnic origins
Gautam Prabhu	Disease risk: pathways, epistatic mutations
Joanna Guo	Disease risk: pathways, graph analysis
David Yang	Disease risk: classifiers
Kavya Tumkur	Disease risk: classifiers
Richard Xu	Disease risk: SVs, classifiers

Project Timeline

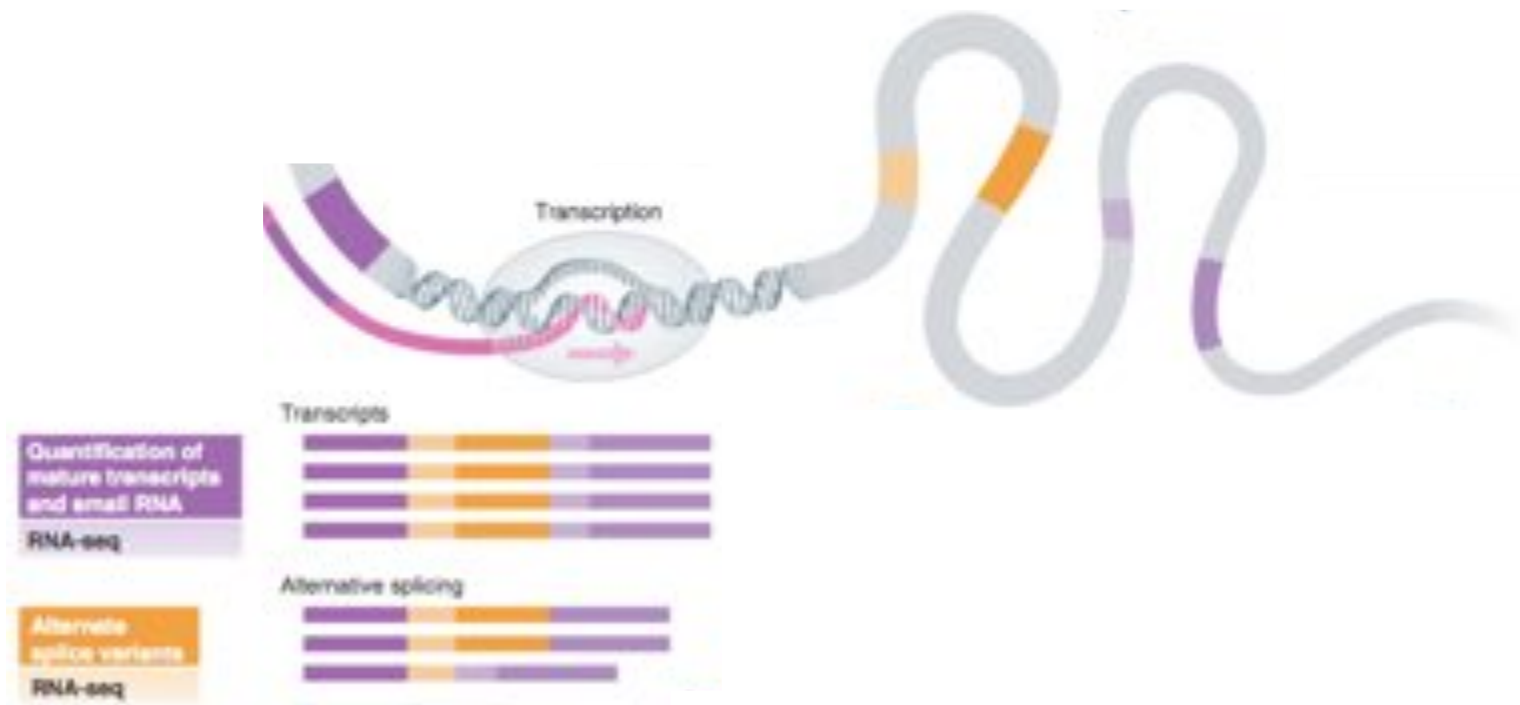
Week	Date	Deliverable
1	Oct 14	Decide teams
2	Oct 21	Abstract + Presentation
3	Oct 28	
4	Nov 4	Intern Report
5	Nov 11	
6	Nov 18	
7	Nov 25	<Thanksgiving>
8	Dec 2	In class presentation
9	Dec 9	
10	Dec 16	Final Report Due

Sequencing techniques

Much of the capacity is used to sequence genomes (or exomes) of individuals...



... but biology is much more than just genomes...



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).

Goal: Genome Annotations

[illegible]

Goal: Genome Annotations

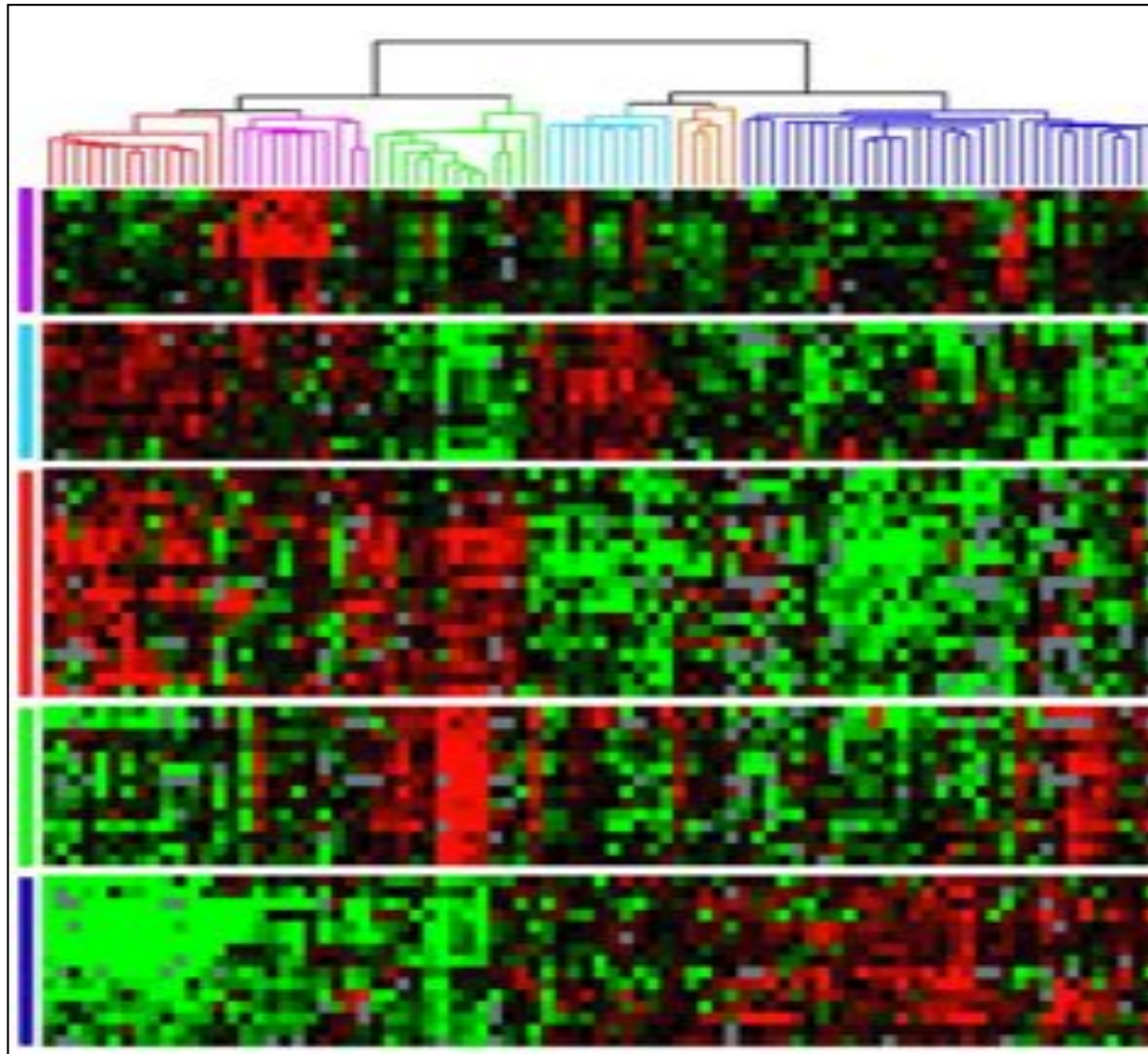
[illegible]



Outline

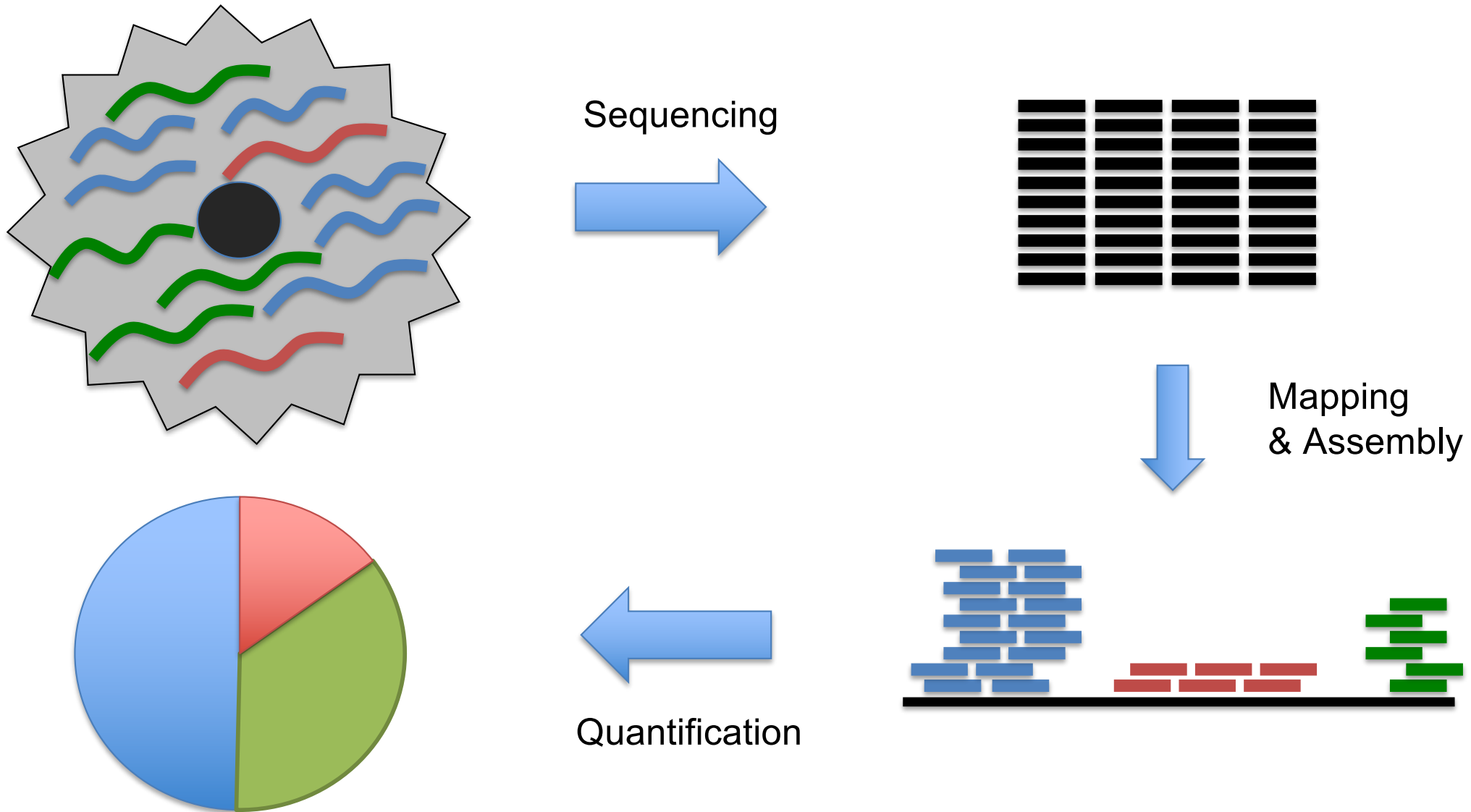
1. Experimental: RNAseq
2. Homology: Alignment to other genomes
3. Prediction: “Gene Finding”

RNA-seq

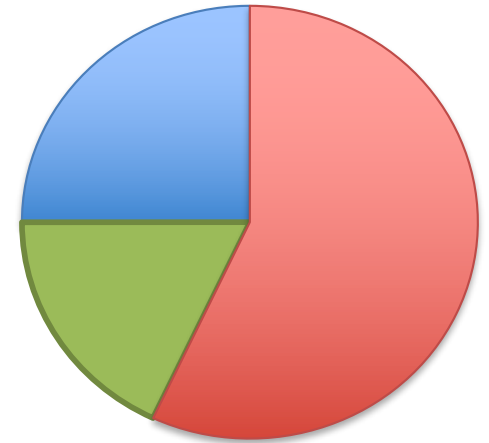
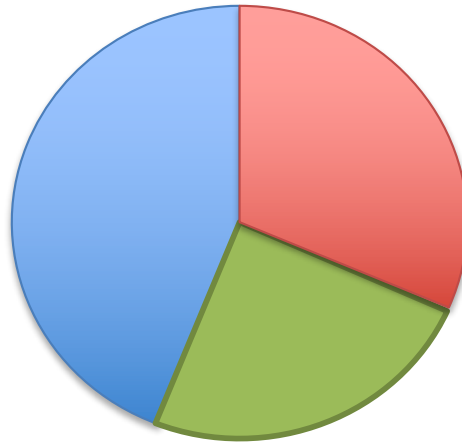
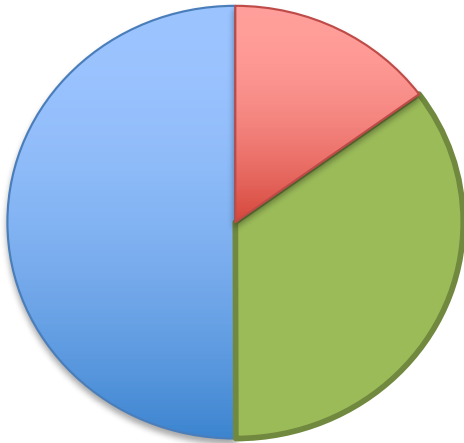
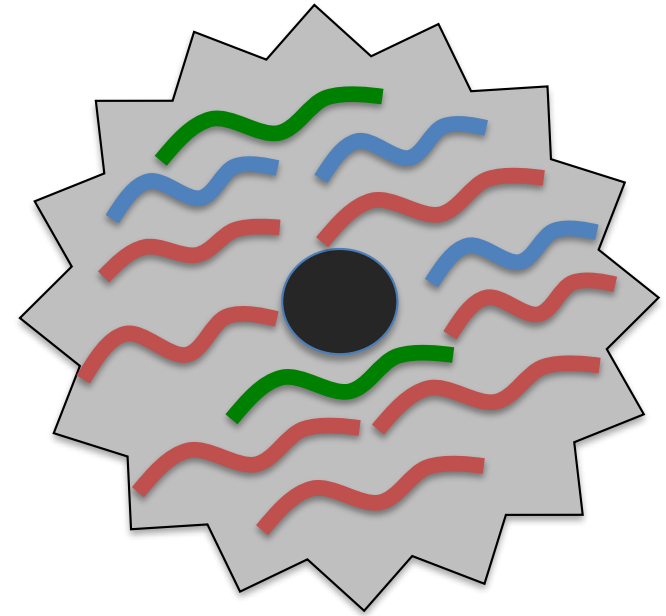
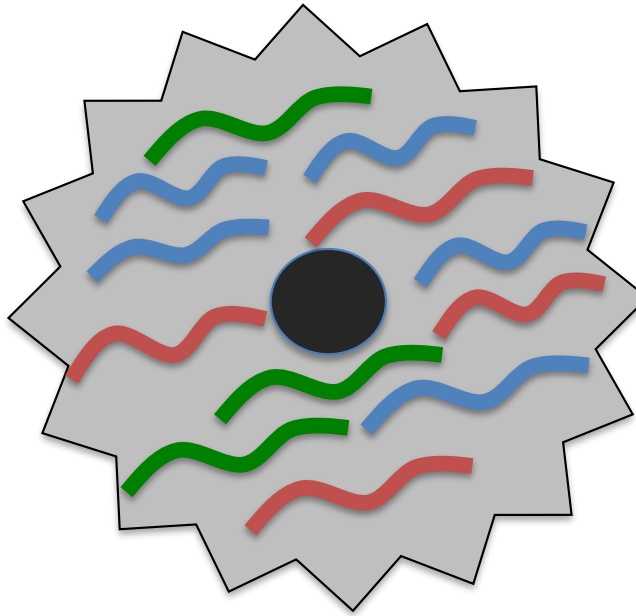
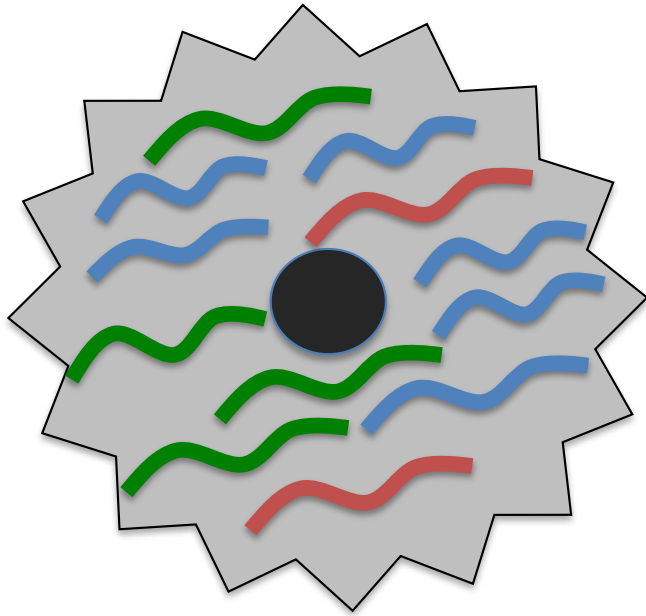


Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications.
Sørli et al (2001) *PNAS*. 98(19):10869-74.

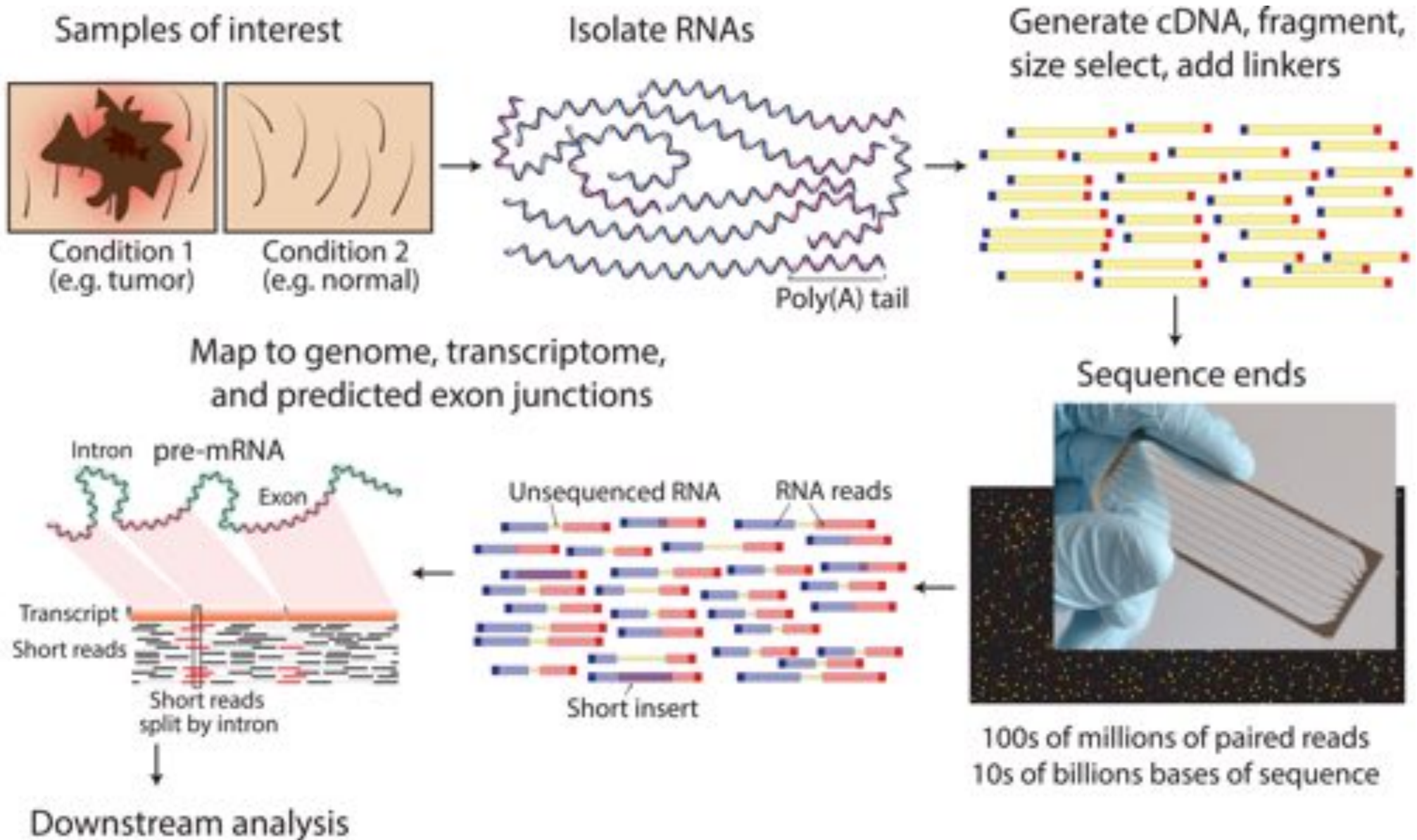
RNA-seq Overview



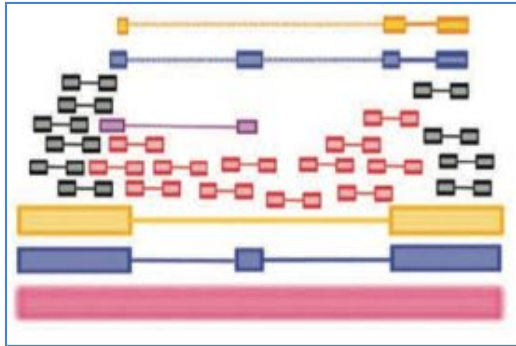
RNA-seq Overview



RNA-seq Overview



RNA-seq Challenges



Challenge 1: Eukaryotic genes are spliced

RNA-Seq Approaches

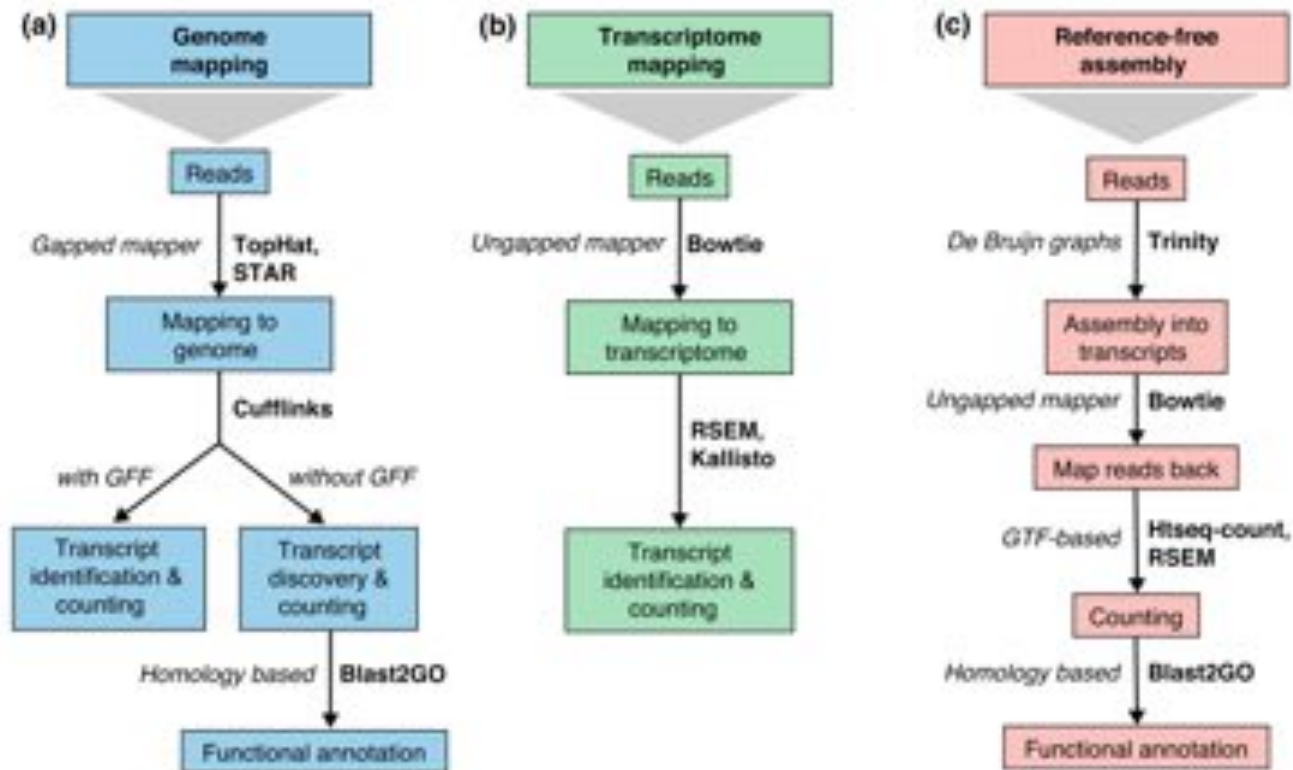


Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome with a gapped mapper. Next (novel) transcript discovery and quantification can proceed with or without an annotation 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

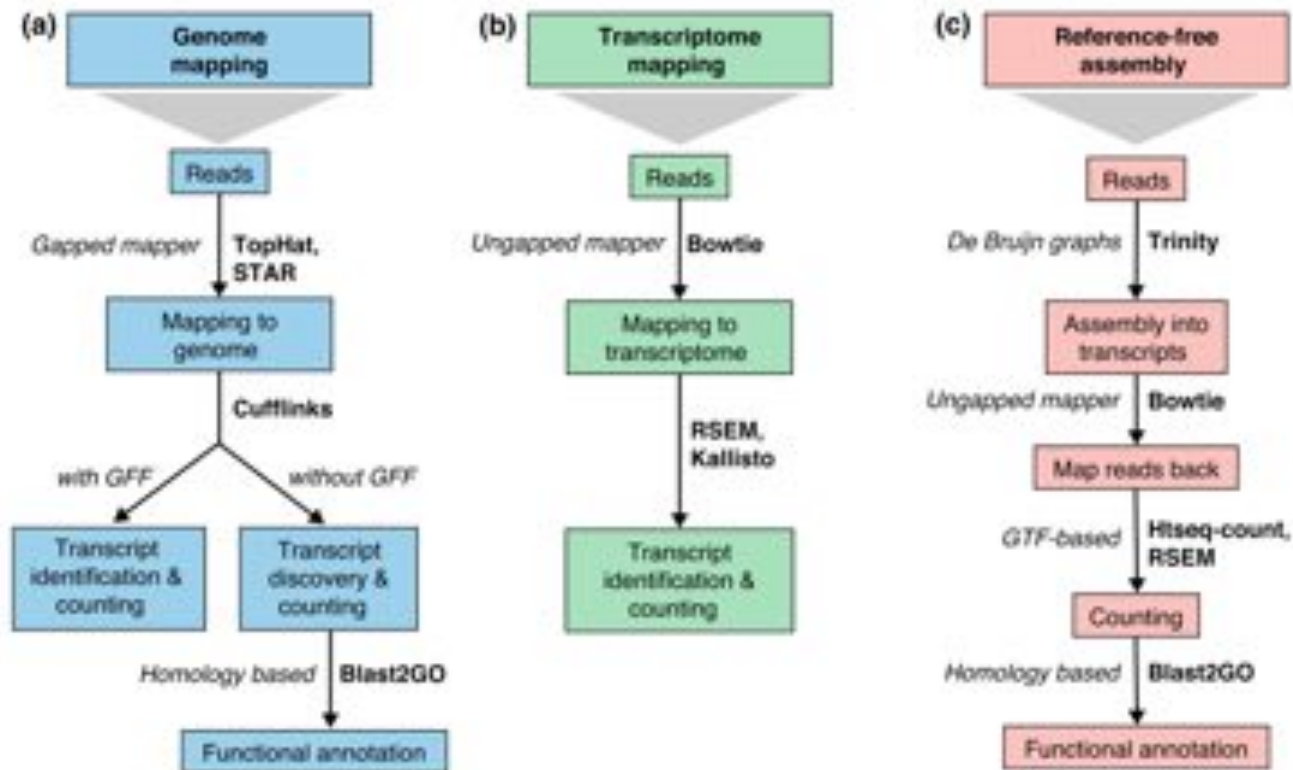


Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the reference genome. **b** Transcript discovery and quantification can proceed with or without an annotated transcriptome. **c** When 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 is 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

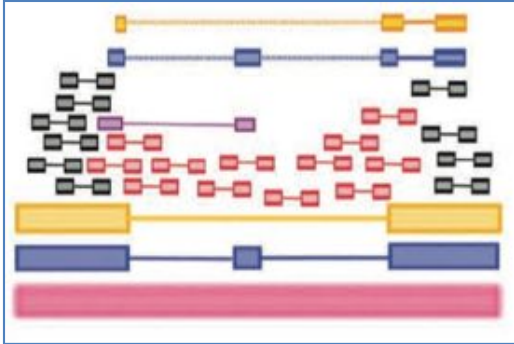
Which approach should we use?

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

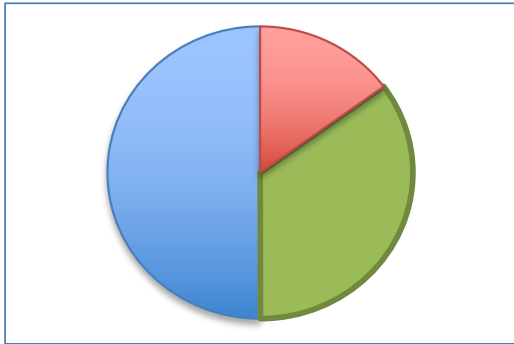


Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

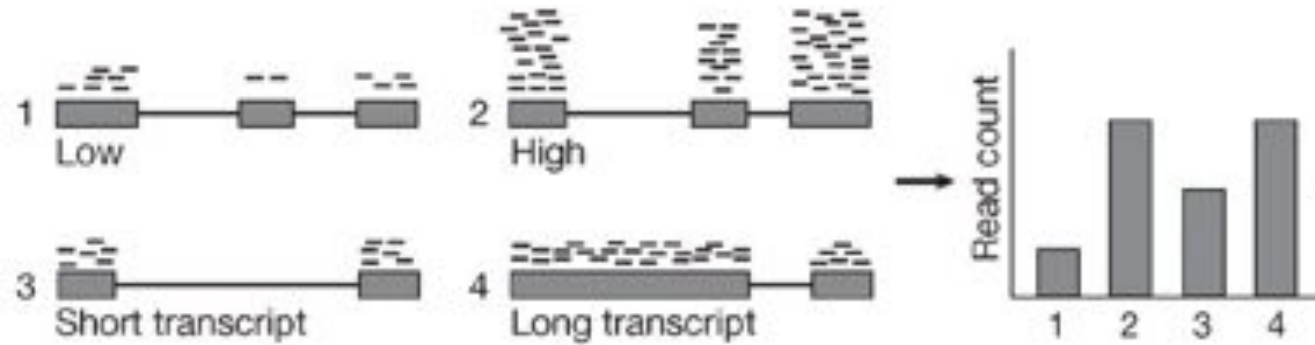
TopHat: discovering spliced junctions with RNA-Seq.

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



Challenge 2: Read Count != Transcript abundance

RPKM, FPKM, TPM

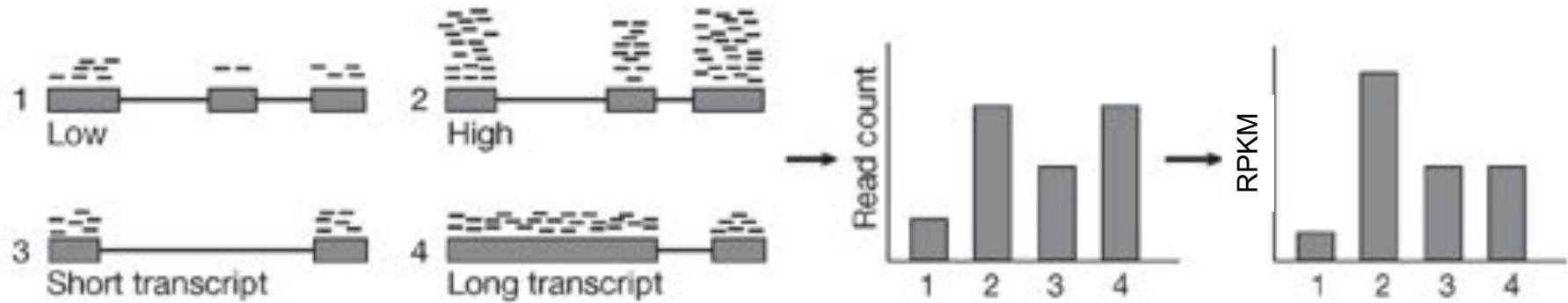


Counting Reads that align to a gene DOESN'T work!

- 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!

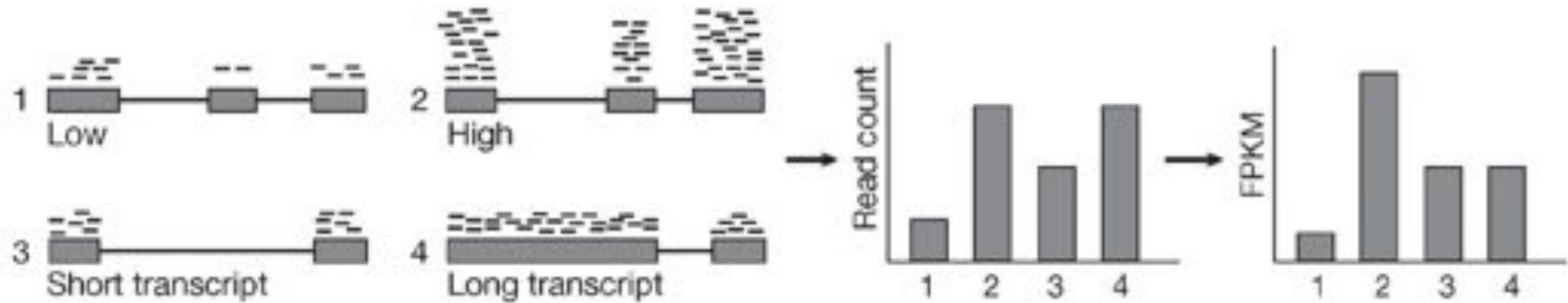
- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

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 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)

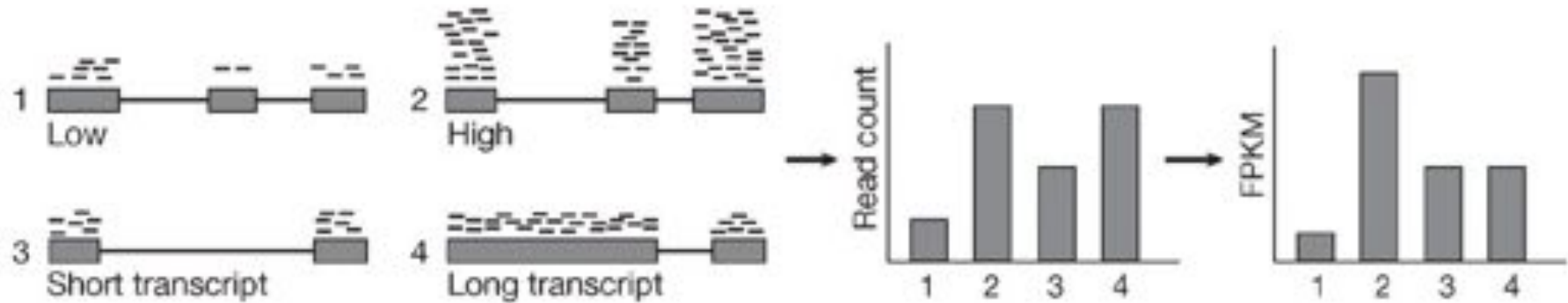
⇒ Wait a second, reads in a pair are independent!

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 are 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)

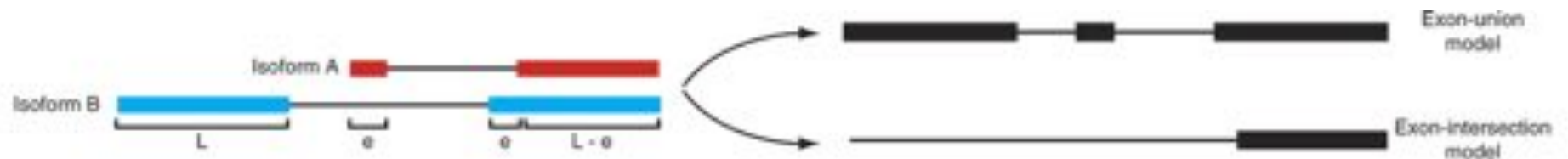
=> If you were to sequence one million full length transcripts, TPM is the number of transcripts you would have seen of type i , given the abundances of the other transcripts in your sample

=> 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?

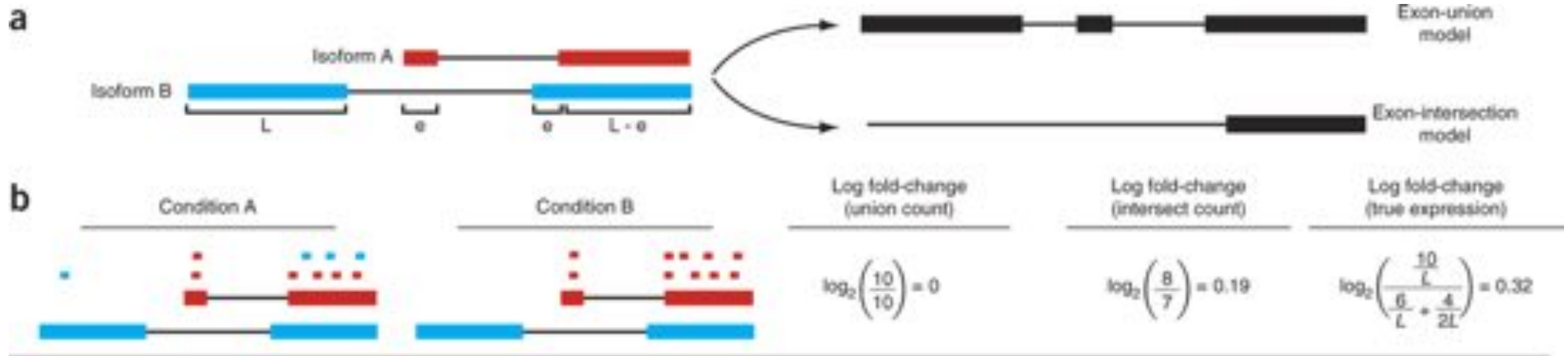
a



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

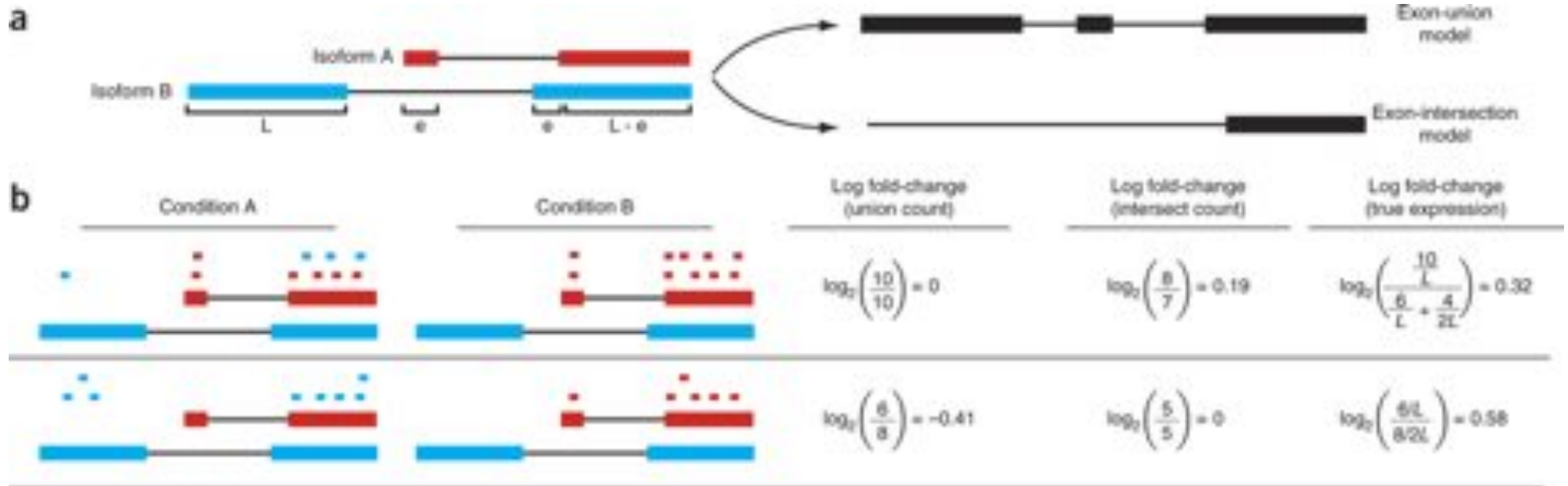
Gene or Isoform Quantification?



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

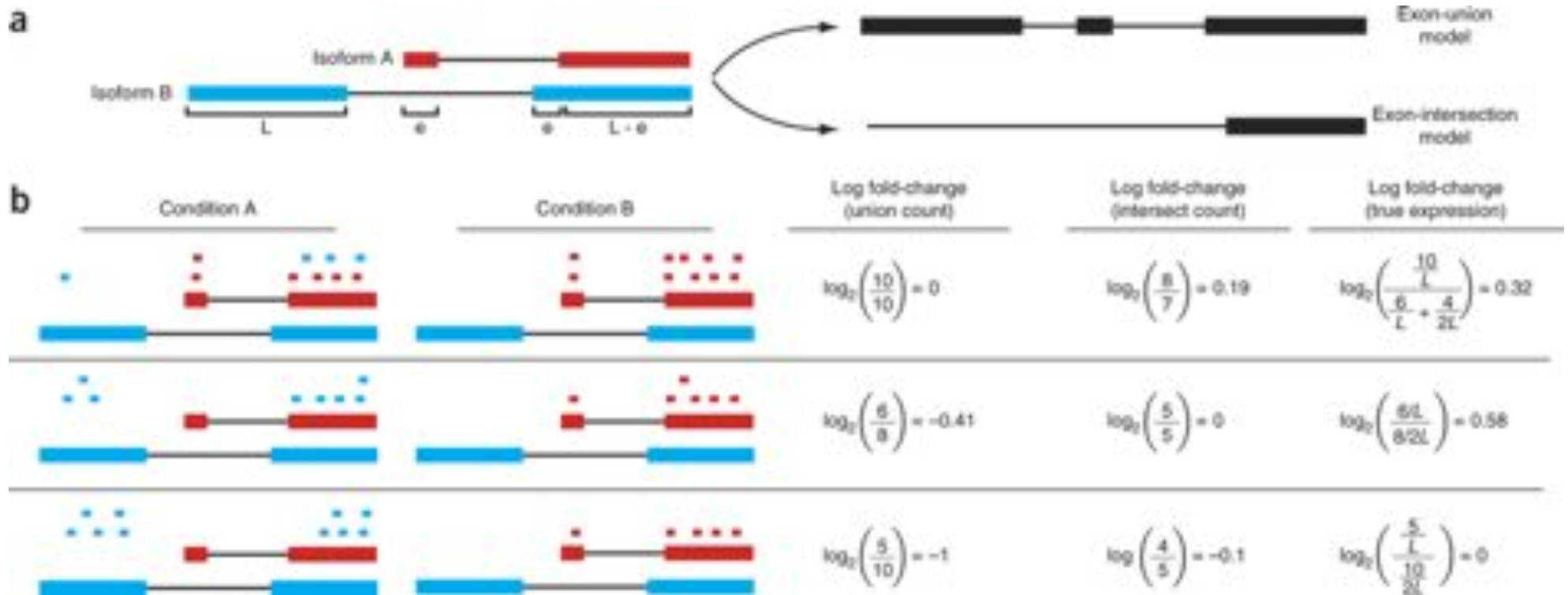
Gene or Isoform Quantification?



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?



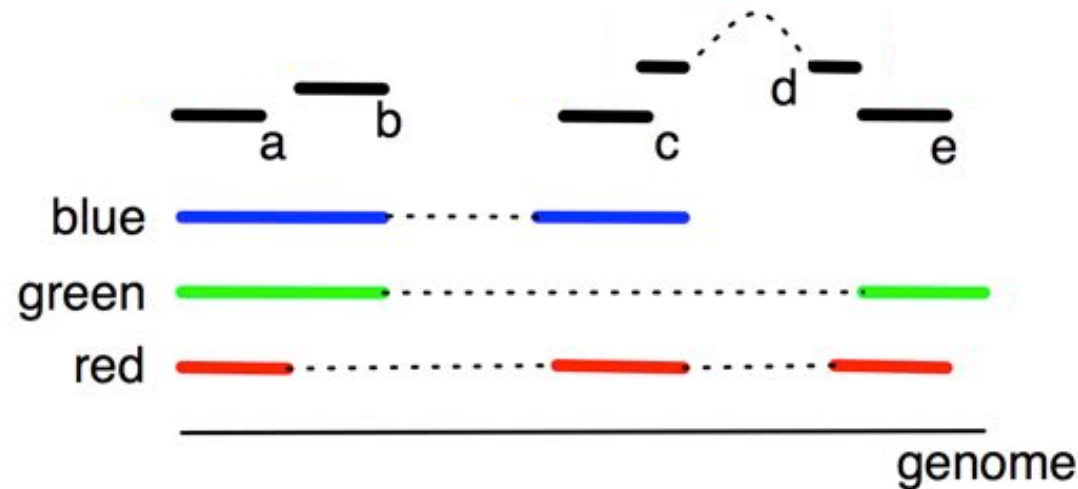
Key point : The length of the actual molecule from which the fragments derive is crucially important to obtaining accurate abundance estimates.

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.

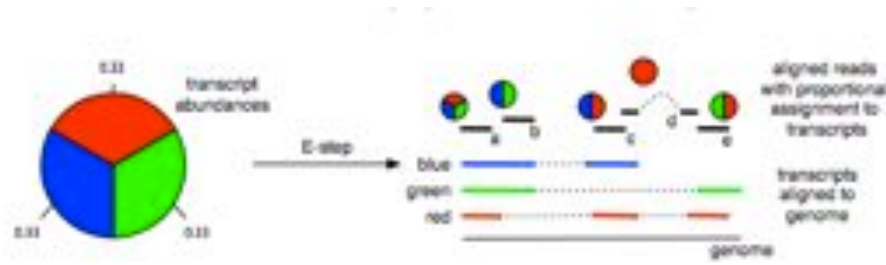
What is the most likely expression level of each isoform?

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



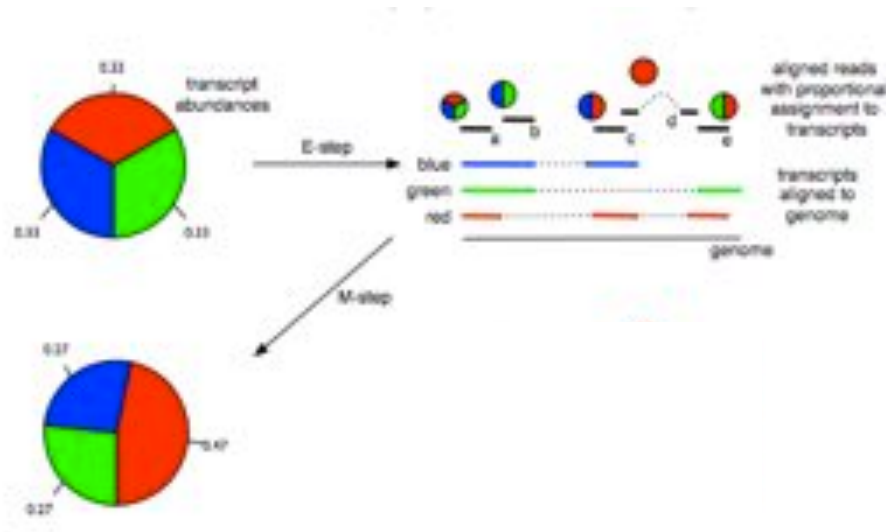
The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

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): $a=(.33,.33,.33)$, $b=(0,.5,.5)$, $c=(.5,.5,0)$, $d=(1,0,0)$, $e=(.5,.5,0)$

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:

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

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

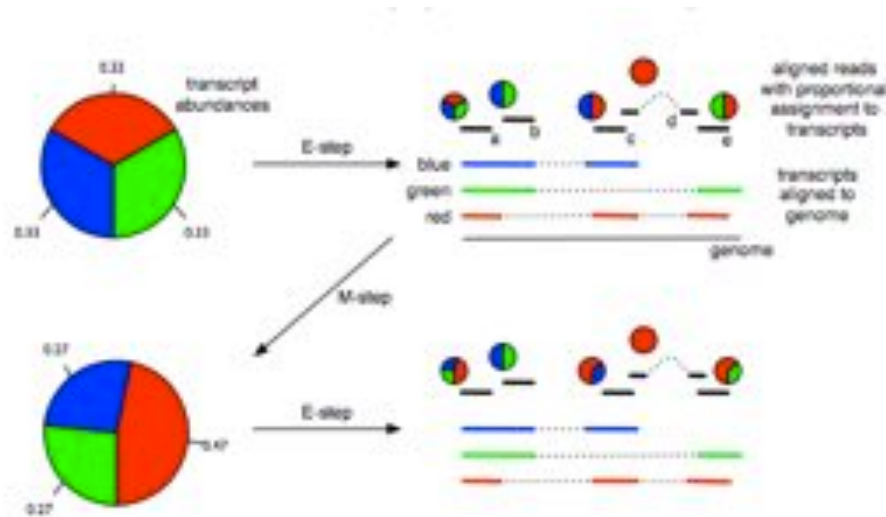
$$\text{green: } 0.27 = (0.33 + 0.5 + 0.5) / (2.33 + 1.33 + 1.33)$$

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:

$$\begin{aligned} \text{red: } 0.47 &= (0.33 + 0.5 + 1 + 0.5) / (2.33 + 1.33 + 1.33) \\ \text{blue: } 0.27 &= (0.33 + 0.5 + 0.5) / (2.33 + 1.33 + 1.33) \\ \text{green: } 0.27 &= (0.33 + 0.5 + 0.5) / (2.33 + 1.33 + 1.33) \end{aligned}$$

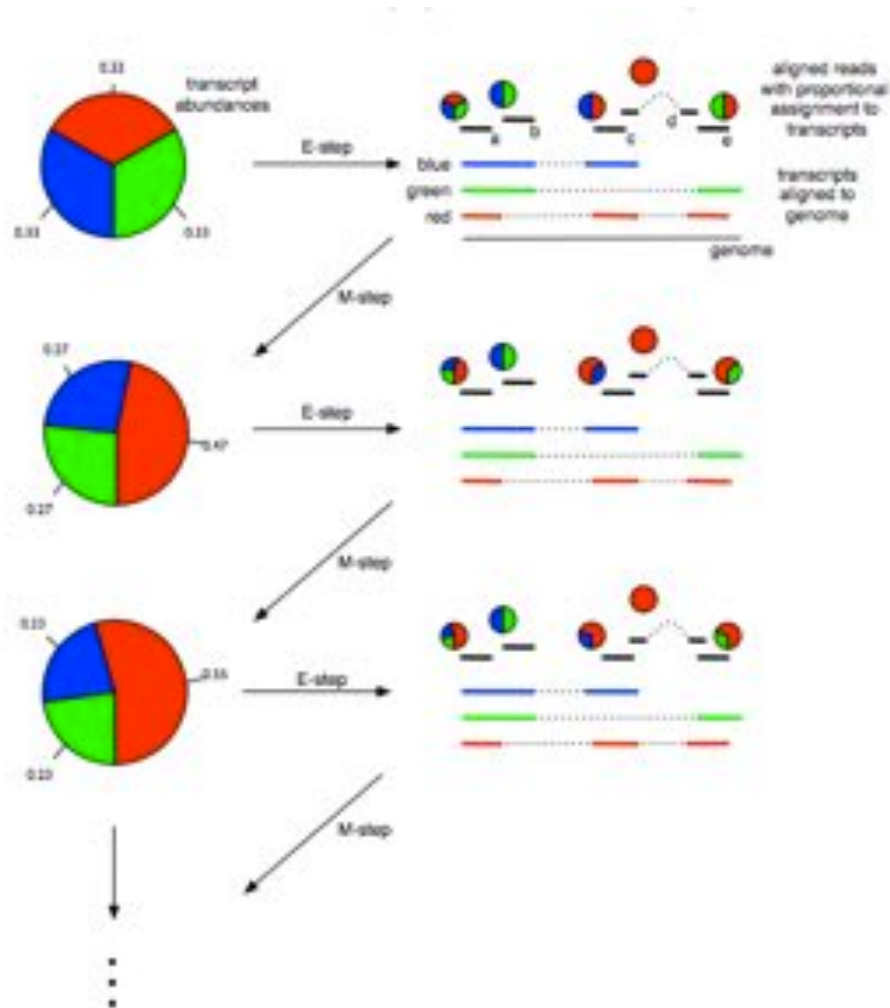
Repeat until convergence!

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): $a=(.33,.33,.33)$, $b=(0,.5,.5)$, $c=(.5,.5,0)$, $d=(1,0,0)$, $e=(.5,.5,0)$

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

$$\begin{aligned} \text{red: } 0.47 &= (0.33 + 0.5 + 1 + 0.5) / (2.33 + 1.33 + 1.33) \\ \text{blue: } 0.27 &= (0.33 + 0.5 + 0.5) / (2.33 + 1.33 + 1.33) \\ \text{green: } 0.27 &= (0.33 + 0.5 + 0.5) / (2.33 + 1.33 + 1.33) \end{aligned}$$

Repeat until convergence!

Models for transcript quantification from RNA-seq

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

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 SKTAVAPIERVKLLQLVQHASKQITADKQYKGIIDCVVRIPKEQGV
FLIDLASGGTAAAV SKTAVAPIERVKLLQLVQDASKAIAVDKRYKGIMDVLIRVPKEQGV

- 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 SKTAVAPIERVKLL LQVQHASKQITADKQYKGI IDCVVRI PKEQGV
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
FLIDLASGGTAAAV SKTAVAPIERVKLL LQVQDASKAIAVDKRYKGIMDVLIRVPKEQGV
```

- 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

Genetic Code

		1st base									
		U		C		A		G			
2nd base	U	UUU	Phenylalanine	UCU	Serine	UAU	Tyrosine	UGU	Cysteine	U	3rd base
		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	
		Nonpolar, aliphatic		Polar, uncharged		Aromatic		Positively charged		Negatively charged	

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  KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPA 115
          K HGKKV  A ++ +AH+D++      + LS+LH  KL VDP NF+LL + L+  LA H
Sbjct   61  KAHGKKVLGAFSDGLAHL DNLKGTFATLSELHCDKLHVDPENFRLLGNVLCVLAHHFGK 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	LSPADKTNVKA	AWGKVGA	HAGEYGA	EALERMFL	SFPTTKTY	FPHF-----	DLSHGSAQV	55
		LS	+	V	WGKV	A	+G E L R+F	P T F F D S +	
Sbjct	3	LSDGEWQLVL	NVWGKVE	ADIPGHGQ	EVLR	LFK	GHPE	TKLEKFDKFKHLKSEDEM	KASEDL 62

Query	56	KGHGKKVAD	ALTNAVA	HVDDMPN	ALSALSD	LHAHKLR	VDPVNF	KLLSHCLL	VTLA	AHLPA 115
		K	HG	V	AL	+	+	L+	HA	K ++ + +S C++ L + P
Sbjct	63	KKHGATVLT	ALGGILK	KKKGH	HEAEIK	PLAQSH	ATKHKI	PVKYLE	FISECII	QVLQSKHPG 122

Query	116	EFTPAVHAS	LDKFLAS	VSTVLT	SKYR	141
		+F	+++K	L	+ S Y+	
Sbjct	123	DFGADAQ	GAMNKAL	ELFRK	DMA	SNYK 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%)

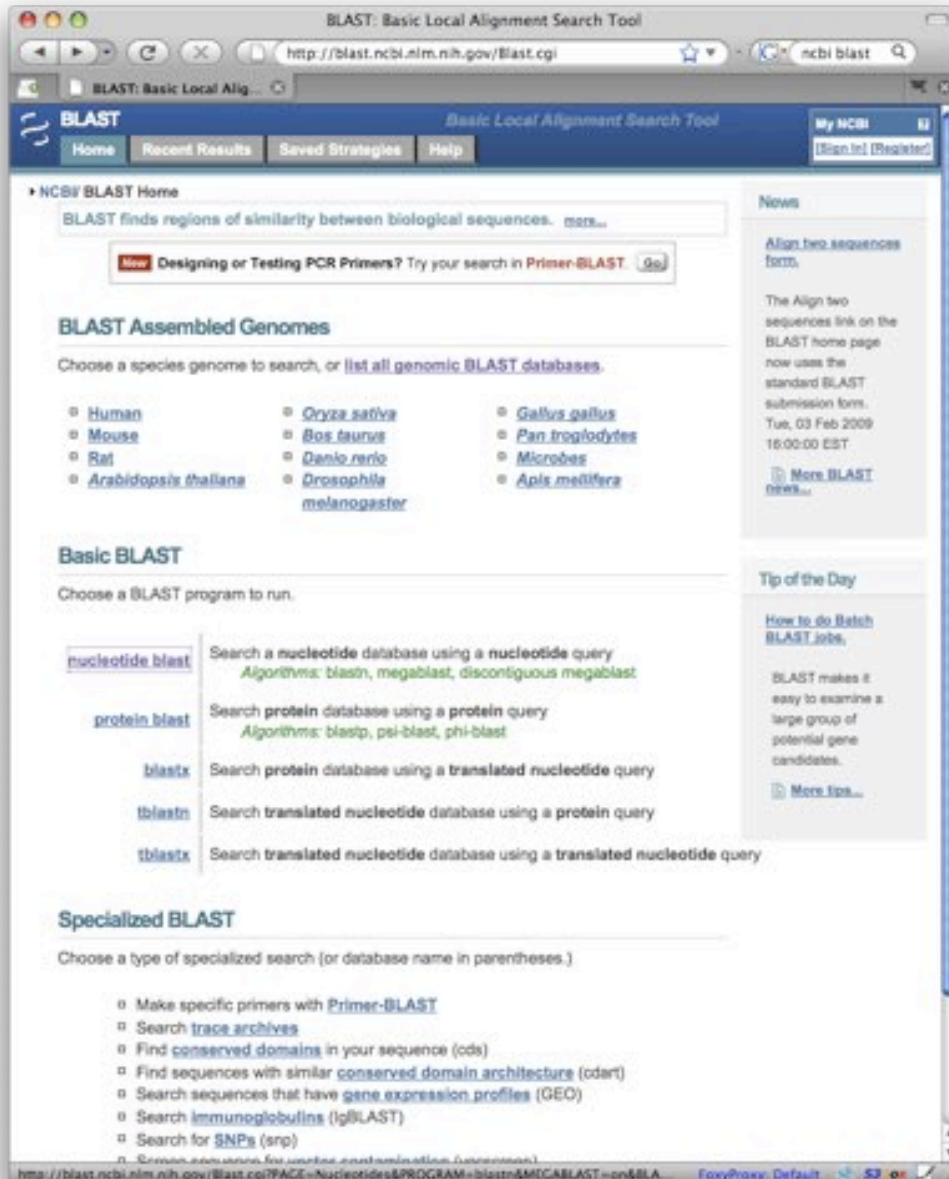
```
Query   30  ERMFLSFPTTKTYFPHFDSLHGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAH   89
        ++M  ++P      P+F+ +H  +      + +A AL N  ++DD+  +LSA  D
Sbjct   59  QKMLGNYPEV---LPYFNKAHQISL--SQPRILAFALLNYAKNIDDL-TSLSAFMDQIVV  112

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

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