

Lecture 13. RNAseq

Michael Schatz

March 9, 2020

Applied Comparative Genomics



Assignment 5: Due Wed Mar 11

Assignment 5: Annotations and RNA-seq

Assignment Date: Wednesday, March 4, 2020

Due Date: Wednesday, March 11, 2020 @ 11:59pm

Assignment Overview

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

Question 1. Gene Annotation Preliminaries [10 pts]

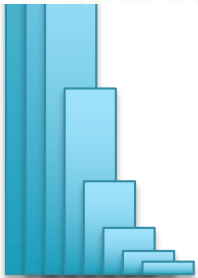
Download the annotation of build 38 of the human genome from here: ftp://ftp.ensembl.org/pub/release-87/gtf/homo_sapiens/Homo_sapiens.GRCh38.87.gtf.gz

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

Question 2. Sampling Simulation [10 pts]

A typical human cell has ~250,000 transcripts, and a typical bulk RNA-seq experiment may involve millions of cells. Consequently in an RNAseq experiment you may start with trillions of RNA molecules, although your sequencer will only give a few million to billions of reads. Therefore your RNAseq experiment will be a small sampling of the full composition. We hope the sequences will be a representative sample of the total population, but if your sample is very unlucky or biased it may not represent the true distribution. We will explore this concept by sampling a small subset of transcripts (500 to 50000) out of a much larger set (1M) so that you can evaluate this bias.

In `data1.txt` with 1,000,000 lines we provide an abstraction of RNA-seq data where normalization has been performed and the number of times a gene name occurs corresponds to the number of transcripts in the sample.





Project Proposal

Project Proposal

Assignment Date: Monday March 9, 2020

Due Date: Monday, March 16 2020 @ 11:59pm

Review the [Project Ideas](#) page

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

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 will 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!



Goal: Genome Annotations

[illegible]

Goal: Genome Annotations

[illegible]



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



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    LSPADKTNVKAAWGKVGAGHAGEYGAELERMFLSFPTTKTYFPHF-----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    KAHGKKVLGAFSDGLAHLNLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGK 120
```

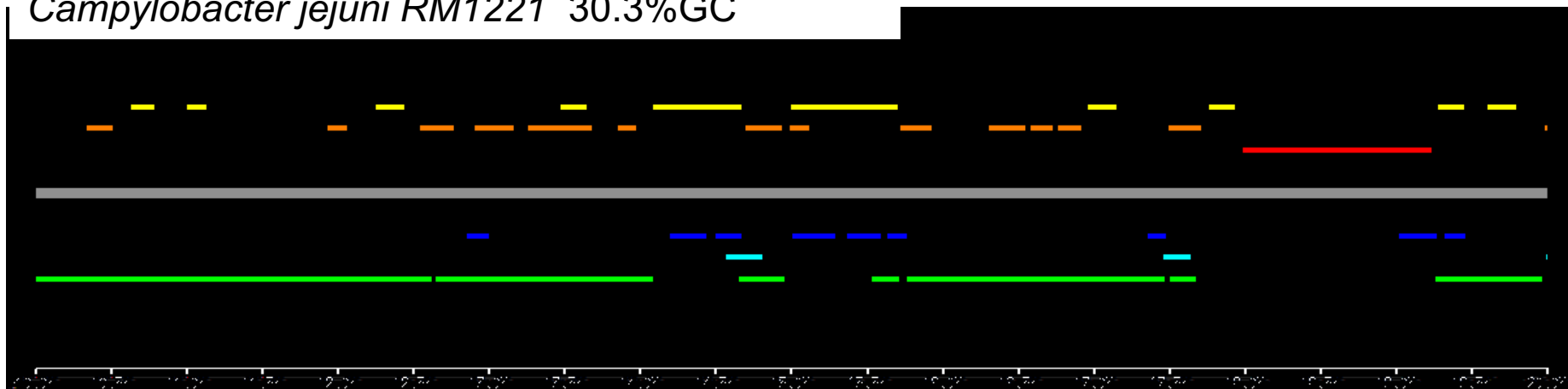
```
Query    116   EFTPAVHASLDKFLASVSTVLTSKY 140
          EFTP V A+  K +A V+  L  KY
Sbjct    121   EFTPPVQAAYQKVVAGVANALAHKY 145
```



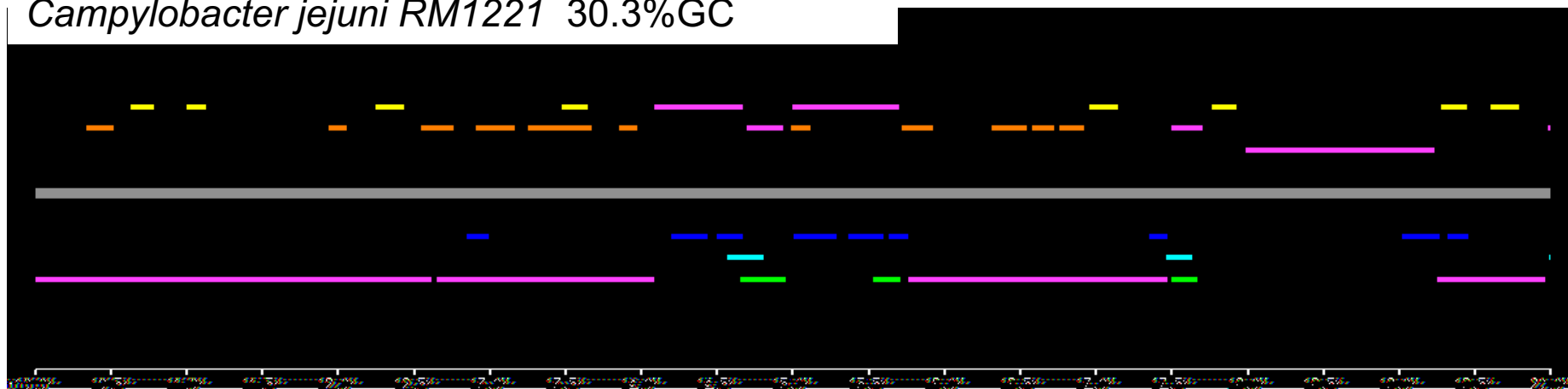

Outline

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

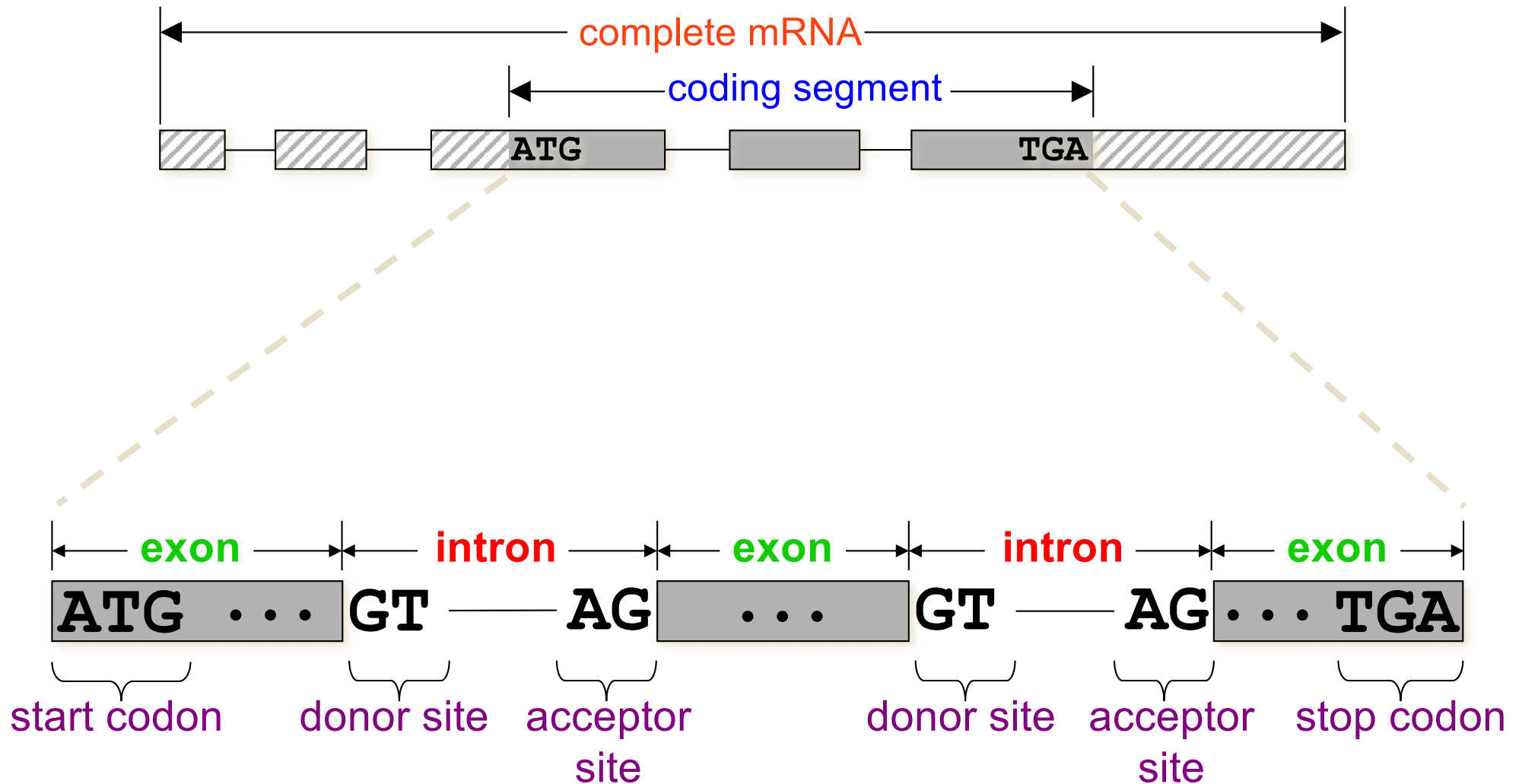
Campylobacter jejuni RM1221 30.3%GC



Campylobacter jejuni RM1221 30.3%GC

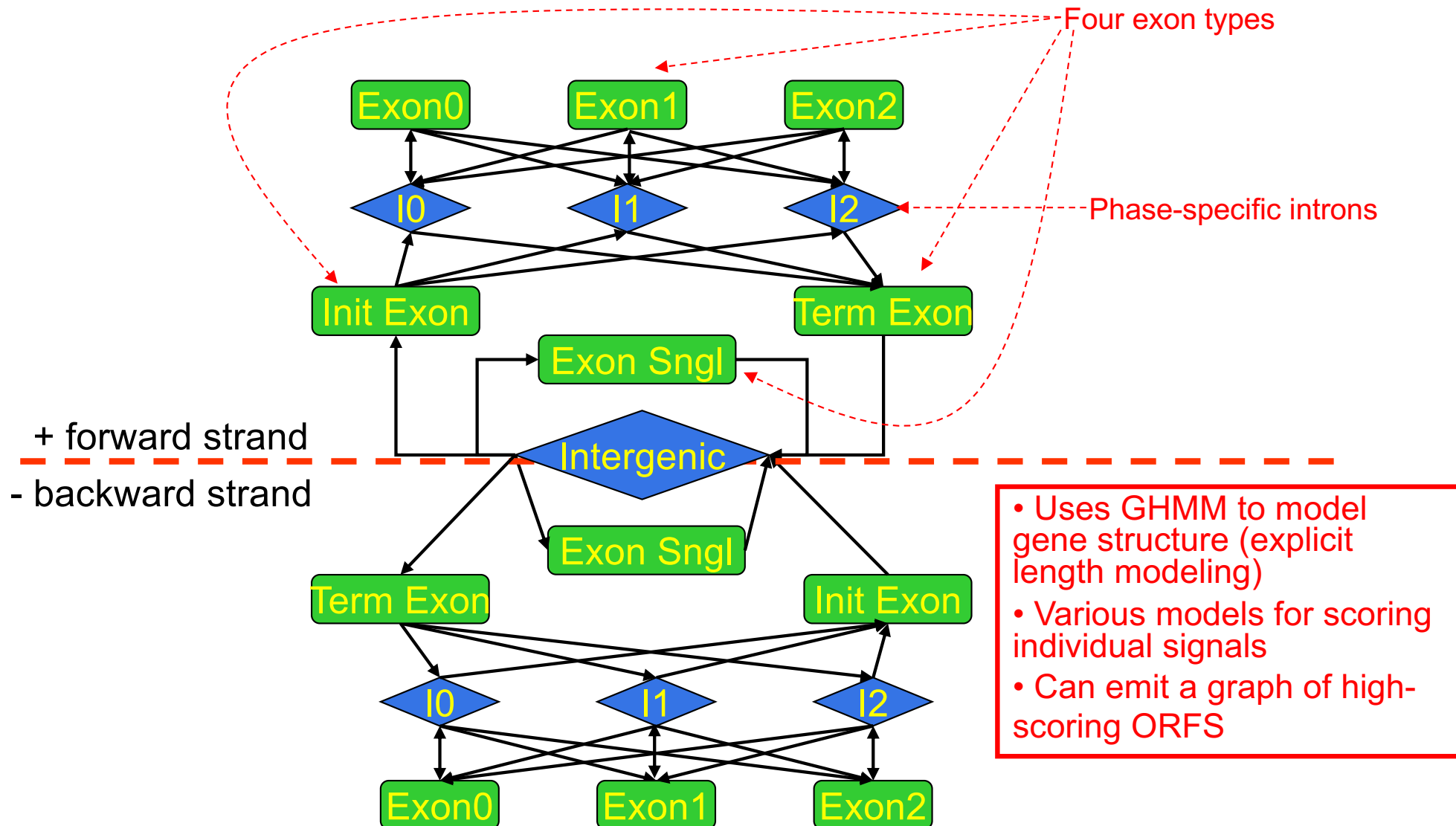


Eukaryotic Gene Syntax



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.

GlimmerHMM architecture



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

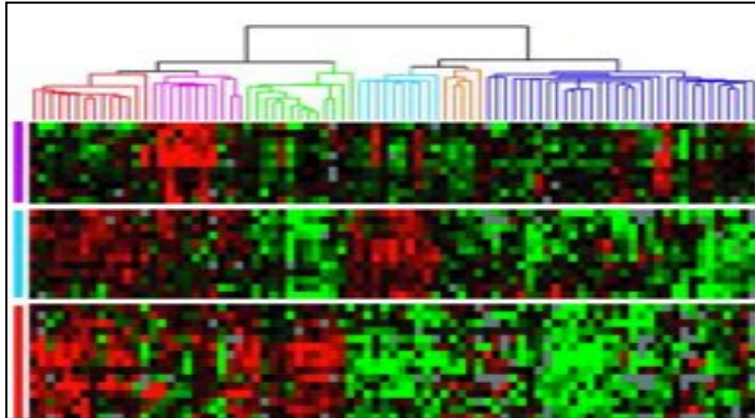


Outline

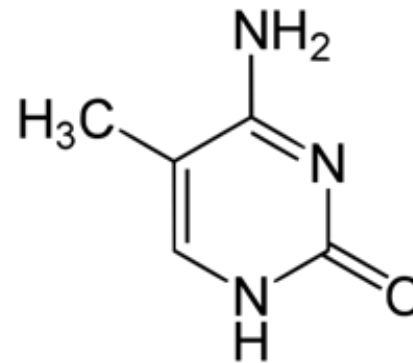
1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. **Experimental & Functional Assays**

*-seq in 4 short vignettes

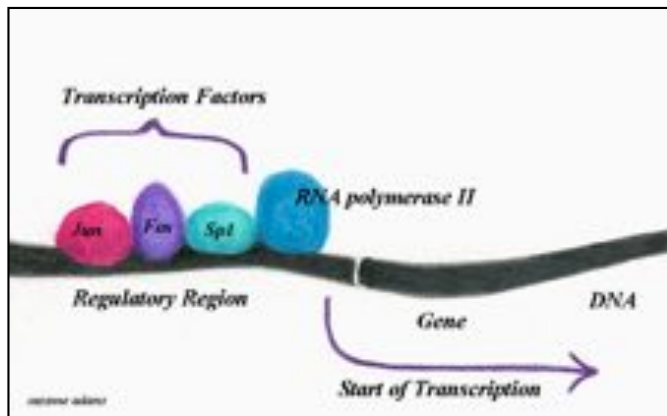
RNA-seq



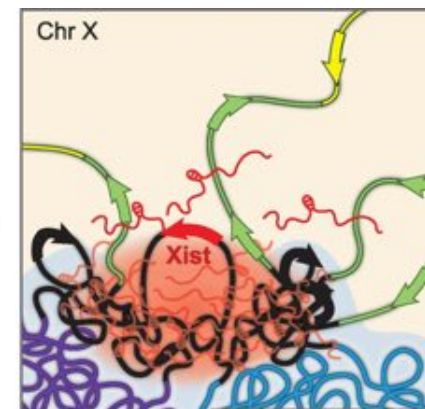
Methyl-seq



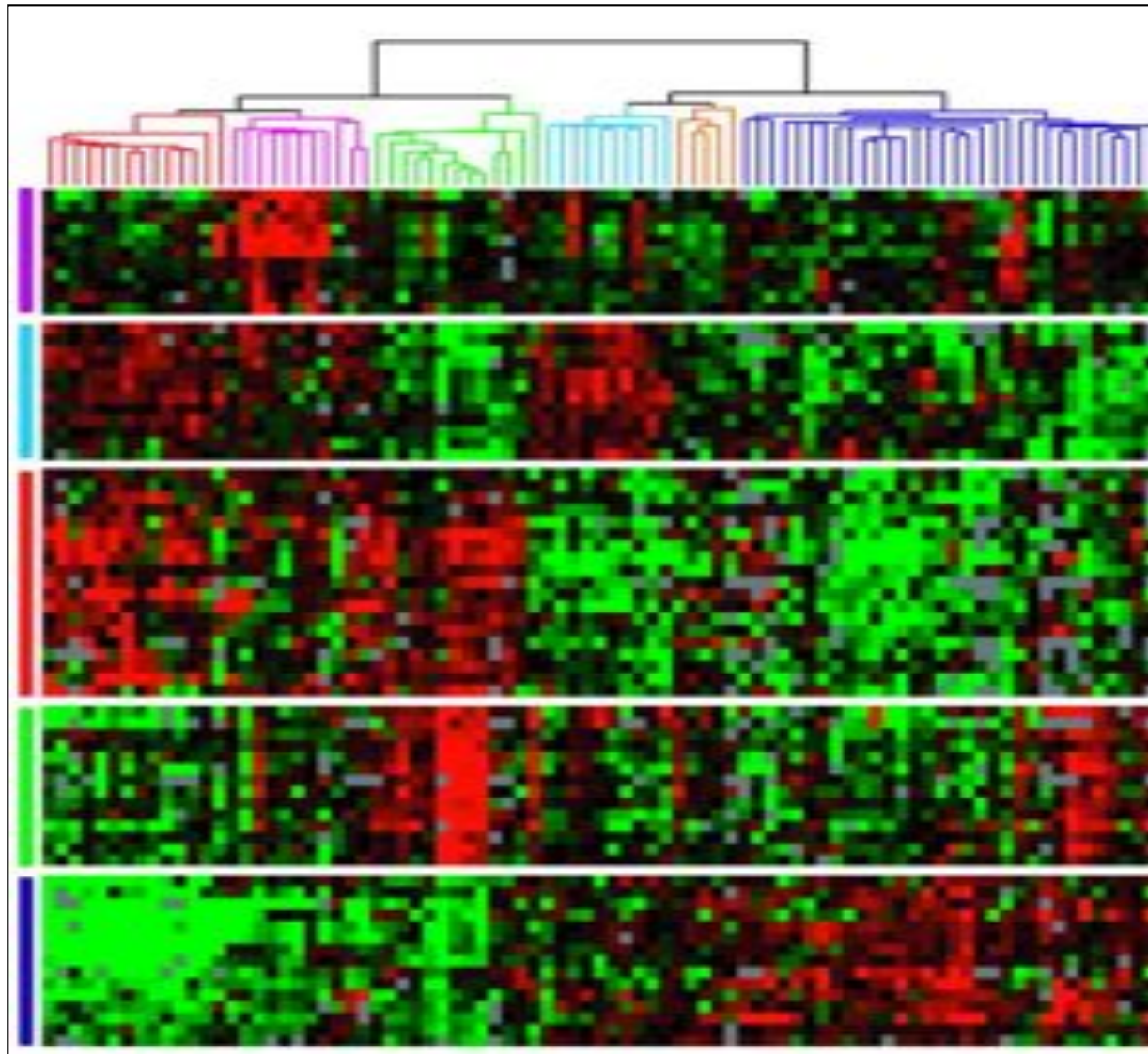
ChIP-seq



Hi-C

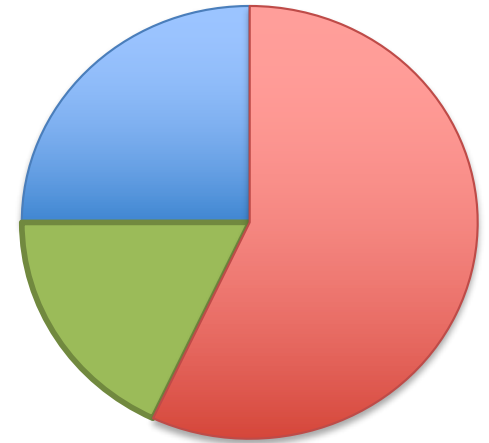
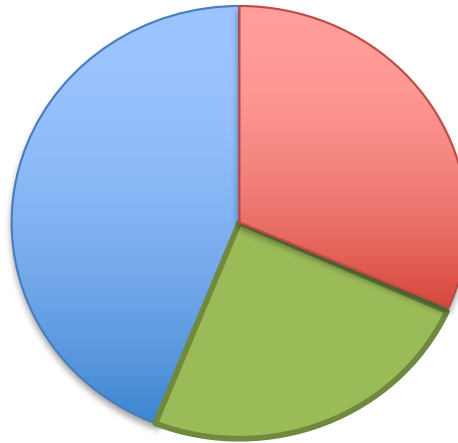
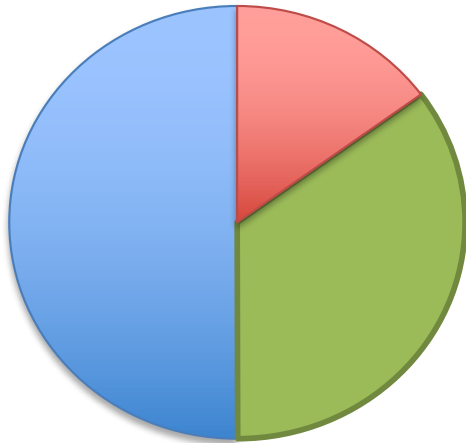
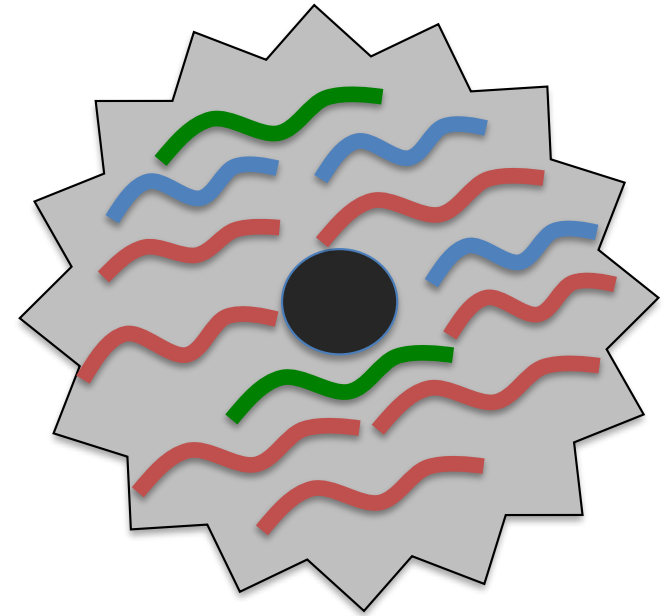
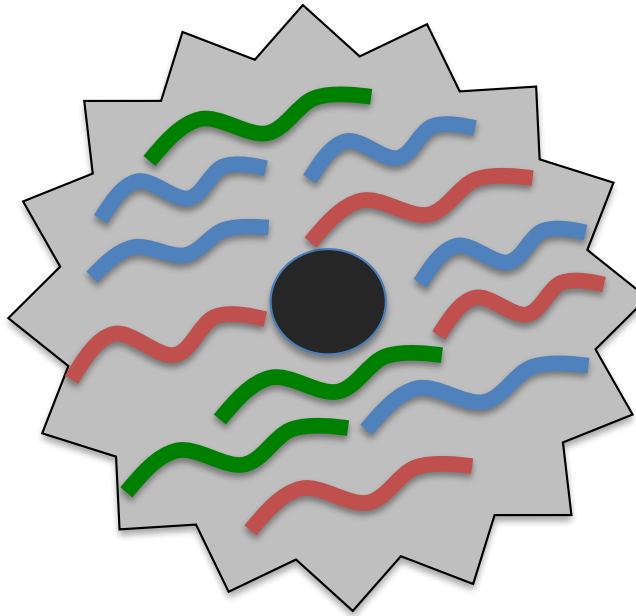
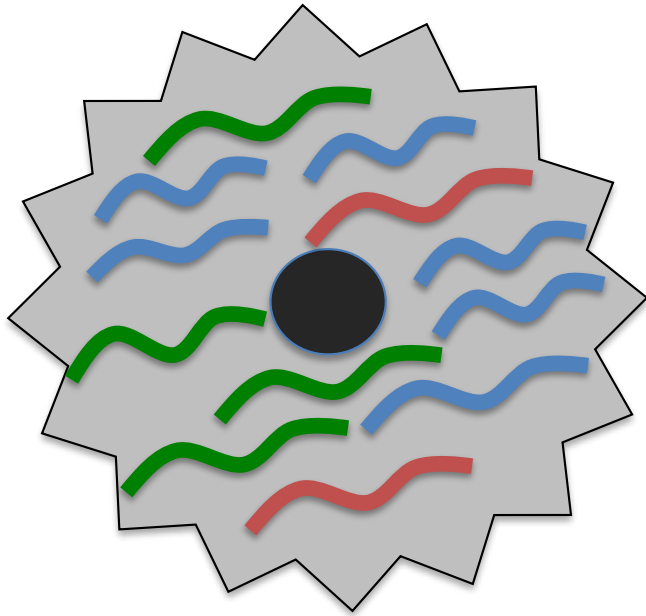


RNA-seq

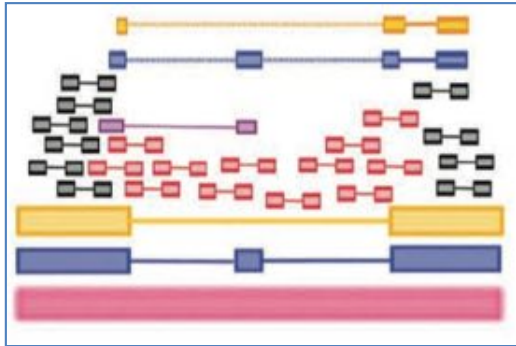


Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications.
Sørli et al (2001) *PNAS*. 98(19):10869-74.

RNA-seq Overview

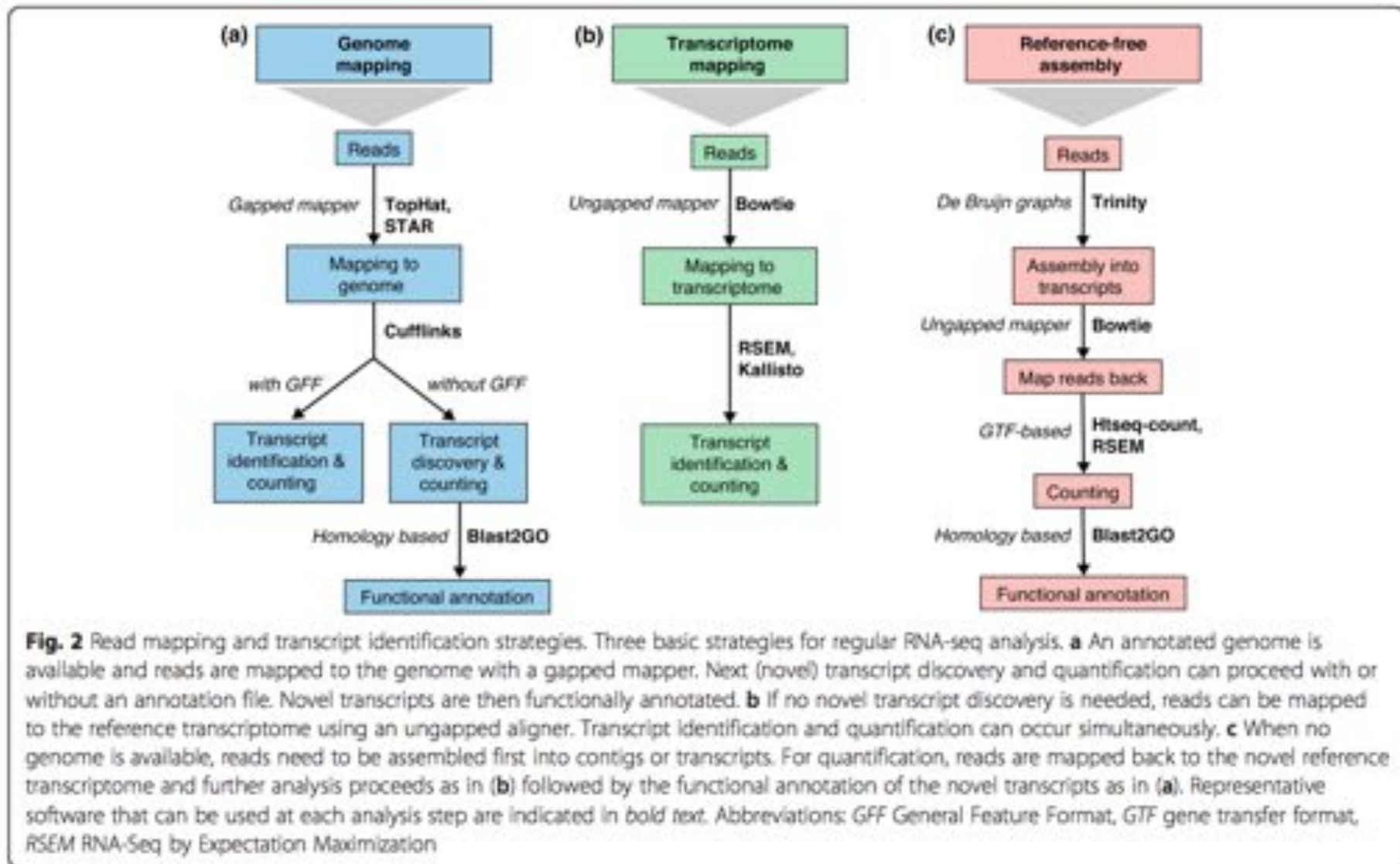


RNA-seq Challenges



Challenge 1: Eukaryotic genes are spliced

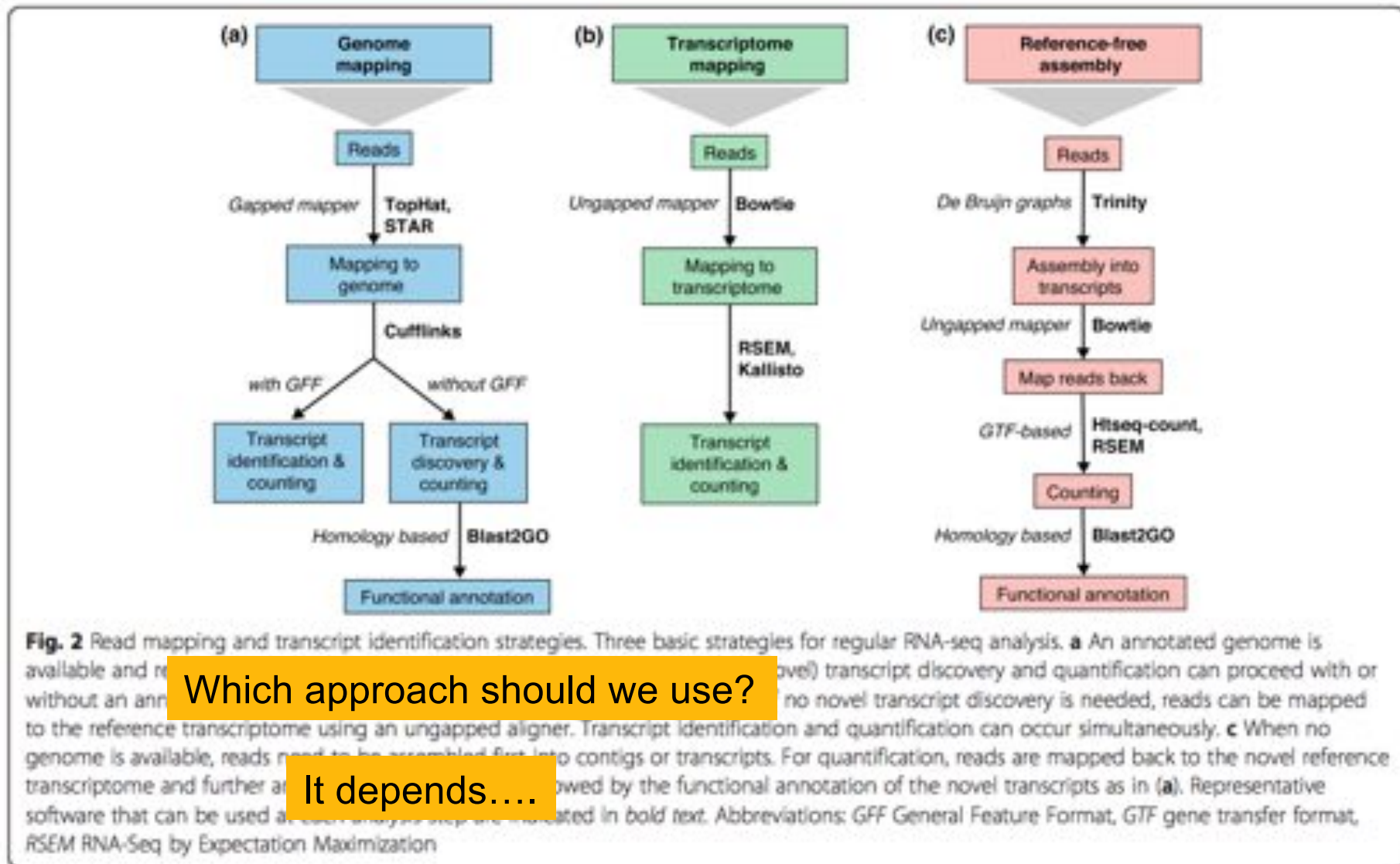
RNA-Seq Approaches



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



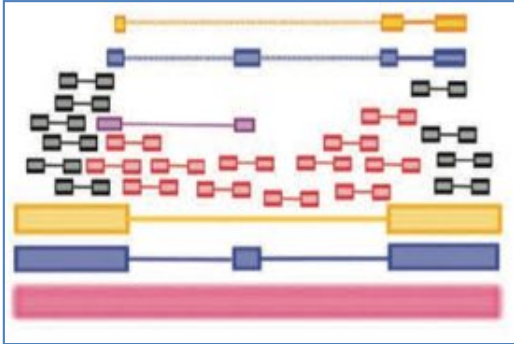
Which approach should we use?

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

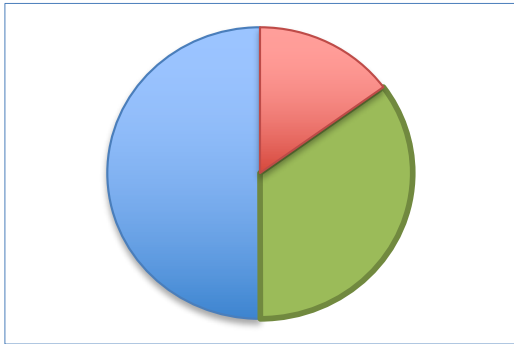


Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

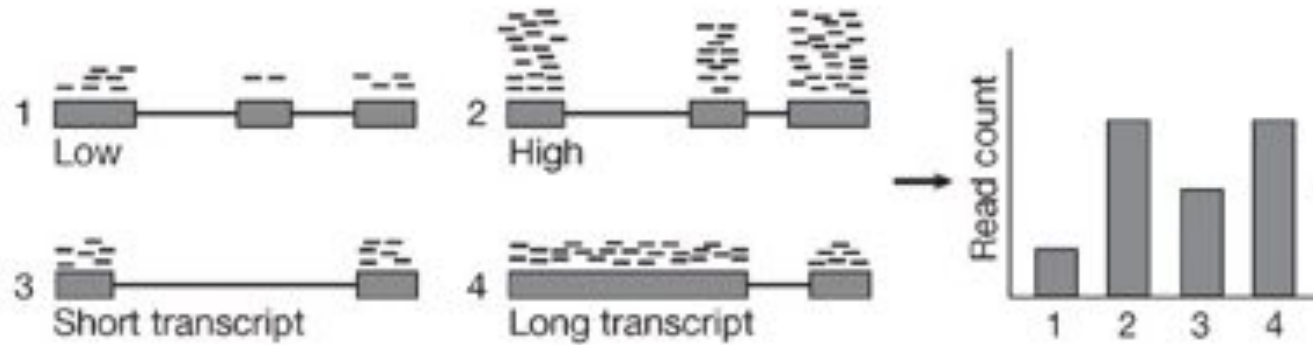
TopHat: discovering spliced junctions with RNA-Seq.

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



Challenge 2: Read Count != Transcript abundance

RPKM, FPKM, TPM

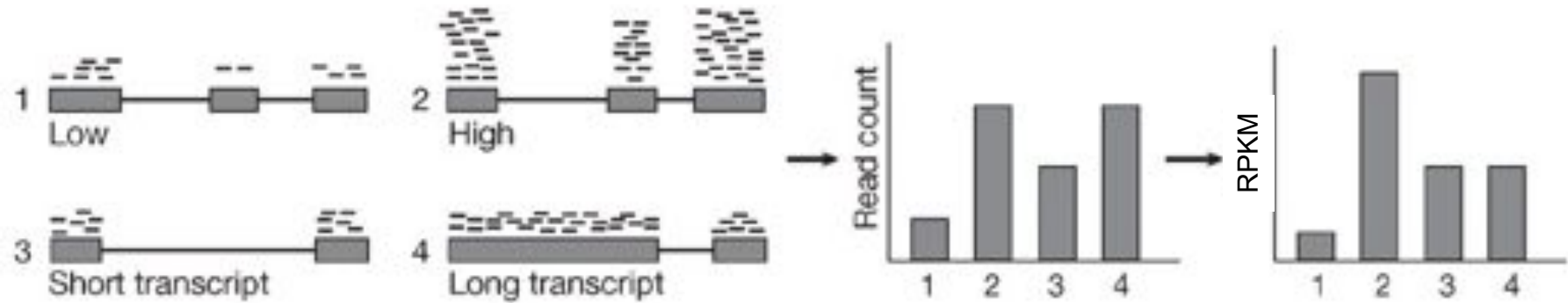


Counting Reads that align to a gene DOESN'T work!

- 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!

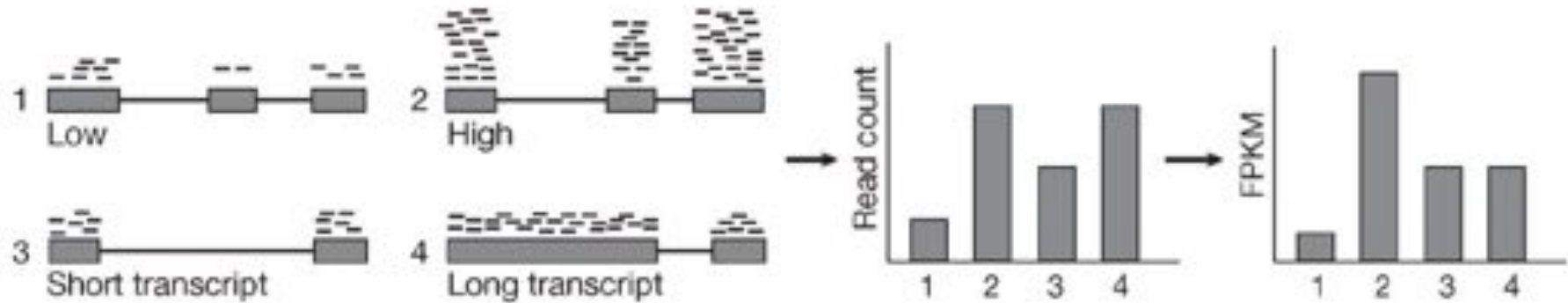
- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

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)

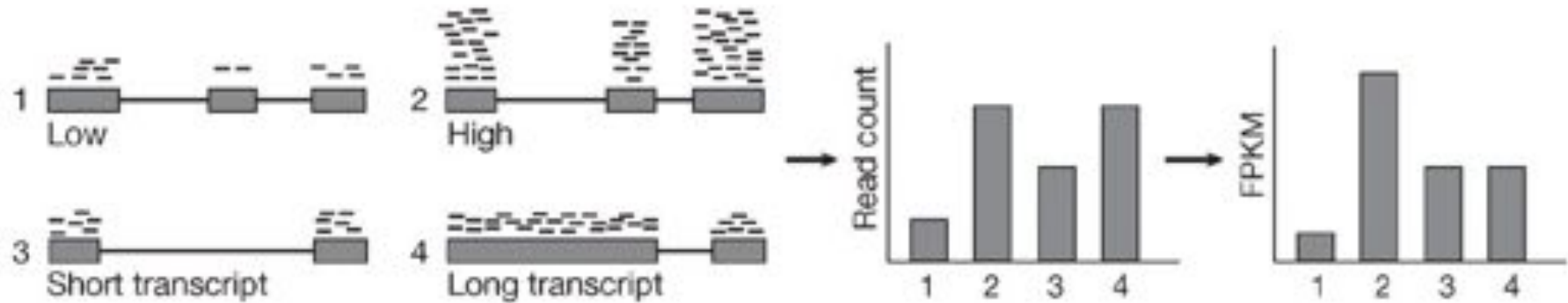
⇒ Wait a second, reads in a pair are not independent!

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)

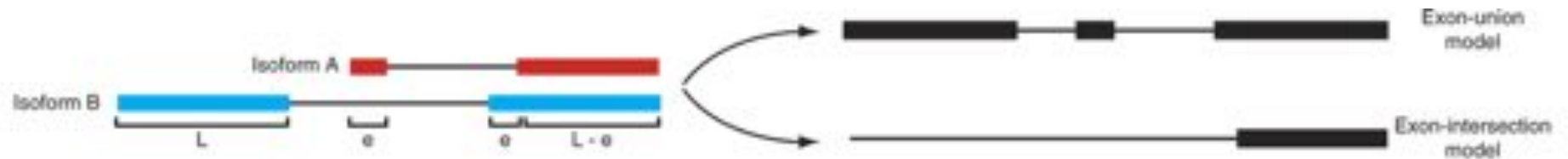
=> If you were to sequence one million full length transcripts, TPM is the number of transcripts you would have seen of type i , given the abundances of the other transcripts in your sample

=> 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?

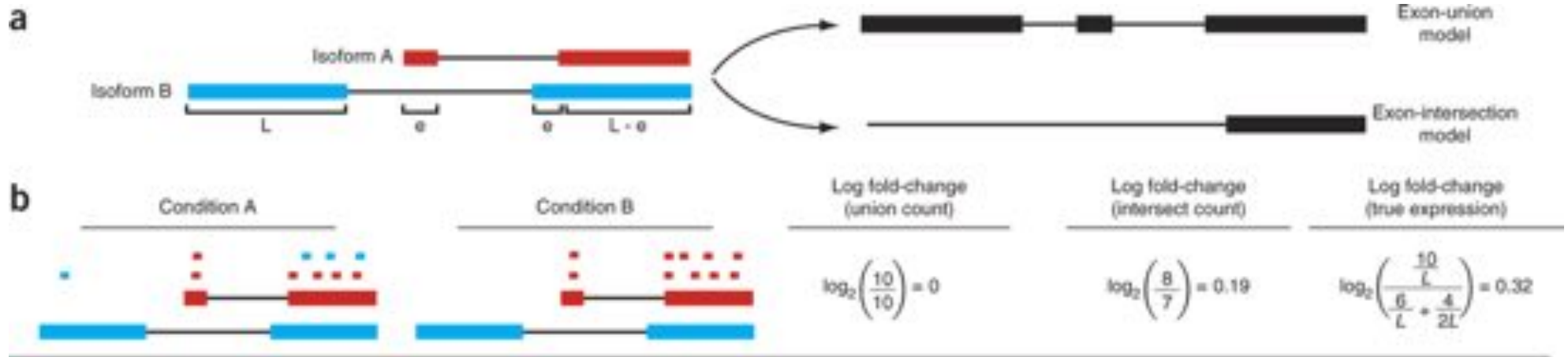
a



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

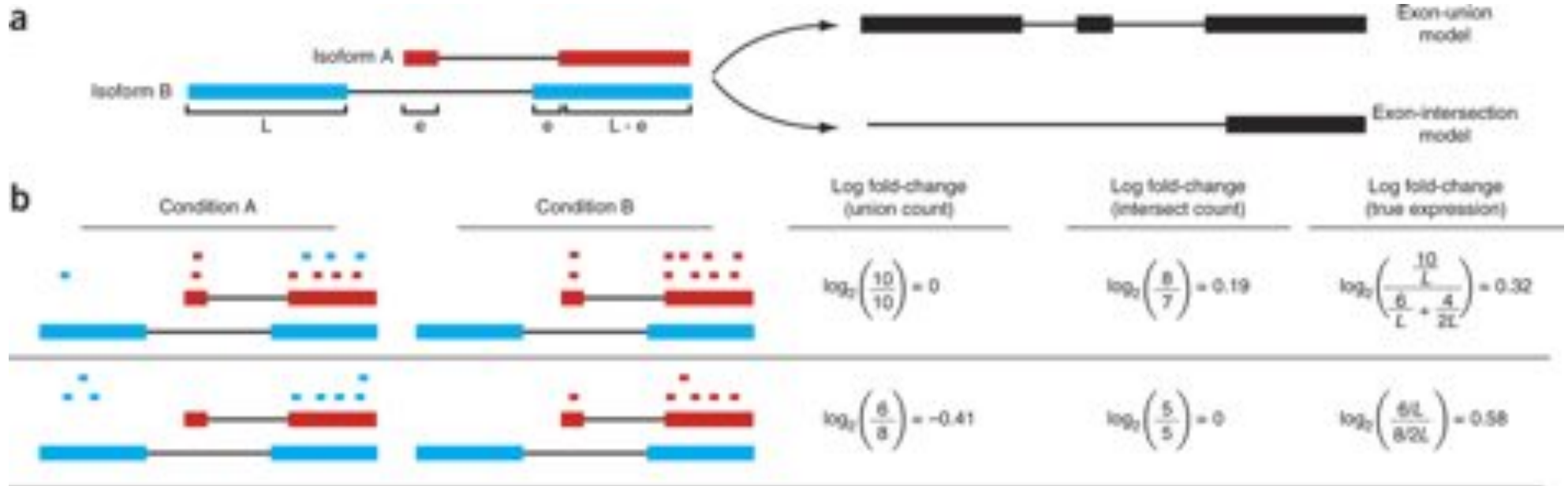
Gene or Isoform Quantification?



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

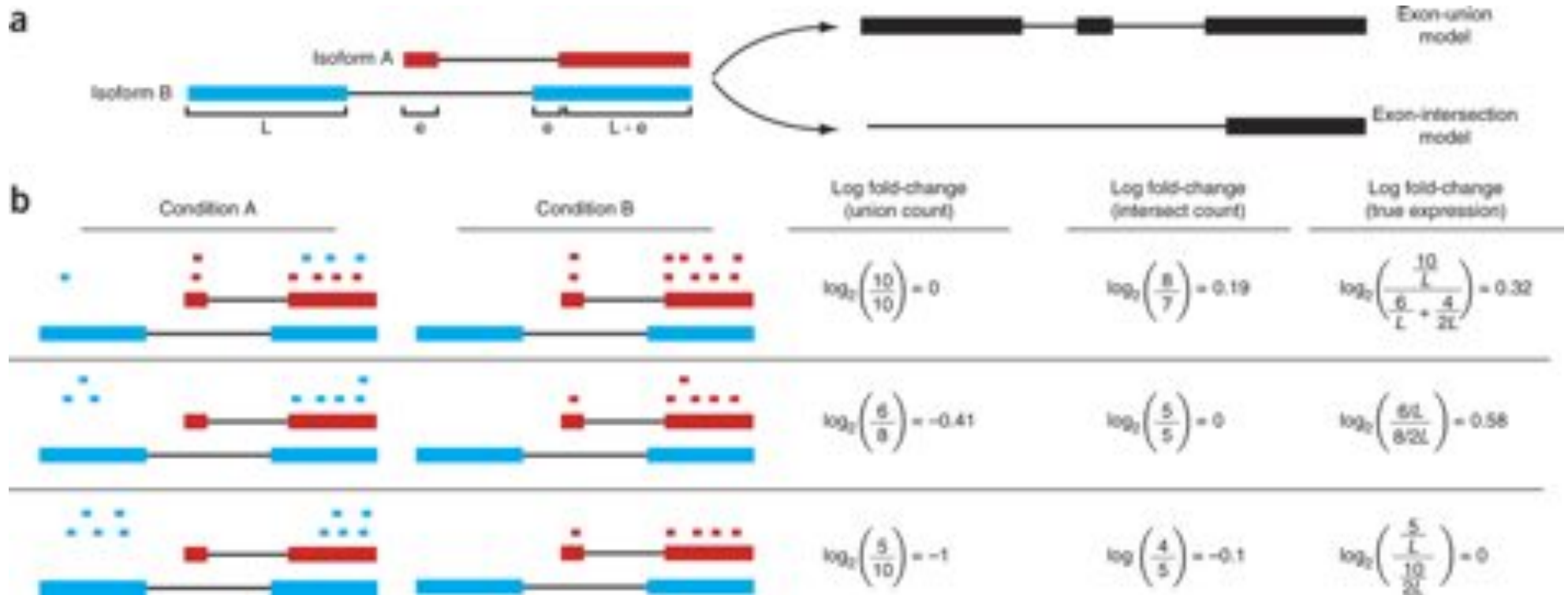
Gene or Isoform Quantification?



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?



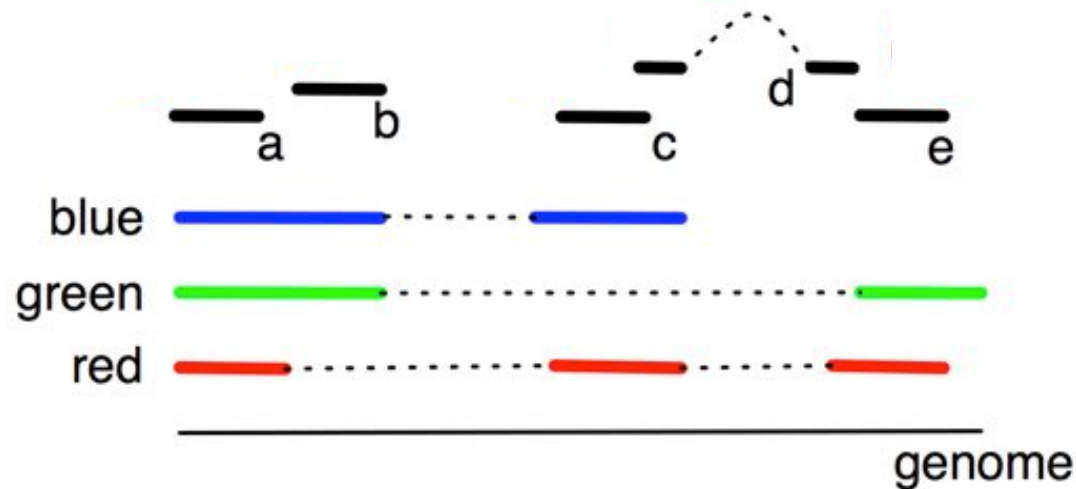
Key point : The length of the actual molecule from which the fragments derive is crucially important to obtaining accurate abundance estimates.

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.

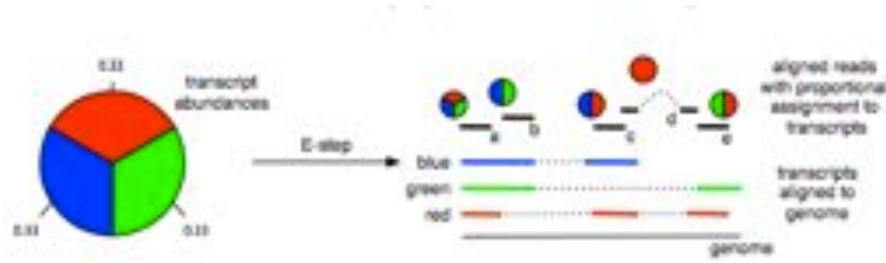
What is the most likely expression level of each isoform?

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



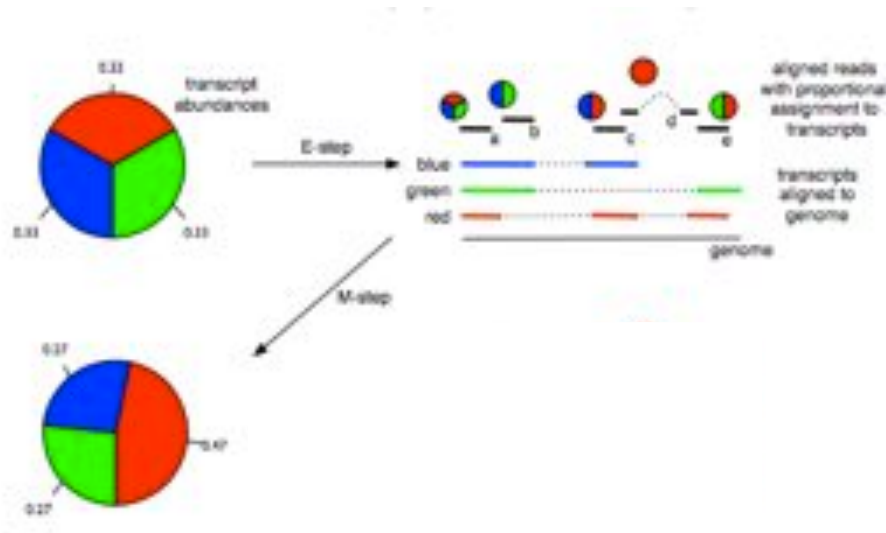
The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

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:

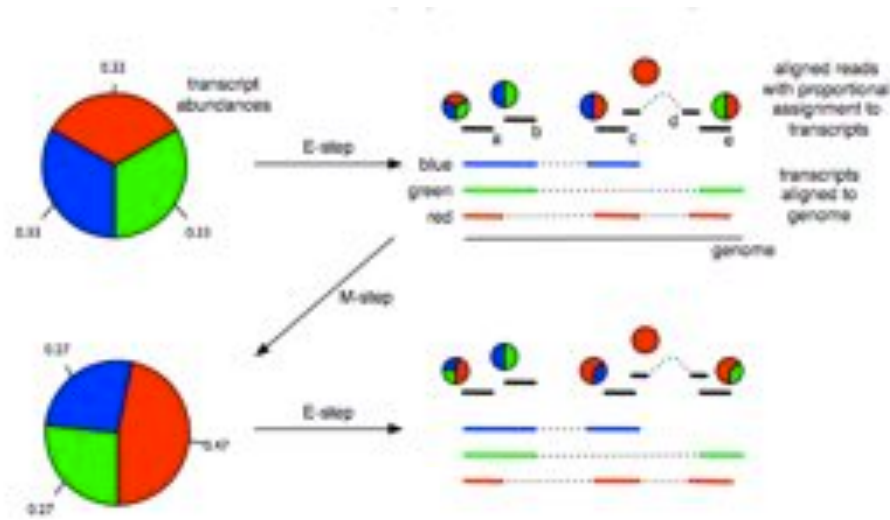
$$\text{red: } 0.47 = (0.33 + 0.5 + 1 + 0.5) / (2.33 + 1.33 + 1.33)$$

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

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

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:

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: $0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$

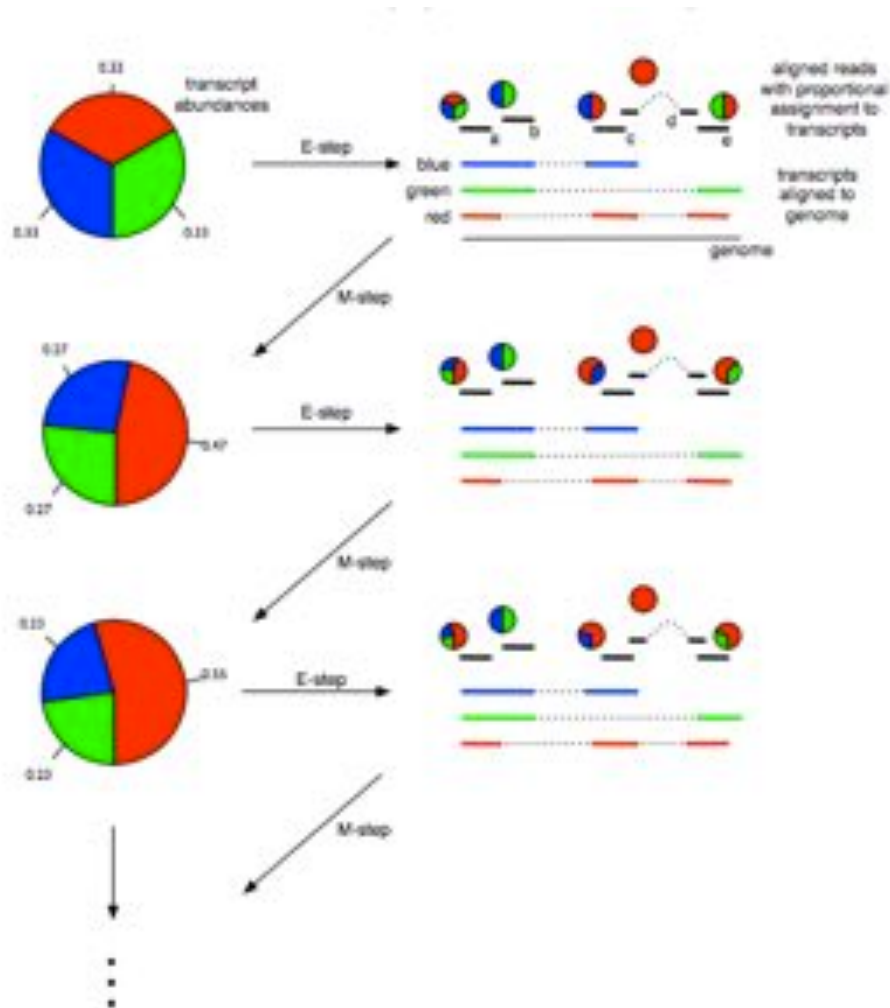
Repeat until convergence!

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:

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

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

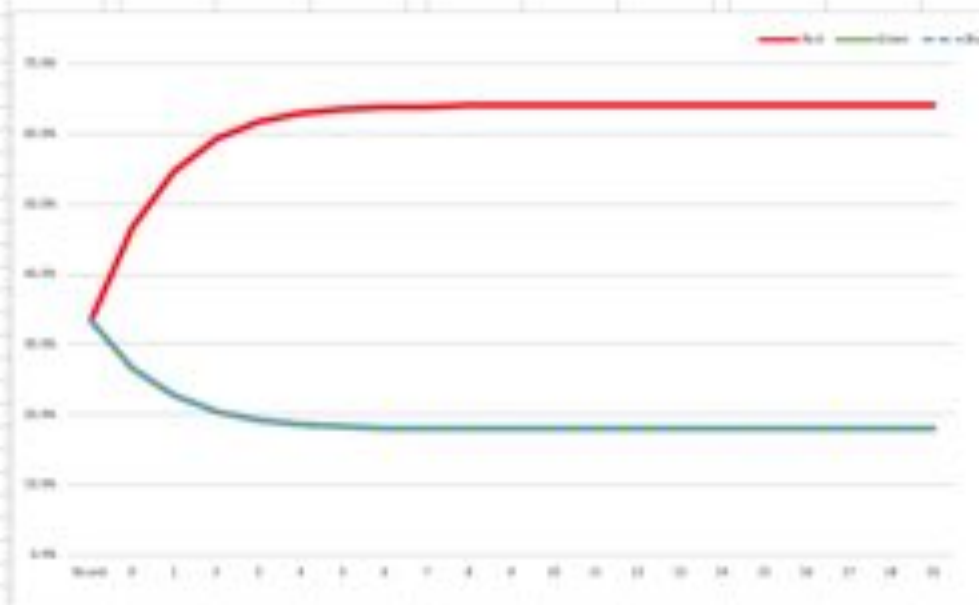
$$\text{green: } 0.27 = (0.33 + 0.5 + 0.5) / (2.33 + 1.33 + 1.33)$$

Repeat until convergence!

Models for transcript quantification from RNA-seq

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

Round	Transcript			Read A			Read B			Read C			Read D			Read E		
	Red	Green	Blue	Red	Agreen	ABlue	Red	Bgreen	BBlue	Red	Cgreen	CBlue	Red	Dgreen	DBlue	Red	Egreen	EBlue
0	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%	0.0%	50.0%	50.0%	50.0%	0.0%	0.0%	0.0%	100.0%	0.0%	0.0%	50.0%	50.0%
1	46.7%	26.7%	26.7%	46.7%	26.7%	26.7%	0.0%	50.0%	50.0%	43.3%	0.0%	56.7%	100.0%	0.0%	0.0%	43.3%	56.7%	0.0%
2	54.8%	22.6%	22.6%	54.8%	22.6%	22.6%	0.0%	50.0%	50.0%	70.8%	0.0%	29.2%	100.0%	0.0%	0.0%	70.8%	29.2%	0.0%
3	59.3%	20.4%	20.4%	59.3%	20.4%	20.4%	0.0%	50.0%	50.0%	76.4%	0.0%	23.6%	100.0%	0.0%	0.0%	76.4%	23.6%	0.0%
4	61.6%	19.2%	19.2%	61.6%	19.2%	19.2%	0.0%	50.0%	50.0%	78.9%	0.0%	21.1%	100.0%	0.0%	0.0%	78.9%	21.1%	0.0%
5	62.8%	18.8%	18.8%	62.8%	18.8%	18.8%	0.0%	50.0%	50.0%	77.2%	0.0%	22.8%	100.0%	0.0%	0.0%	77.2%	22.8%	0.0%
6	63.4%	18.3%	18.3%	63.4%	18.3%	18.3%	0.0%	50.0%	50.0%	77.6%	0.0%	22.4%	100.0%	0.0%	0.0%	77.6%	22.4%	0.0%
7	63.7%	18.1%	18.1%	63.7%	18.1%	18.1%	0.0%	50.0%	50.0%	77.9%	0.0%	22.1%	100.0%	0.0%	0.0%	77.9%	22.1%	0.0%
8	63.9%	18.1%	18.1%	63.9%	18.1%	18.1%	0.0%	50.0%	50.0%	78.0%	0.0%	22.0%	100.0%	0.0%	0.0%	78.0%	22.0%	0.0%
9	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.0%	0.0%	22.0%	100.0%	0.0%	0.0%	78.0%	22.0%	0.0%
10	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.0%	0.0%	22.0%	100.0%	0.0%	0.0%	78.0%	22.0%	0.0%
11	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
12	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
13	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
14	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
15	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
16	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
17	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
18	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
19	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%
20	64.0%	18.0%	18.0%	64.0%	18.0%	18.0%	0.0%	50.0%	50.0%	78.1%	0.0%	21.9%	100.0%	0.0%	0.0%	78.1%	21.9%	0.0%

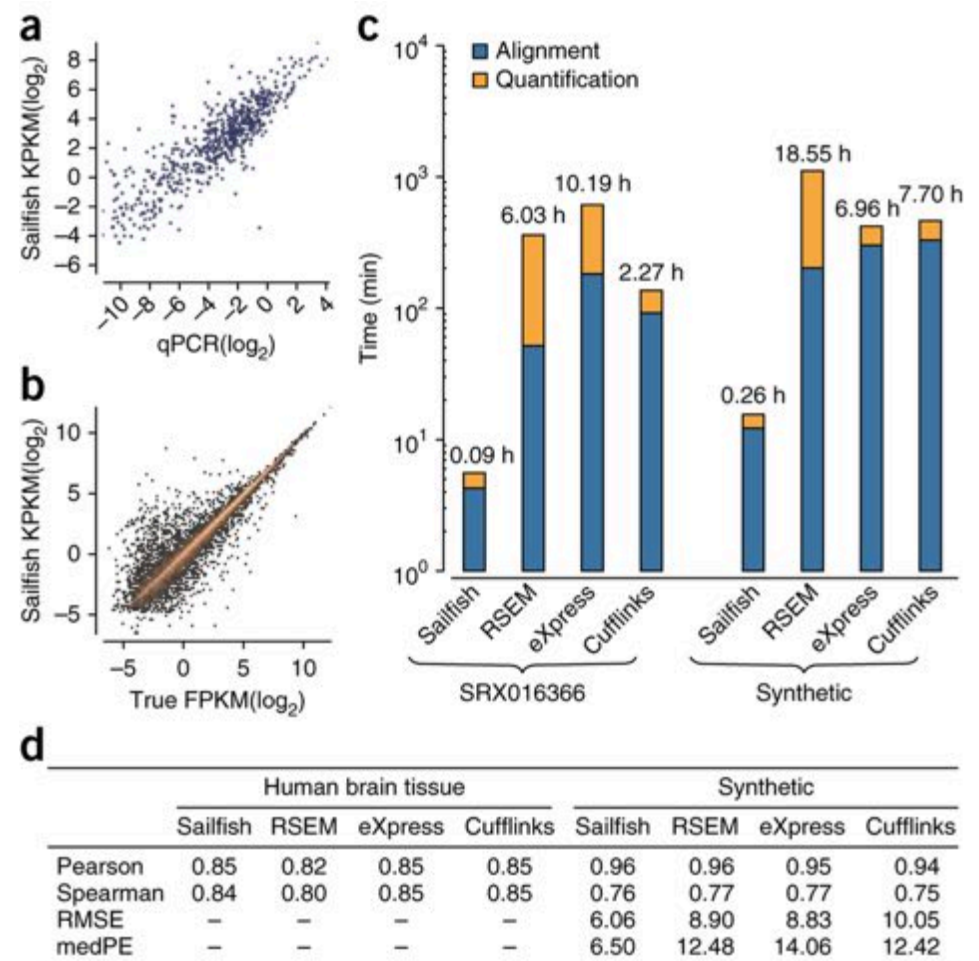
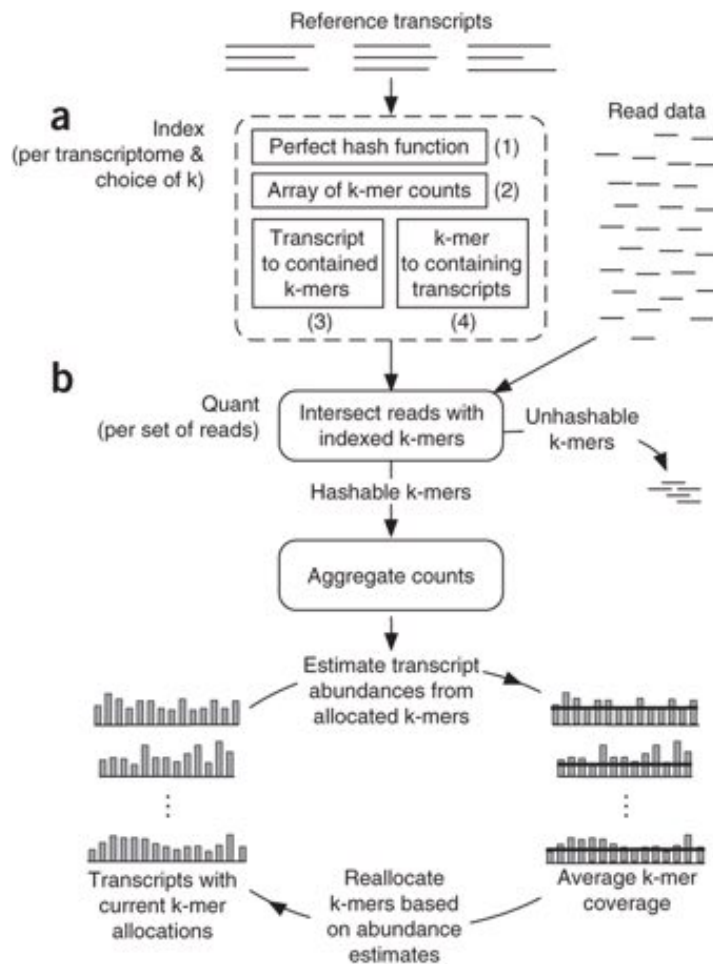
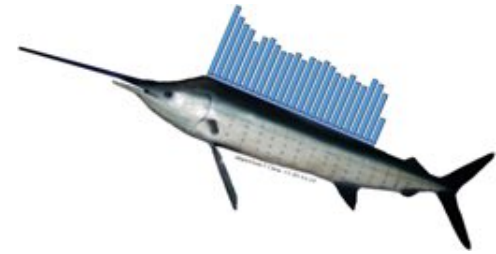


RNAseqExpectationMaximization.xlsx

Models for transcript quantification from RNA-seq

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

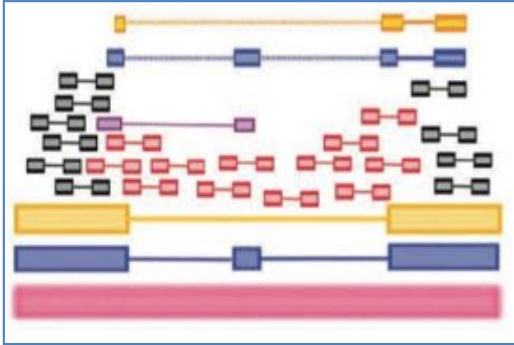
Sailfish: Fast & Accurate RNA-seq Quantification



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

RNA-seq Challenges

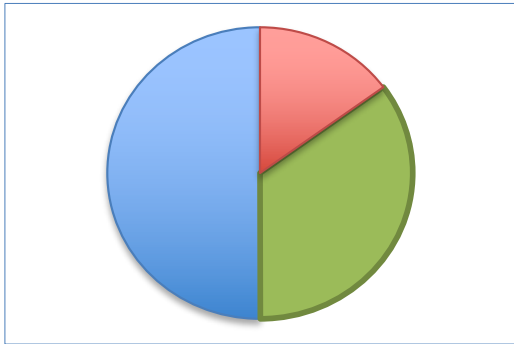


Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

TopHat: discovering spliced junctions with RNA-Seq.

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

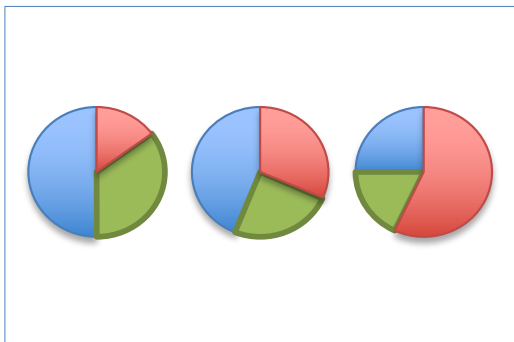


Challenge 2: Read Count \neq Transcript abundance

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

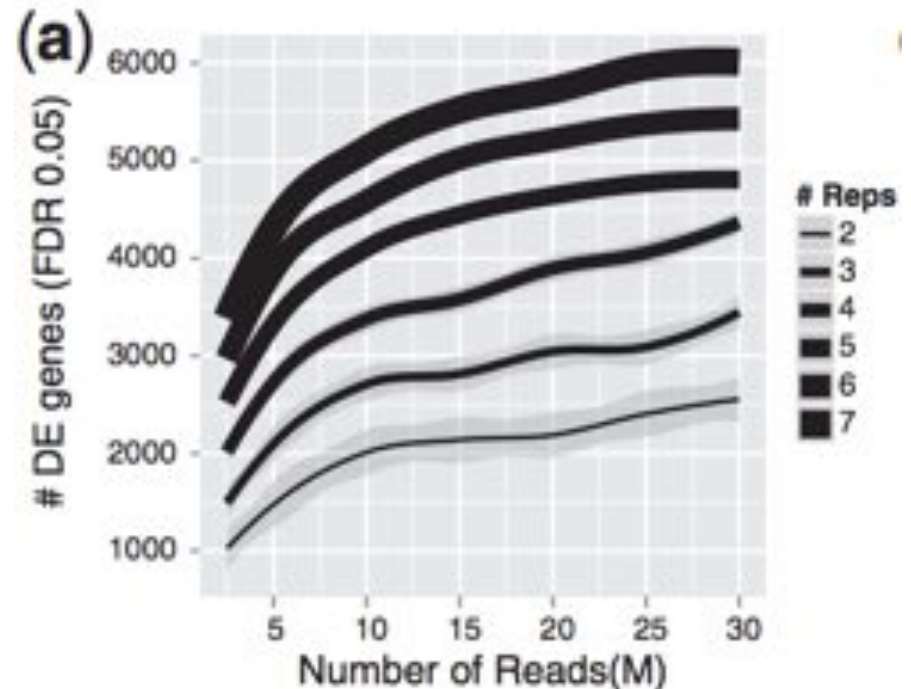
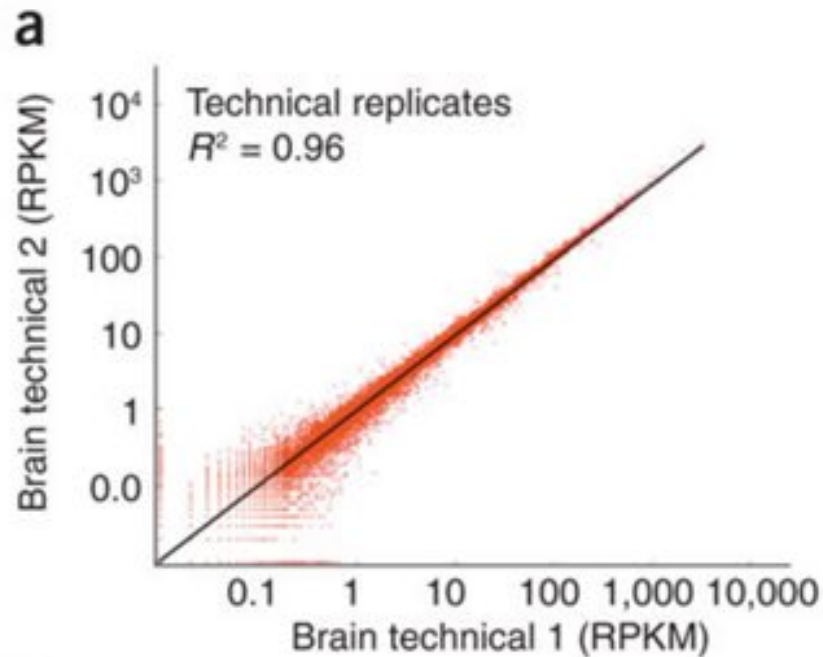
Transcript assembly and quantification by RNA-seq

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



Challenge 3: Transcript abundances are stochastic

How Many Replicates?



Why don't we have perfect replicates?

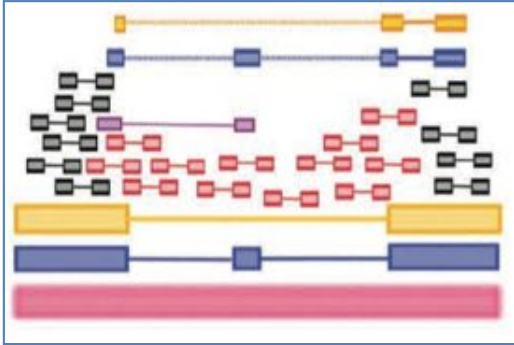
Mapping and quantifying mammalian transcriptomes by RNA-Seq

Mortazavi et al (2008) Nature Methods. 5, 62-628

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

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

RNA-seq Challenges

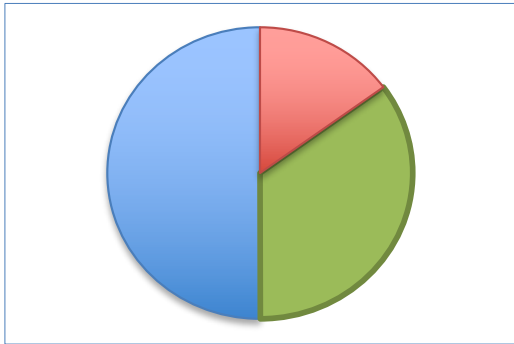


Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

TopHat: discovering spliced junctions with RNA-Seq.

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

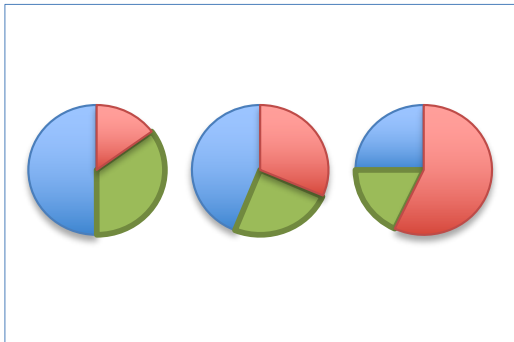


Challenge 2: Read Count != Transcript abundance

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

Transcript assembly and quantification by RNA-seq

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



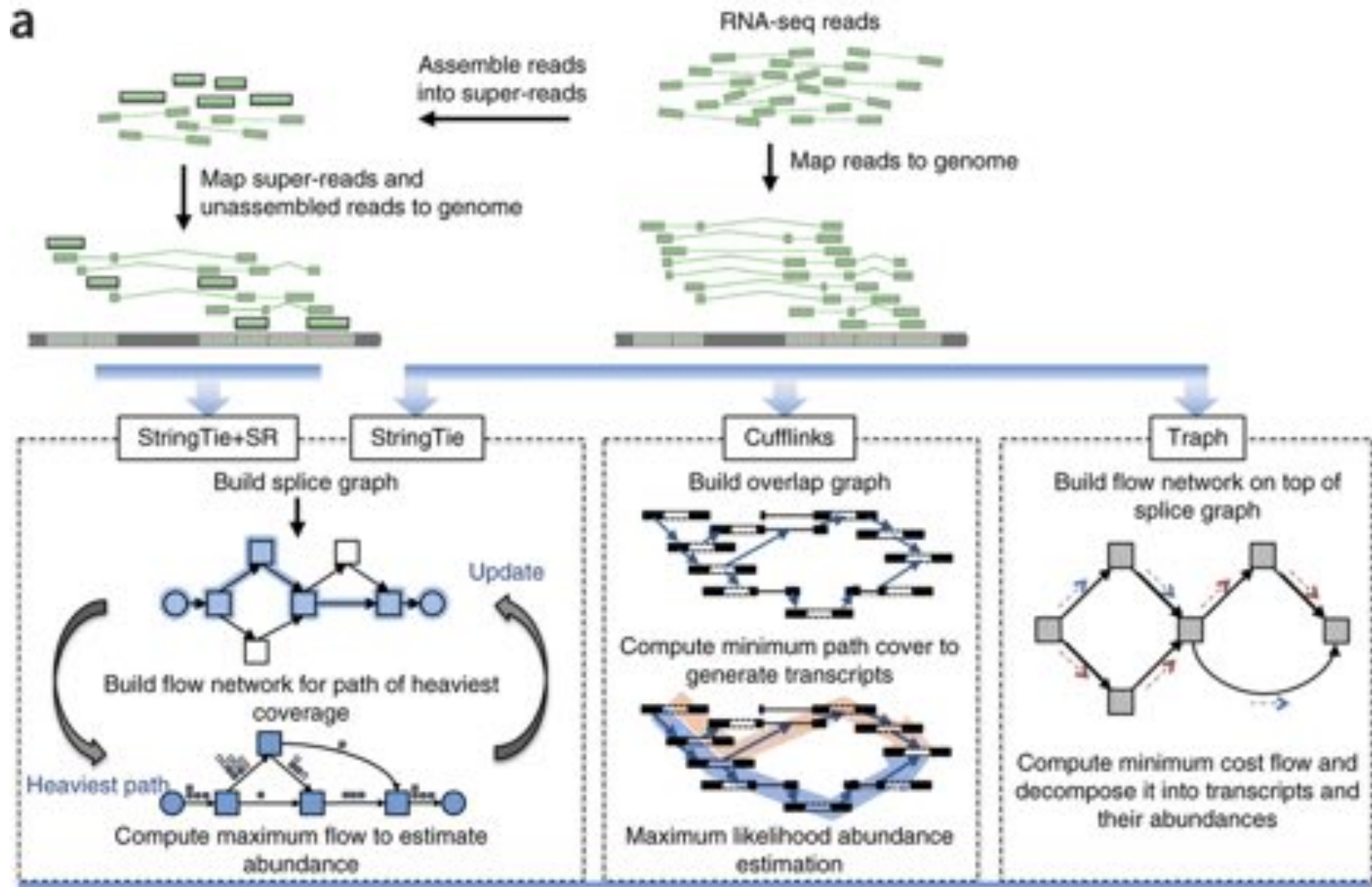
Challenge 3: Transcript abundances are stochastic

Solution: Replicates, replicates, and more replicates

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

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

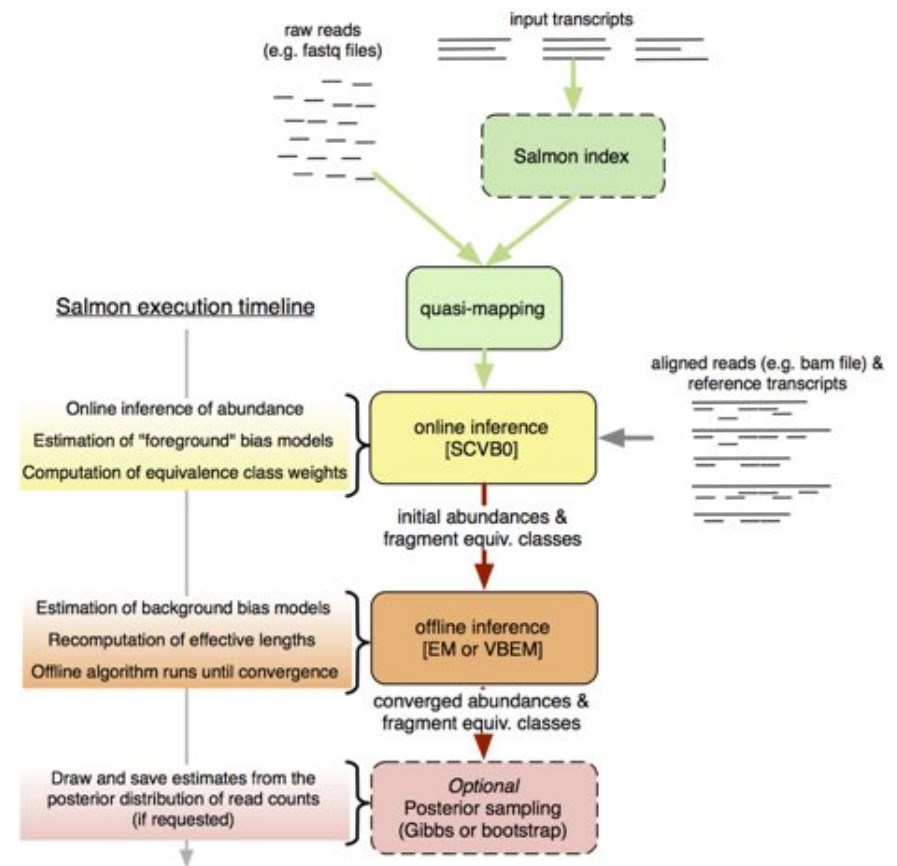
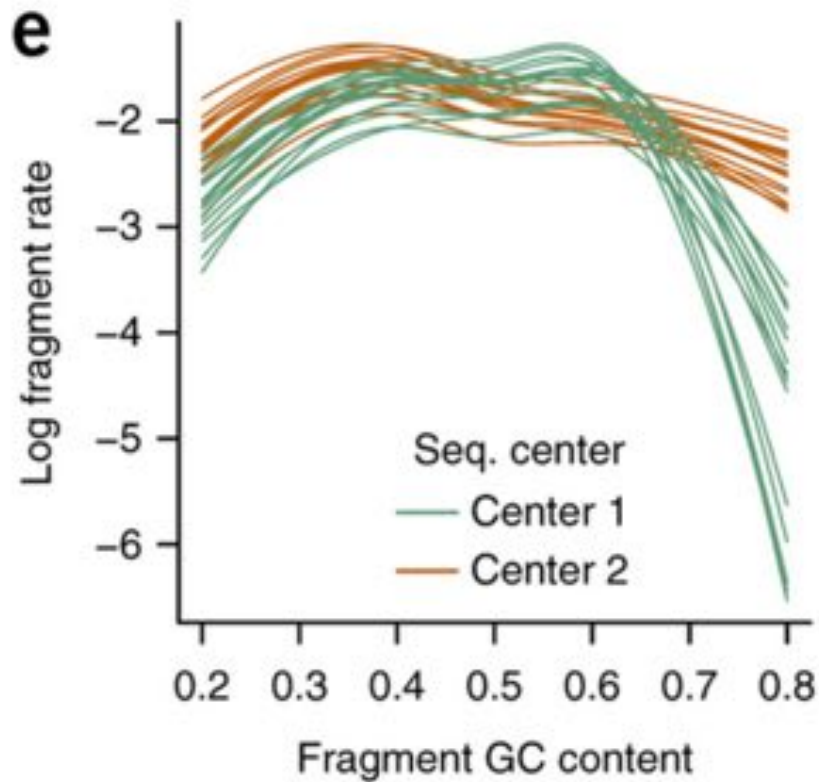
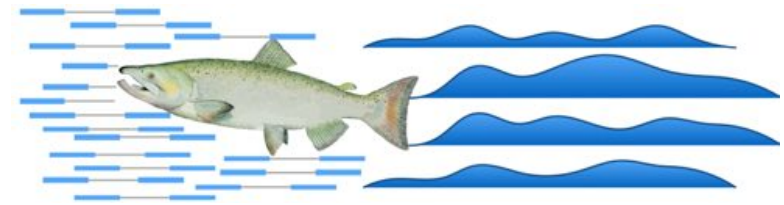
Isoform Quantification Approaches



StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.

Pertea M, et al. (2015) Nature Biotechnology. doi: 10.1038/nbt.3122.

Salmon: The ultimate RNA-seq Pipeline?



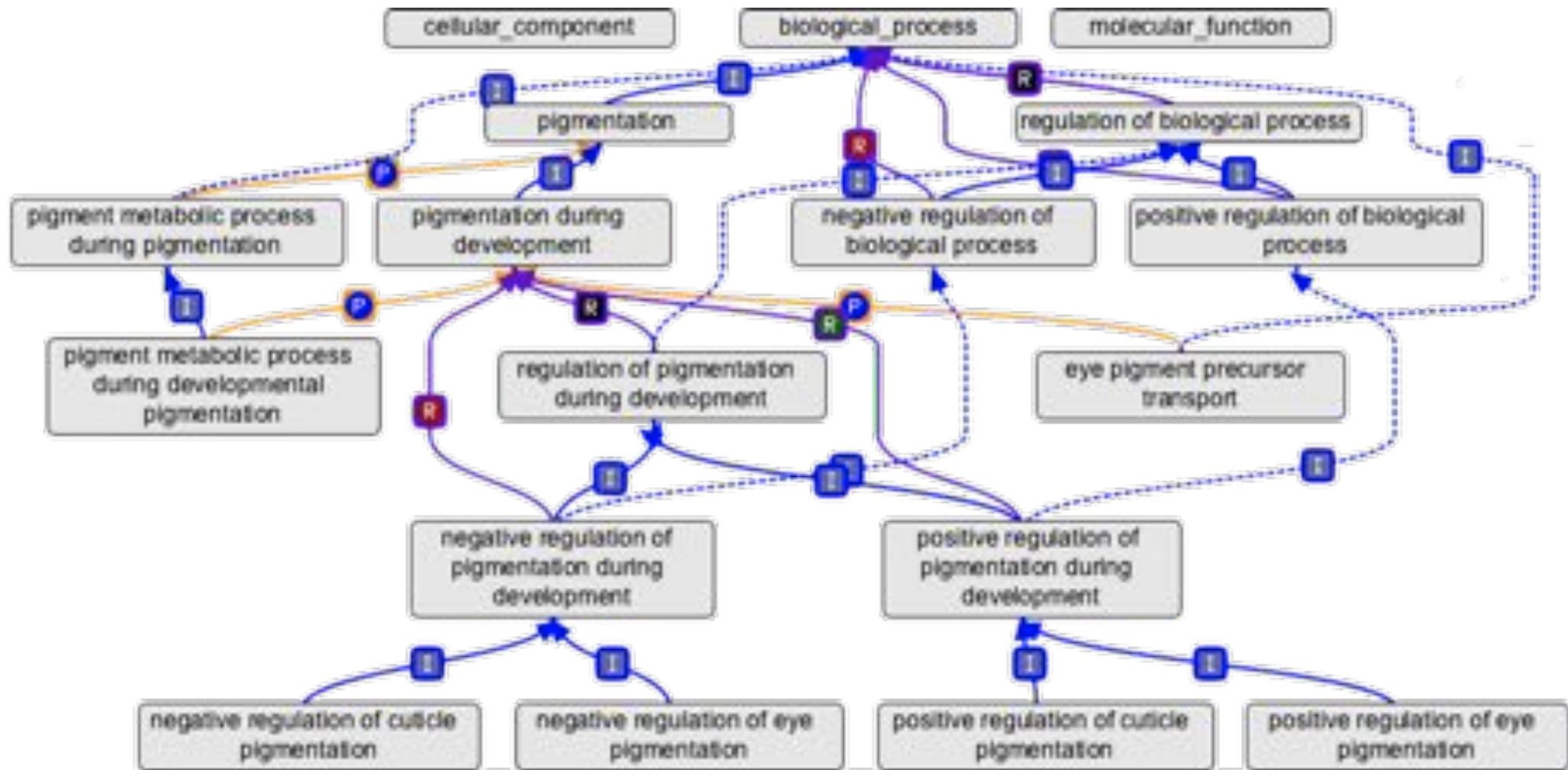
Modeling of RNA-seq fragment sequence bias reduces systematic errors in transcript abundance estimation

Love et al (2016) Nature Biotechnology 34, 1287–1291 (2016) doi:10.1038/nbt.3682

Salmon provides fast and bias-aware quantification of transcript expression

Patro et al (2017) Nature Methods (2017) doi:10.1038/nmeth.4197

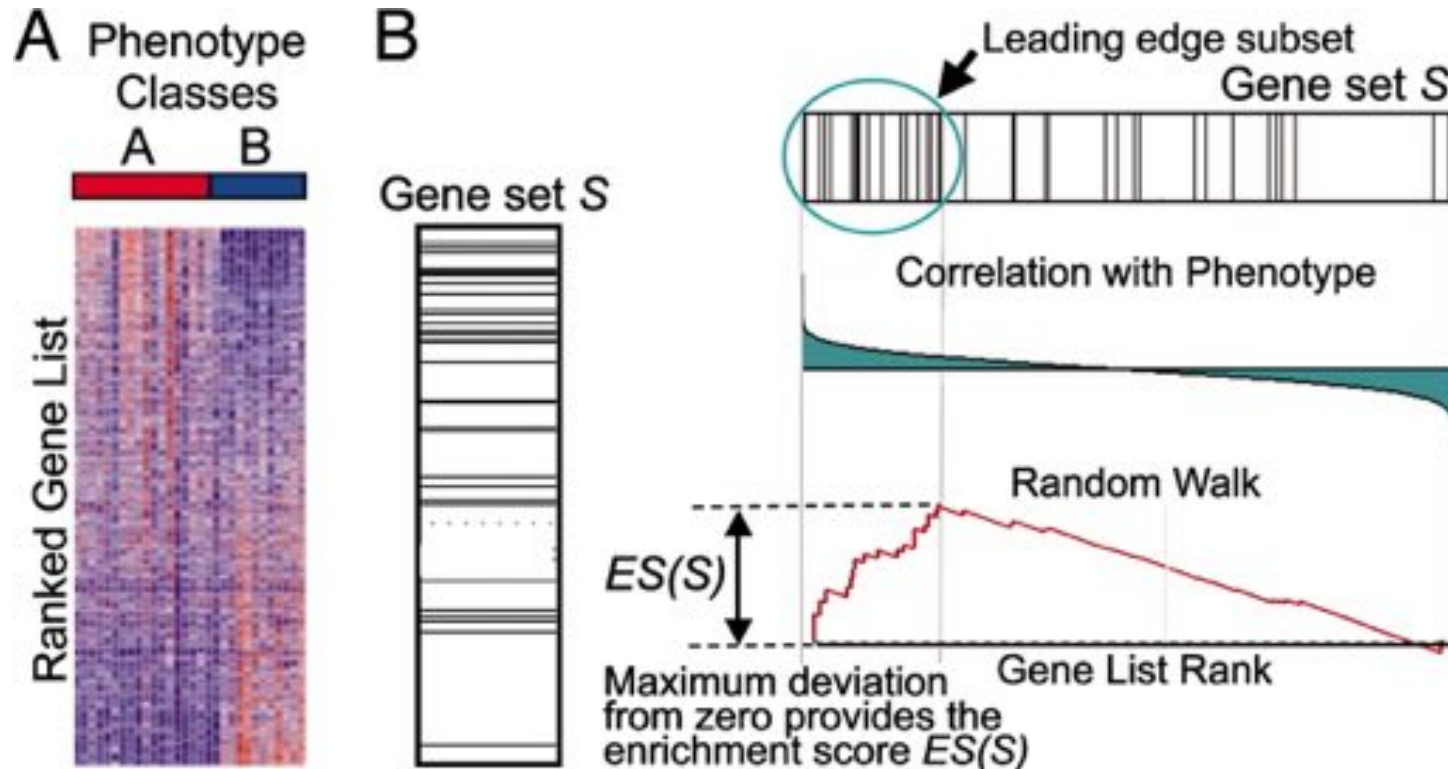
Gene Ontology (GO)



AmiGO: online access to ontology and annotation data

Carbon et al (2009) Bioinformatics doi:10.1093/bioinformatics/btn615

GSEA Overview



Collections

The MSigDB gene sets are divided into 8 major collections:

- H** **hallmark gene sets** are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes.
- C1** **positional gene sets** for each human chromosome and cytogenetic band.
- C2** **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts.
- C3** **motif gene sets** based on conserved cis-regulatory motifs from a comparative analysis of the human, mouse, rat, and dog genomes.
- C4** **computational gene sets** defined by mining large collections of cancer-oriented microarray data.
- C5** **GO gene sets** consist of genes annotated by the same GO terms.
- C6** **oncogenic gene sets** defined directly from microarray gene expression data from cancer gene perturbations.
- C7** **immunologic gene sets** defined directly from microarray gene expression data from immunologic studies.

Aravind Subramanian et al. PNAS 2005;102:43:15545-15550

PNAS

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 😊

