

RNAseq

Michael Schatz

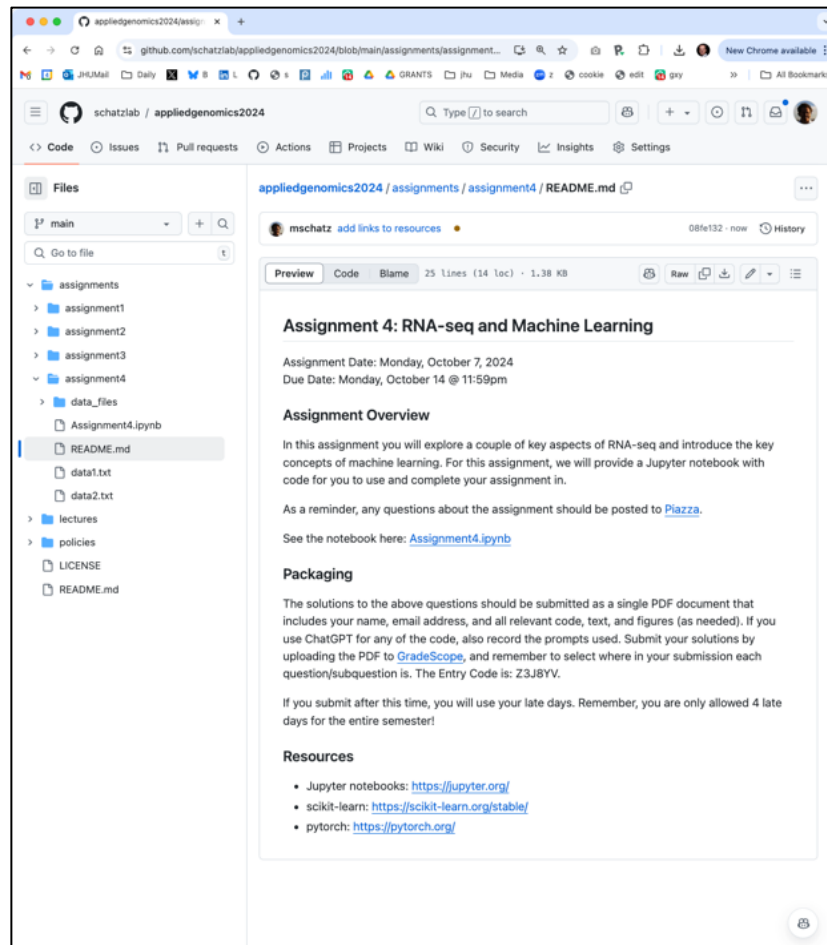
October 7, 2024

Lecture 12. Applied Comparative Genomics



Assignment 4

Due: Monday Oct 14, 2024 by 11:59pm



The screenshot shows the GitHub repository page for 'appliedgenomics2024' by 'mschatz'. The file 'Assignment4/README.md' is selected. The README content includes the assignment title, dates, overview, packaging instructions, and resources.

Assignment 4: RNA-seq and Machine Learning

Assignment Date: Monday, October 7, 2024
Due Date: Monday, October 14 @ 11:59pm

Assignment Overview

In this assignment you will explore a couple of key aspects of RNA-seq and introduce the key concepts of machine learning. For this assignment, we will provide a Jupyter notebook with code for you to use and complete your assignment in.

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

See the notebook here: [Assignment4.ipynb](#)

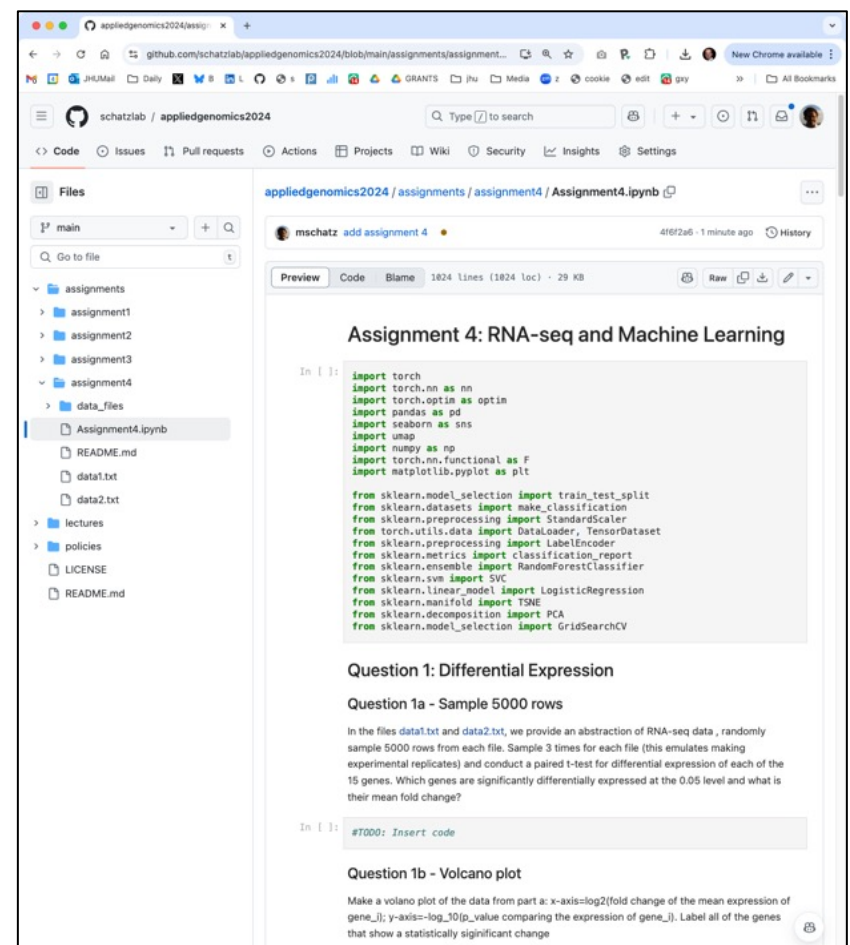
Packaging

The solutions to the above questions should be submitted as a single PDF document that includes your name, email address, and all relevant code, text, and figures (as needed). If you use ChatGPT for any of the code, also record the prompts used. Submit your solutions by uploading the PDF to [GradeScope](#), and remember to select where in your submission each question/subquestion is. The Entry Code is: Z3J8YV.

If you submit after this time, you will use your late days. Remember, you are only allowed 4 late days for the entire semester!

Resources

- Jupyter notebooks: <https://jupyter.org/>
- scikit-learn: <https://scikit-learn.org/stable/>
- pytorch: <https://pytorch.org/>



The screenshot shows the GitHub repository page for 'appliedgenomics2024' by 'mschatz'. The file 'Assignment4/Assignment4.ipynb' is selected. The Jupyter notebook content includes the assignment title, a list of imports, and the start of Question 1: Differential Expression.

Assignment 4: RNA-seq and Machine Learning

```
In [ ]: import torch
import torch.nn as nn
import torch.optim as optim
import pandas as pd
import seaborn as sns
import umap
import numpy as np
import torch.nn.functional as F
import matplotlib.pyplot as plt

from sklearn.model_selection import train_test_split
from sklearn.datasets import make_classification
from sklearn.preprocessing import StandardScaler
from torch.utils.data import DataLoader, TensorDataset
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import classification_report
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression
from sklearn.manifold import TSNE
from sklearn.decomposition import PCA
from sklearn.model_selection import GridSearchCV
```

Question 1: Differential Expression

Question 1a - Sample 5000 rows

In the files `data1.txt` and `data2.txt`, we provide an abstraction of RNA-seq data, randomly sample 5000 rows from each file. Sample 3 times for each file (this emulates making experimental replicates) and conduct a paired t-test for differential expression of each of the 15 genes. Which genes are significantly differentially expressed at the 0.05 level and what is their mean fold change?

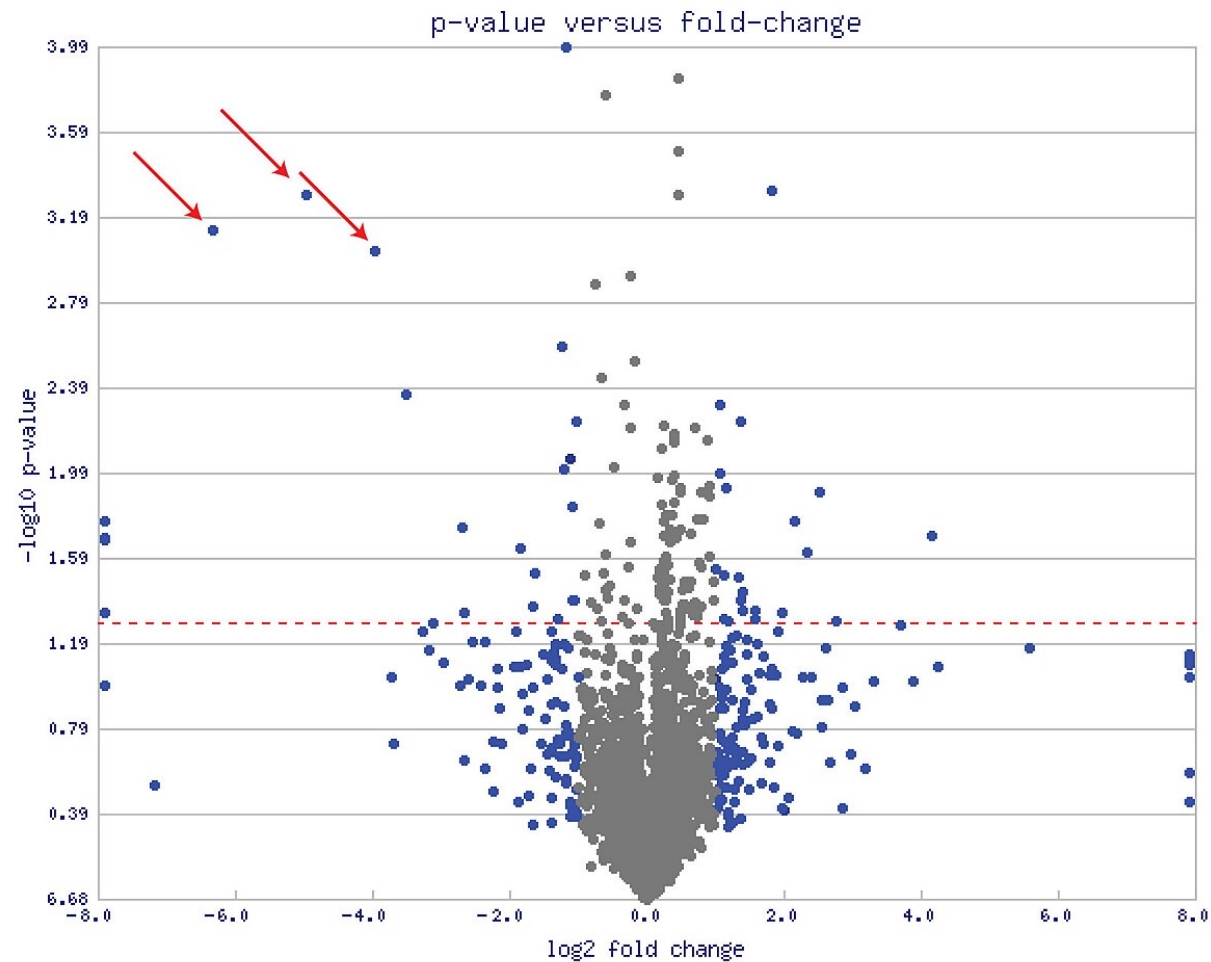
Question 1b - Volcano plot

Make a volcano plot of the data from part a: x-axis= $\log_2(\text{fold change of the mean expression of gene}_i)$; y-axis= $-\log_{10}(p\text{-value comparing the expression of gene}_i)$. Label all of the genes that show a statistically significant change

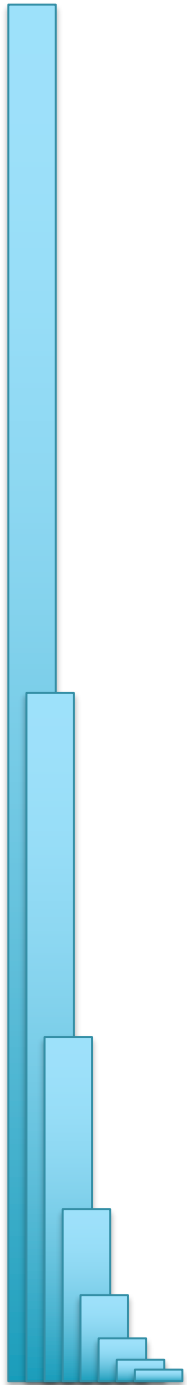
<https://schatz-lab.org/appliedgenomics2024/assignments/assignment4/>

Check Piazza for questions!

Volcano Plot



https://en.wikipedia.org/wiki/Volcano_plot_%28statistics%29

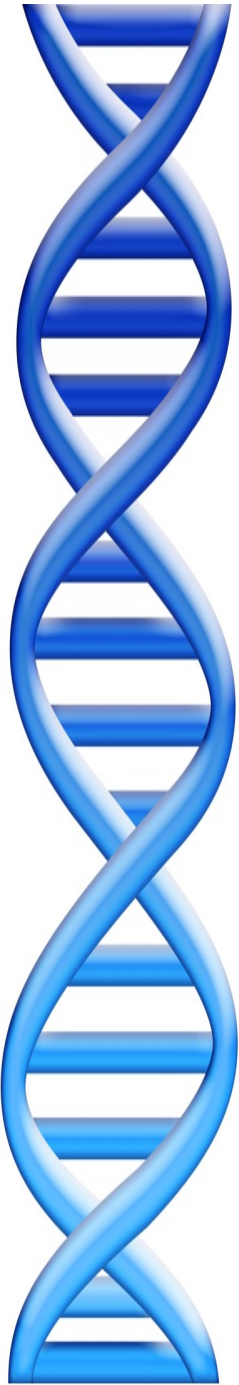


Annotation

Goal: Genome Annotations

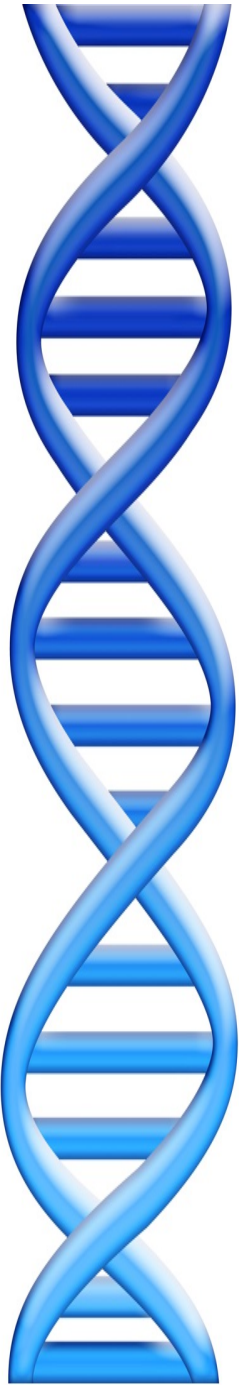
aatgcatgctggctatgctaagcatgctggctatgctaagctgggatccgatgacaatgcatgctggctatgctaag
gcatgctggctatgcaagctgggatccgatgactatgctaagctgggatccgatgacaatgcatgctggctatgct
aatgaatgggtcttgggatttaccttgggaatgctaagctgggatccgatgacaatgcatgctggctatgctaag
tggtcttgggatttaccttgggaatgctaagcatgctggctatgctaagctgggatccgatgacaatgcatgctg
gctatgctaagcatgctggctatgcaagctgggatccgatgactatgctaagctgctggctatgctaagcatgctg
gctatgctaagctgggatccgatgacaatgcatgctggctatgctaagcatgctggctatgcaagctgggatcc
gctggctatgctaagcatgctgggtcttgggatttaccttgggaatgctaagctgggatccgatgacaatgcatgctg
atgctaagcatgctgggtcttgggatttaccttgggaatgctaagcatgctggctatgctaagcatgctg
gctatgctaagctgggatccgatgacaatgcatgctggctatgctaagcatgctggctatgcaagctgggatccg
atgactatgctaagctgctggctatgctaagcatgctggctatgctaagcatgctggctatgctaagctgggaat
gcatgctggctatgctaagctgggatccgatgacaatgcatgctggctatgctaagcatgctggctatgcaagctg
ggatccgatgactatgctaagctgctggctatgctaagcatgctggctatgctaagctgctggctatgctaagcatg
gtcttgggatttaccttgggaatgctaagctgggatccgatgacaatgcatgctggctatgctaagcatgctgg
gatttaccttgggaatgctaagcatgctggctatgctaagctgggaatgcatgctggctatgctaagctgggatc
cgatgacaatgcatgctggctatgctaagcatgctggctatgcaagctgggatccgatgactatgctaagctgctg
gctatgctaagcatgctggctatgctaagctcatgctg

Gene!



Outline

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



Outline

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

Very Similar Sequences

Query: HBA_HUMAN Hemoglobin alpha subunit

Sbjct: HBB_HUMAN Hemoglobin beta subunit

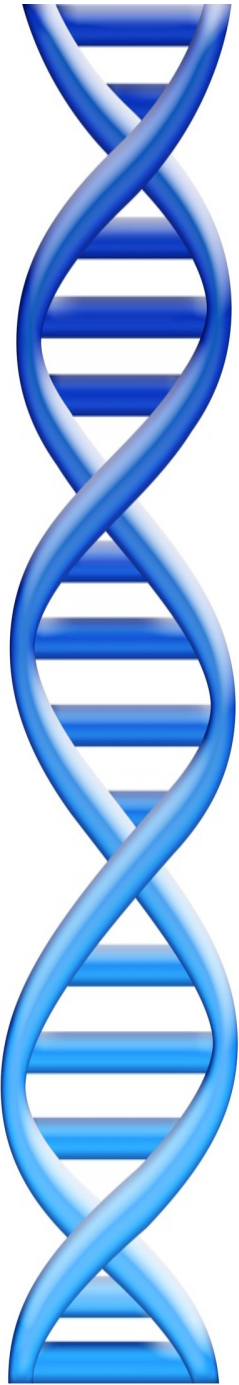
Score = 114 bits (285), Expect = 1e-26

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

```
Query    2    LSPADKTNVKAANGKVGAGHAGEYGAELERMFLSFPTTKTYFPHF-----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    KAHGKKVLGAFSDGLAHLNLRGTFATLSELDKLVDPENFRLLGNVLVCVLAHHFGK 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



Bacterial Gene Finding and Glimmer

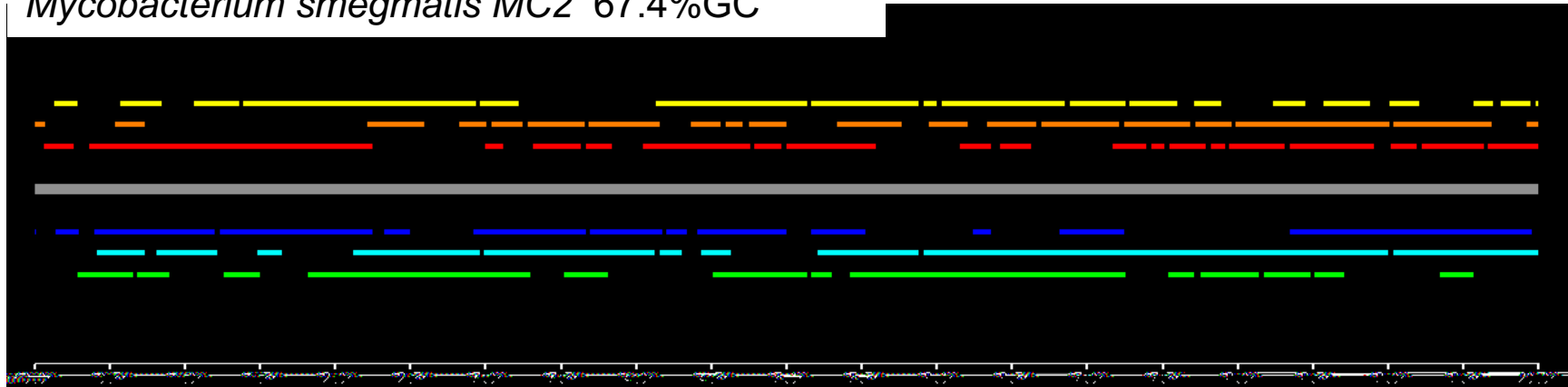
(also Archaeal and viral gene finding)

Arthur L. Delcher and Steven Salzberg

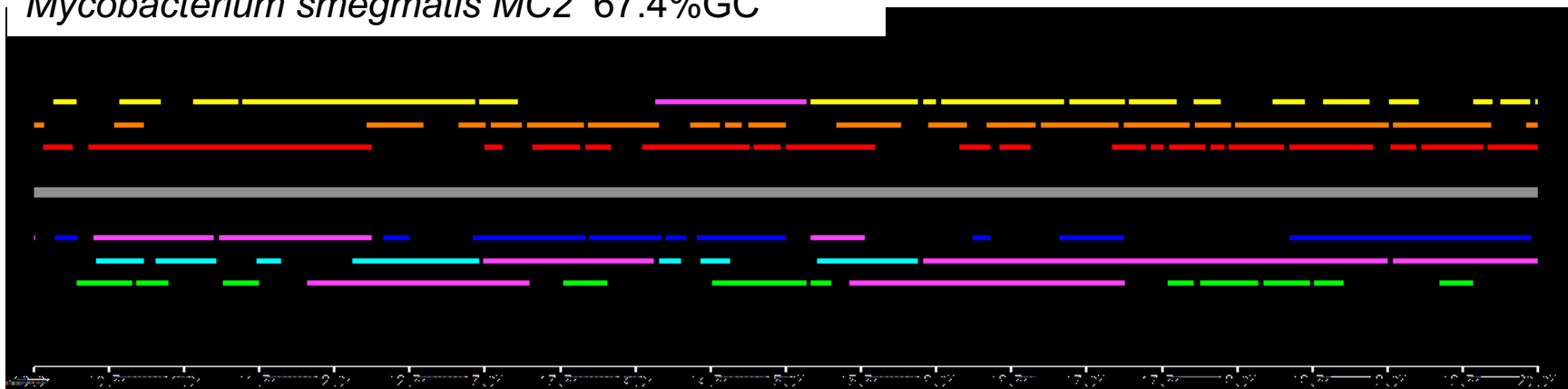
Center for Bioinformatics and Computational Biology

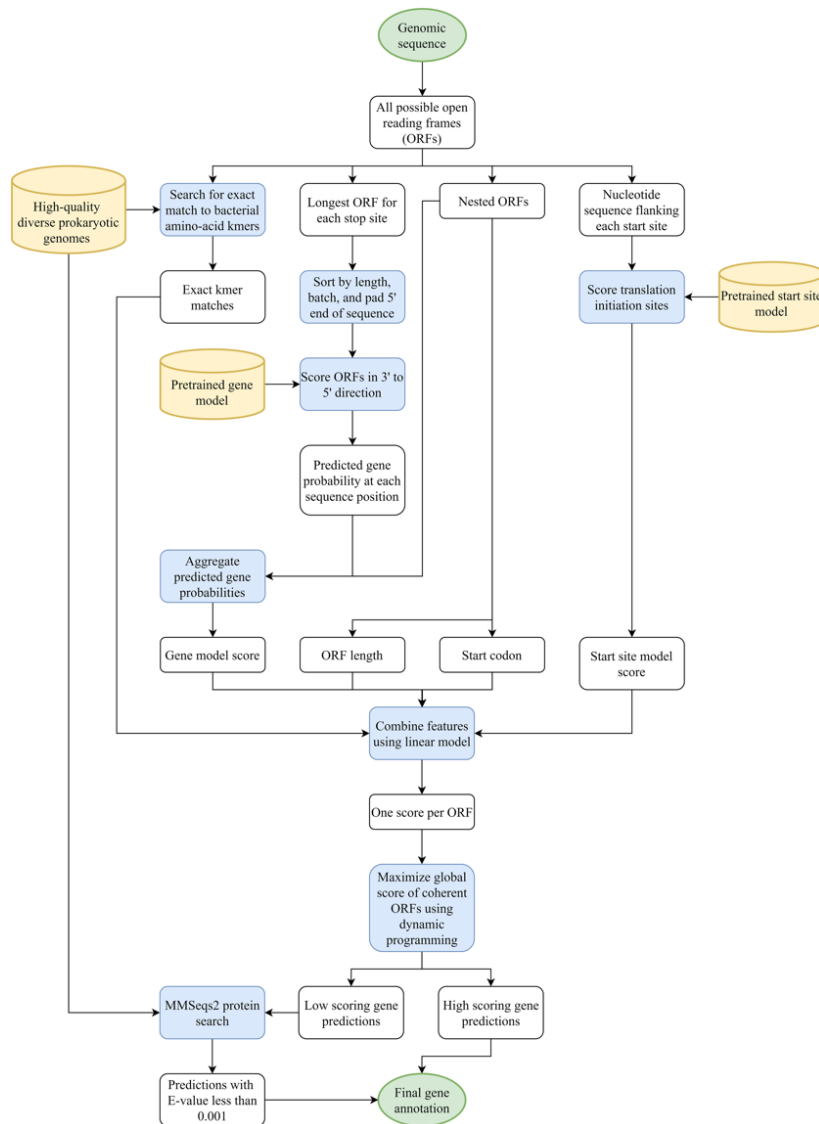
Johns Hopkins University

Mycobacterium smegmatis MC2 67.4%GC

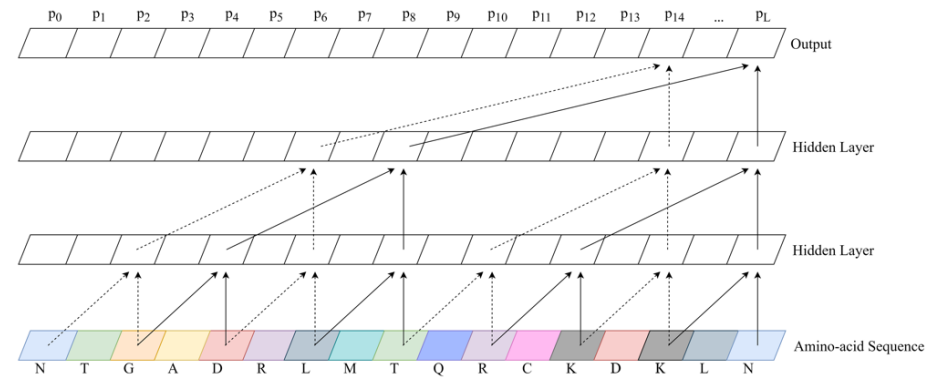


Mycobacterium smegmatis MC2 67.4%GC





Temporal Convolutional Network



Balrog: A universal protein model for prokaryotic gene prediction

Sommer, MJ, Salzberg, SL (2021) PLOS Comp. Bio. doi: 10.1371/journal.pcbi.1008727

Probabilistic Methods

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

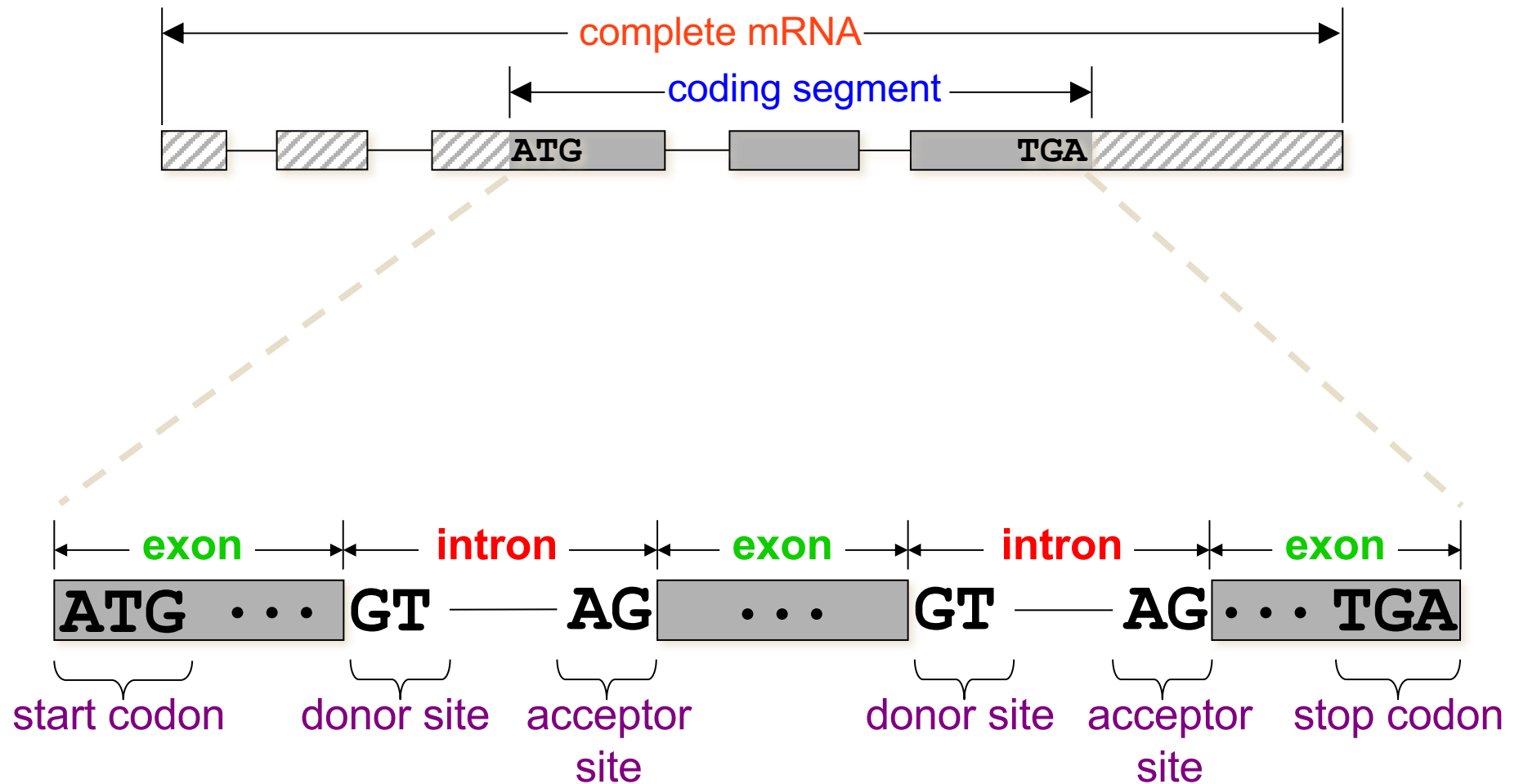


Overview of Eukaryotic Gene Prediction

CBB 231 / COMPSCI 261

W.H. Majoros

Eukaryotic Gene Syntax



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

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

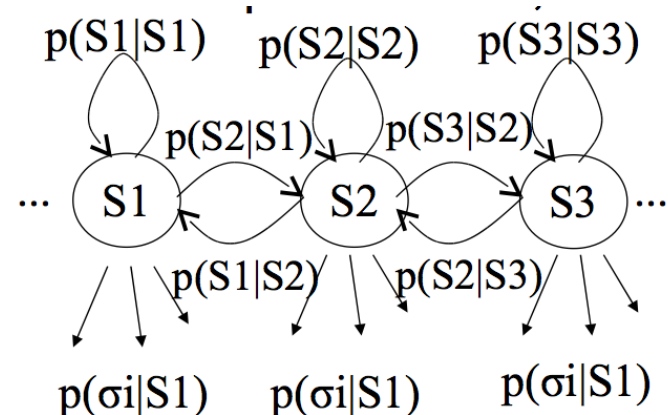
- $E=\{H,T\}$ for coin tossing
 - $E=\{1,2,3,4,5,6\}$ for dice tossing
 - $E=\{A,C,G,T\}$ for DNA

- State-specific emission probabilities

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

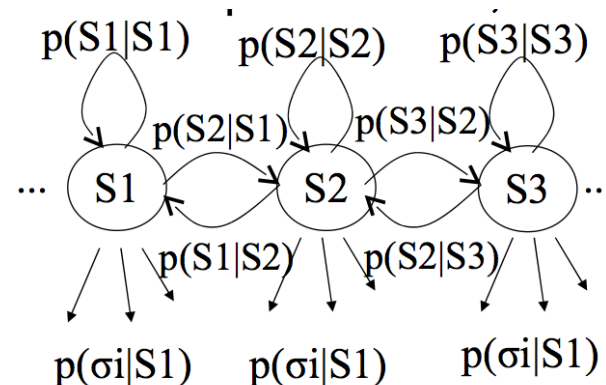
- A probability of taking a transition

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



Why Hidden?

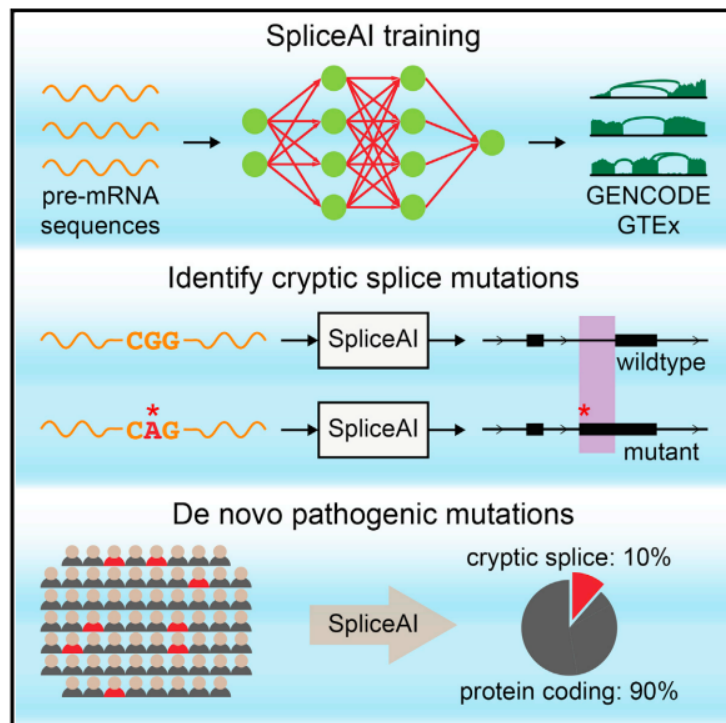
- Similar to Markov models used for prokaryotic gene finding, but system may transition between multiple models called states (gene/non-gene, intergenic/exon/intron)
- Observers can see the emitted symbols of an HMM (i.e., nucleotides) but have no ability to know which state the HMM is currently in.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



AAAGCATGCATTTAACGTGAGCACAAATAGATTACA

Predicting Splicing from Primary Sequence with Deep Learning

Graphical Abstract



Authors

Kishore Jaganathan,
Sofia Kyriazopoulou Panagiotopoulou,
Jeremy F. McRae, ..., Serafim Batzoglou,
Stephan J. Sanders, Kyle Kai-How Farh

Correspondence

kfarh@illumina.com

In Brief

A deep neural network precisely models mRNA splicing from a genomic sequence and accurately predicts noncoding cryptic splice mutations in patients with rare genetic diseases.

Highlights

- SpliceAI, a 32-layer deep neural network, predicts splicing from a pre-mRNA sequence
- 75% of predicted cryptic splice variants validate on RNA-seq
- Cryptic splicing may yield ~10% of pathogenic variants in neurodevelopmental disorders
- Cryptic splice variants frequently give rise to alternative splicing

Genome analysis

Helixer: cross-species gene annotation of large eukaryotic genomes using deep learning

Felix Stiehler¹, Marvin Steinborn¹, Stephan Scholz, Daniela Dey²,
Andreas P. M. Weber¹ and Alisandra K. Denton^{1,*} 

¹Institute of Plant Biochemistry, Faculty of Mathematics and Natural Sciences, Heinrich-Heine-University, Dusseldorf 40225, Germany and ²Institute of Human Genetics, Medical Faculty, RWTH Aachen University, Aachen 52062, Germany

*To whom correspondence should be addressed.

Associate Editor: Inanc Birol

Received on July 31, 2020; revised on November 6, 2020; editorial decision on December 4, 2020; accepted on December 7, 2020

Abstract

Motivation: Current state-of-the-art tools for the *de novo* annotation of genes in eukaryotic genomes have to be specifically fitted for each species and still often produce annotations that can be improved much further. The fundamental algorithmic architecture for these tools has remained largely unchanged for about two decades, limiting learning capabilities. Here, we set out to improve the cross-species annotation of genes from DNA sequence alone with the help of deep learning. The goal is to eliminate the dependency on a closely related gene model while also improving the predictive quality in general with a fundamentally new architecture.

Results: We present Helixer, a framework for the development and usage of a cross-species deep learning model that improves significantly on performance and generalizability when compared to more traditional methods. We evaluate our approach by building a single vertebrate model for the base-wise annotation of 186 animal genomes and a separate land plant model for 51 plant genomes. Our predictions are shown to be much less sensitive to the length of the genome than those of a current state-of-the-art tool. We also present two novel post-processing techniques that each worked to further strengthen our annotations and show in-depth results of an RNA-Seq based comparison of our predictions. Our method does not yet produce comprehensive gene models but rather outputs base pair wise probabilities.

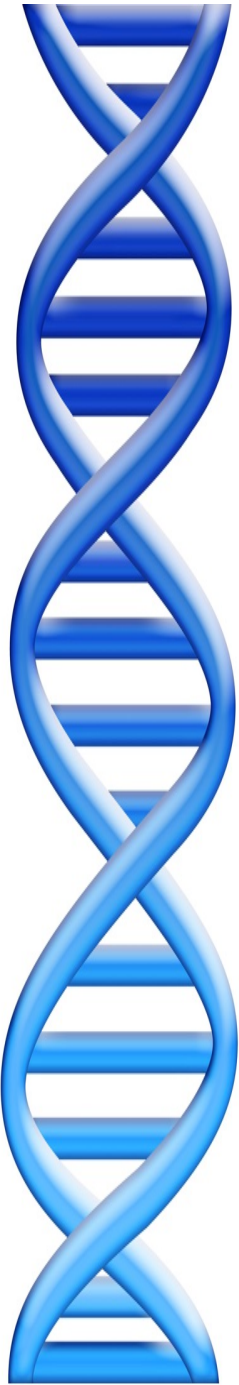
Availability and implementation: The source code of this work is available at <https://github.com/weberlab-hhu/Helixer> under the GNU General Public License v3.0. The trained models are available at <https://doi.org/10.5281/zenodo.3974409>

Contact: alisandra.denton@hhu.de

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

Gene Finding Overview

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



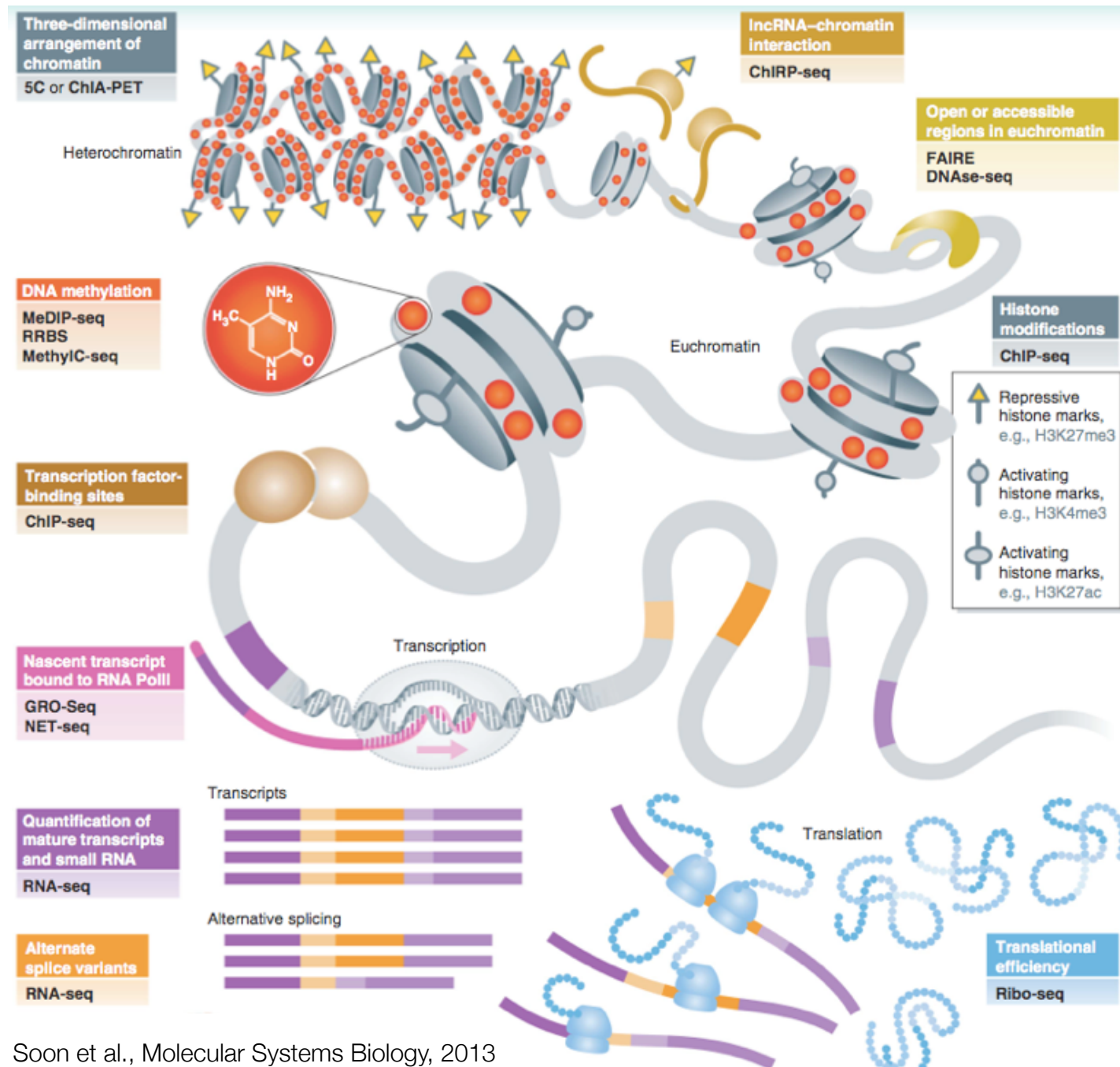
Outline

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

Sequencing Assays

The *Seq List (in chronological order)

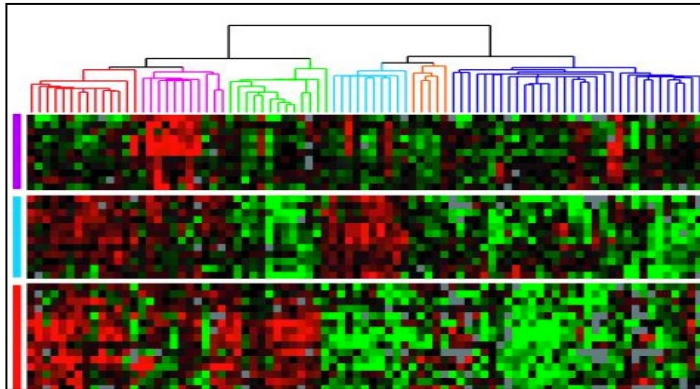
1. Gregory E. Crawford et al., "Genome-wide Mapping of DNase Hypersensitive Sites Using Massively Parallel Signature Sequencing (MPSS)," *Genome Research* 16, no. 1 (January 1, 2006): 123–131, doi:10.1101/gr.4074106.
2. David S. Johnson et al., "Genome-Wide Mapping of in Vivo Protein-DNA Interactions,"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



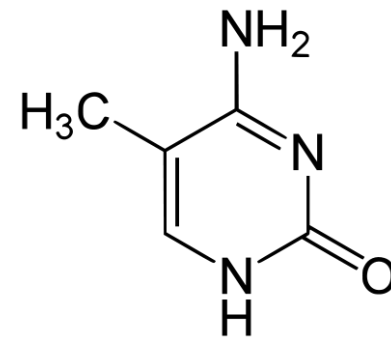
Soon et al., Molecular Systems Biology, 2013

*-seq in 4 short vignettes

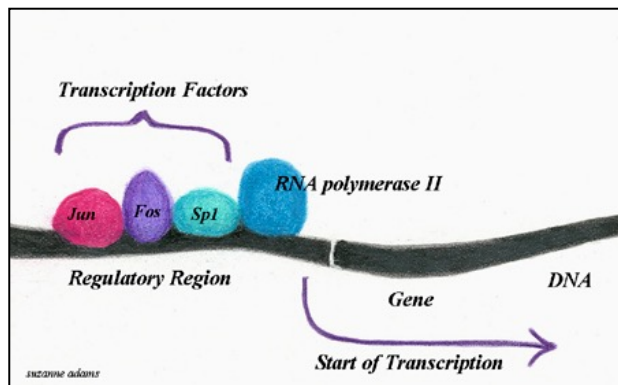
RNA-seq



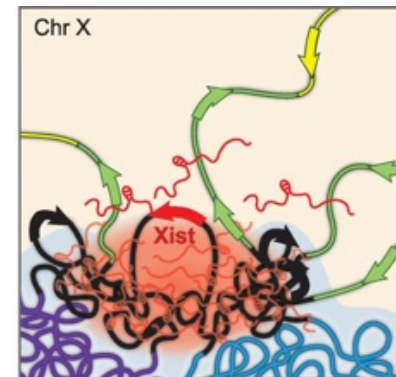
Methyl-seq



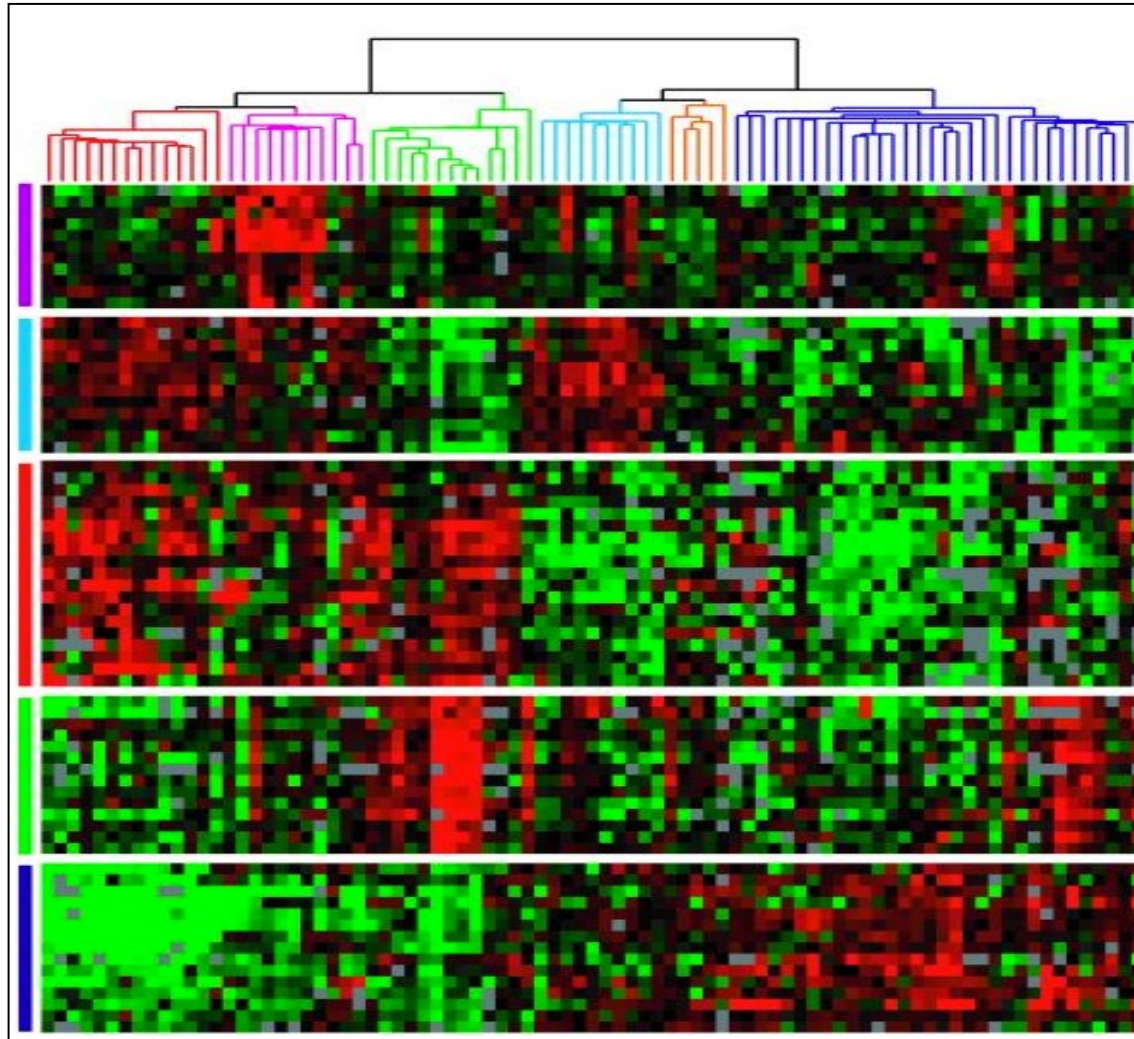
ChIP-seq



Hi-C

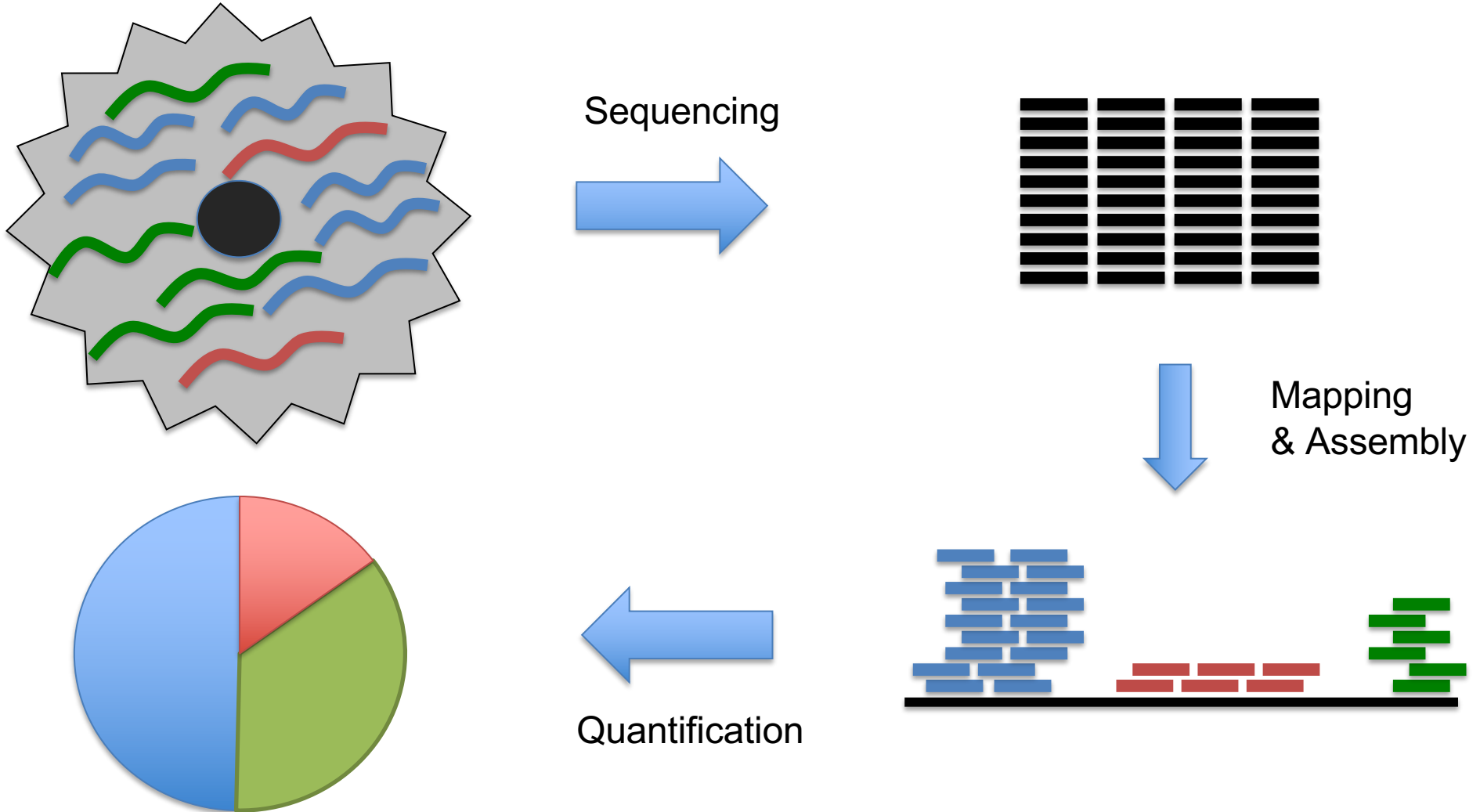


RNA-seq

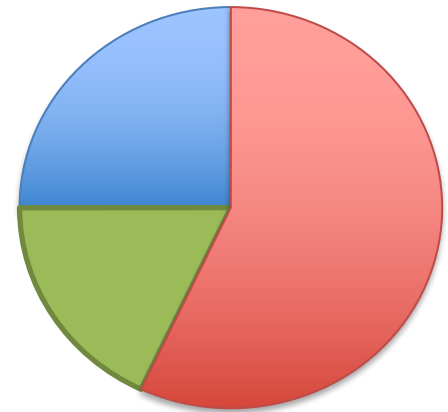
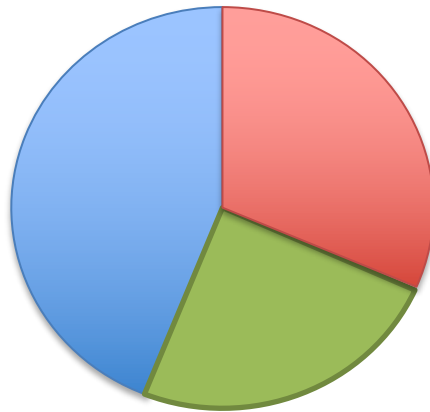
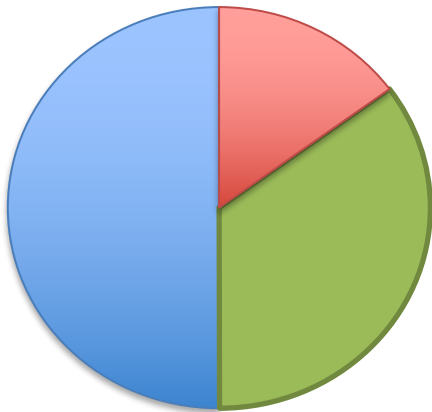
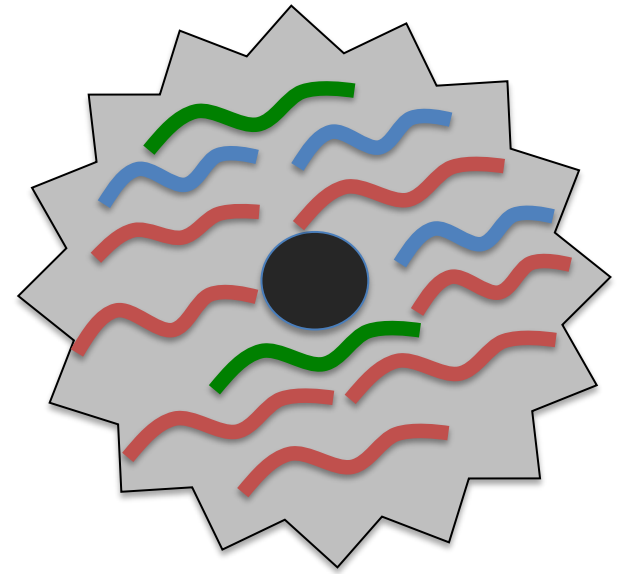
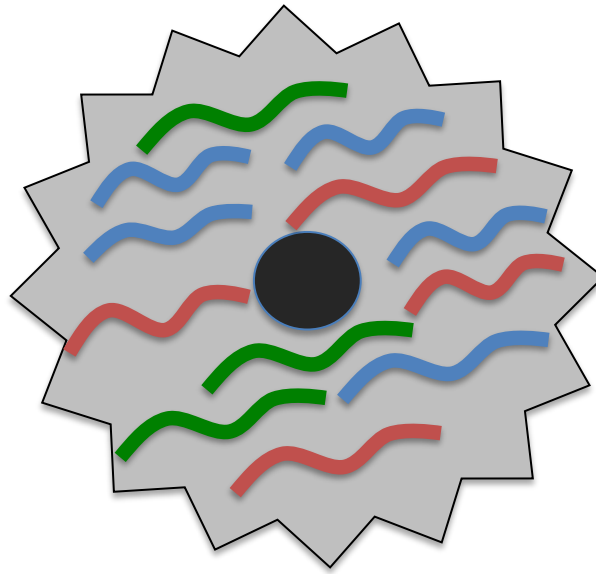
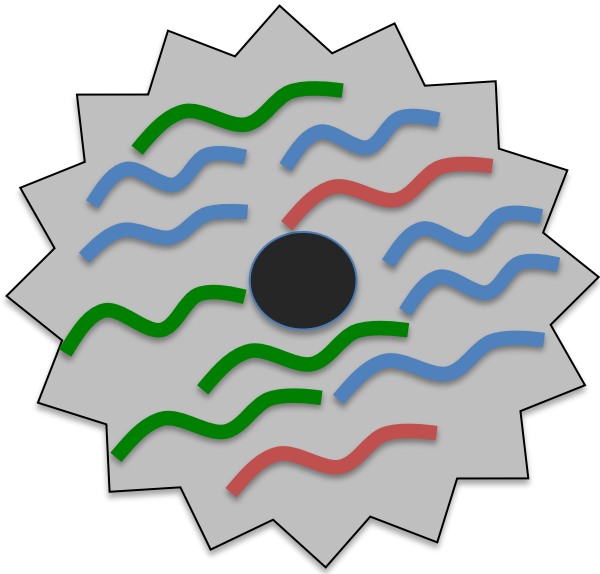


Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications.
Sørli et al (2001) *PNAS*. 98(19):10869-74.

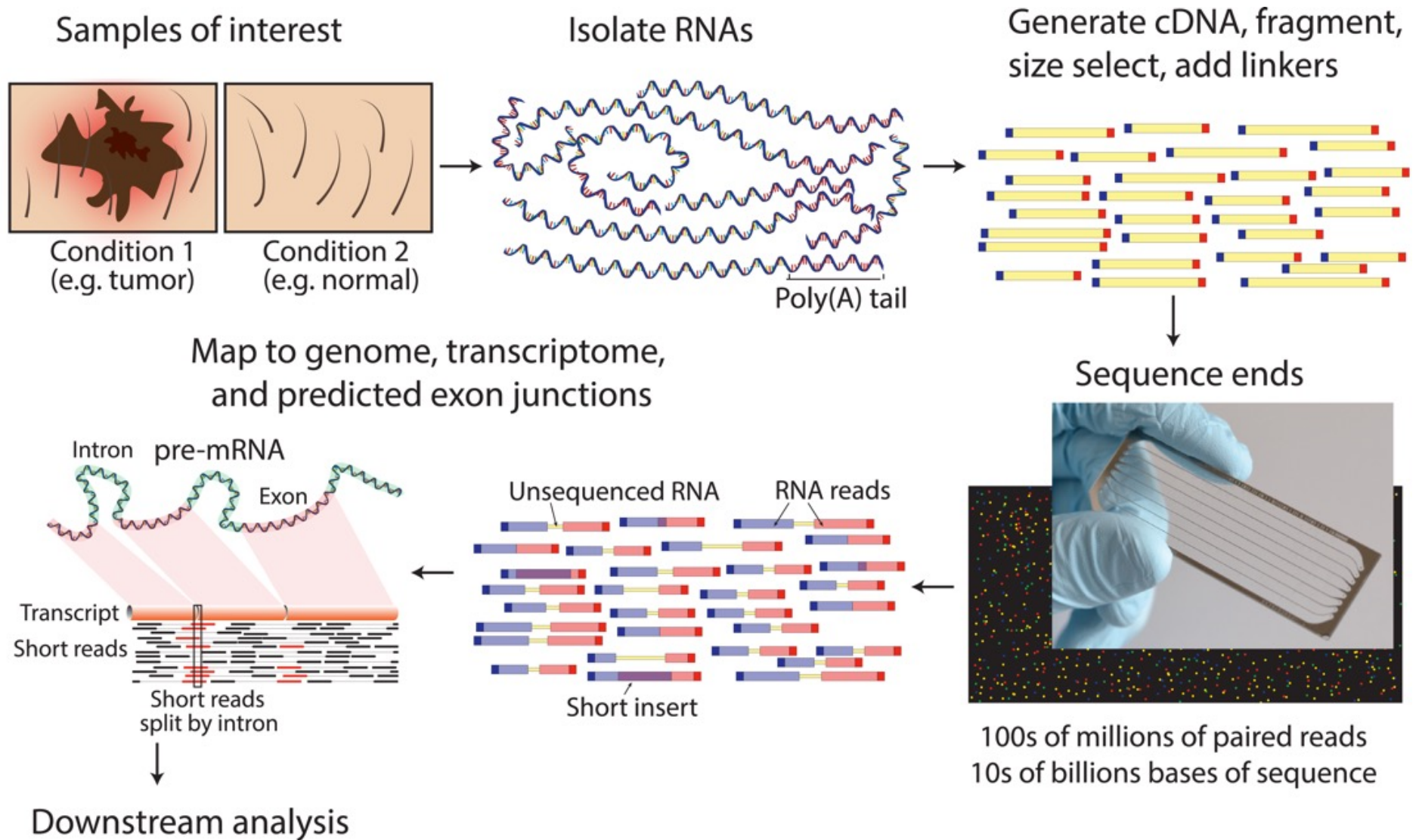
RNA-seq Overview



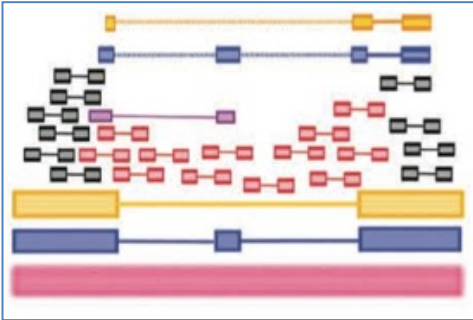
RNA-seq Overview



RNA-seq Overview



RNA-seq Challenges



Challenge I: Eukaryotic genes are spliced

RNA-Seq Approaches

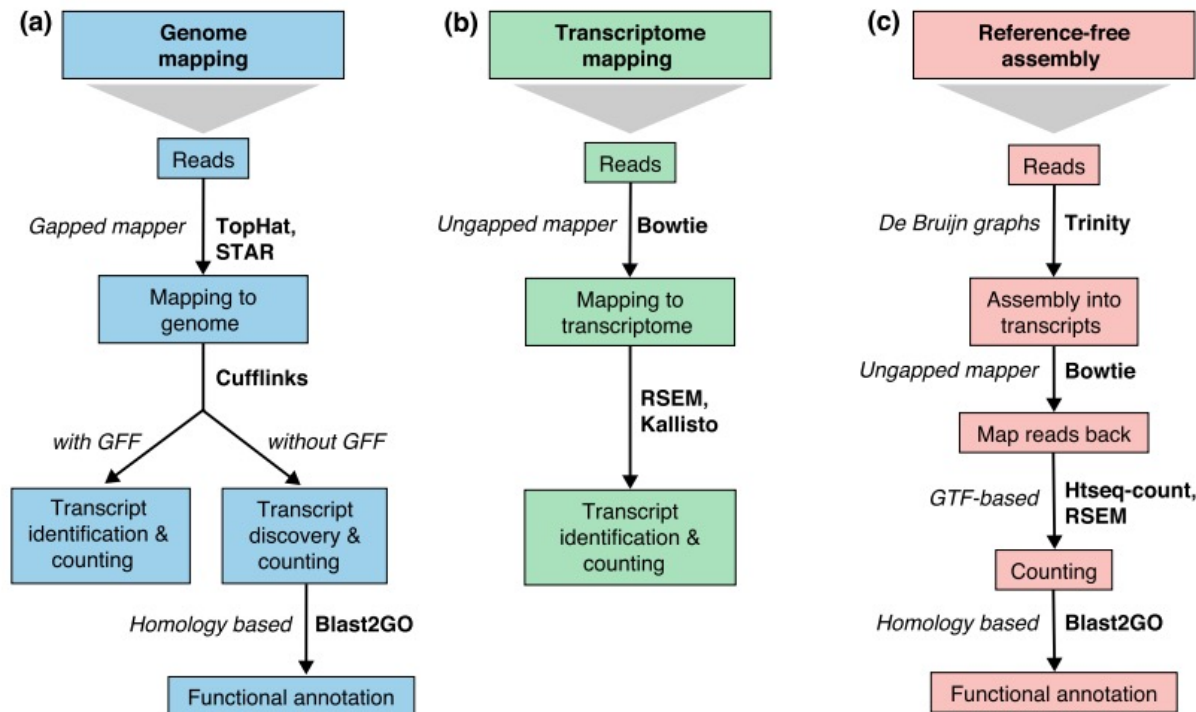


Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome with a gapped mapper. Next (novel) transcript discovery and quantification can proceed with or without an annotation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

A survey of best practices for RNA-seq data analysis

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

RNA-Seq Approaches

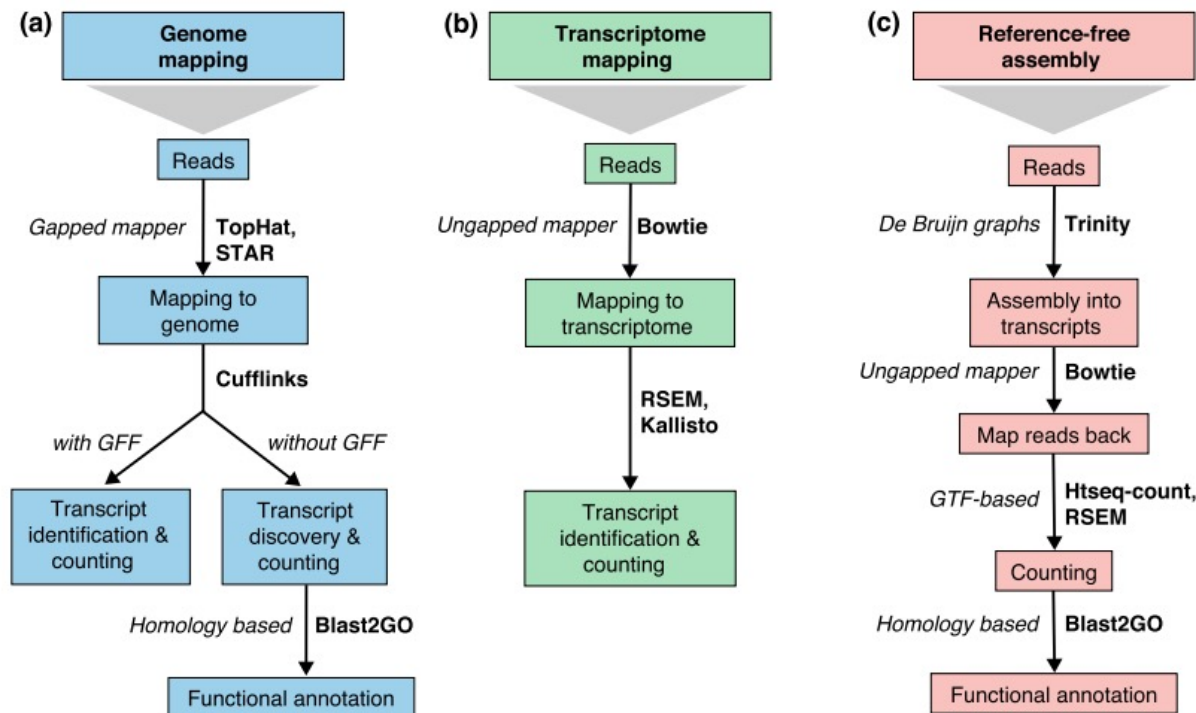


Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads can be mapped to the reference genome. (b) Transcript discovery and quantification can proceed with or without an annotated transcriptome.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

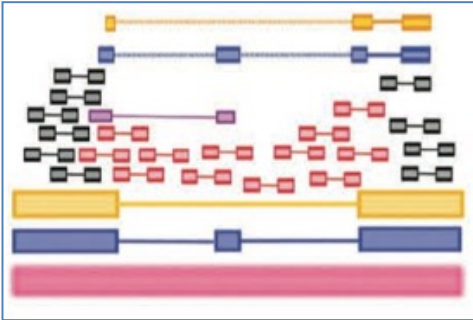
Which approach should we use?

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

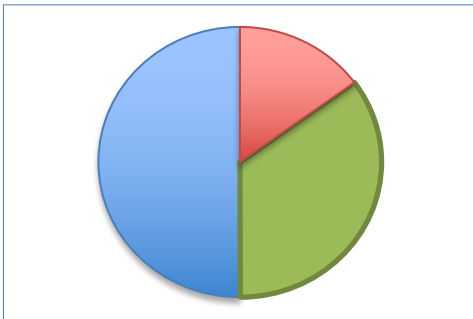


Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

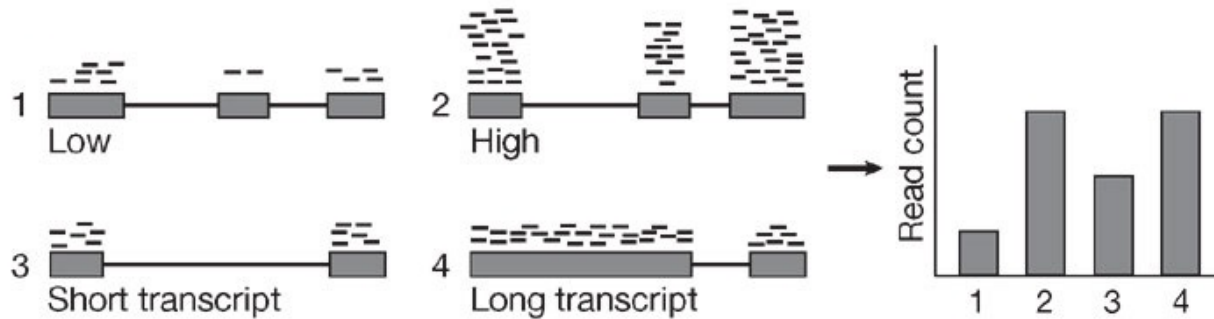
TopHat: discovering spliced junctions with RNA-Seq.

Trapnell et al (2009) *Bioinformatics*. 25:0 ||05-||||



Challenge 2: Read Count != Transcript abundance

RPKM, FPKM, TPM

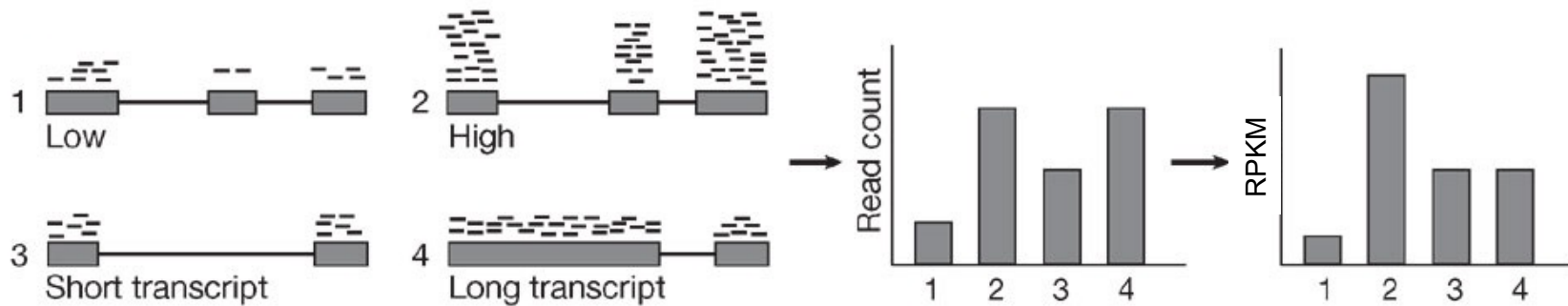


Counting Reads that align to a gene DOESN'T work!

- 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!

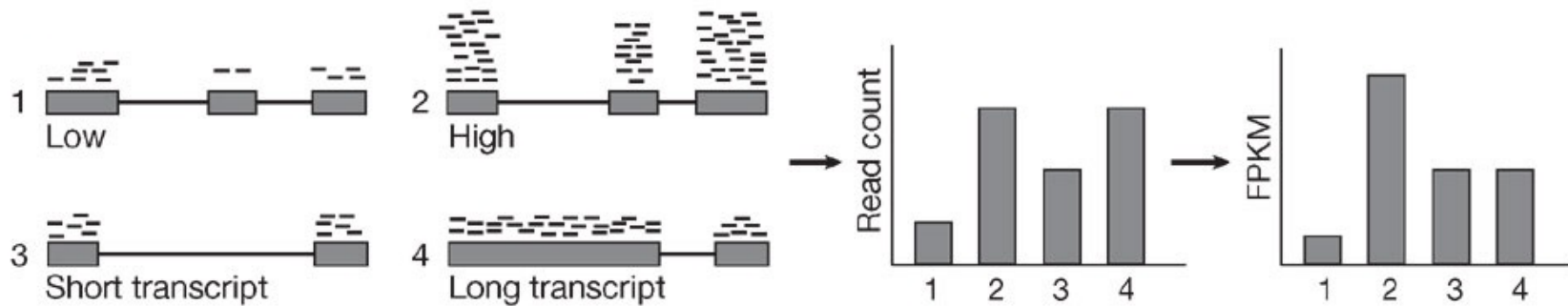
- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

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)

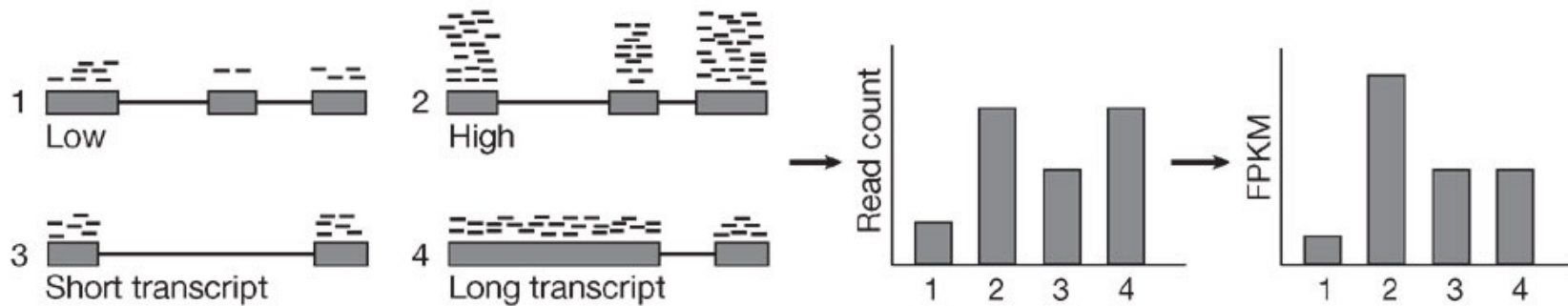
=> Wait a second, reads in a pair are not independent!

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 not 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!

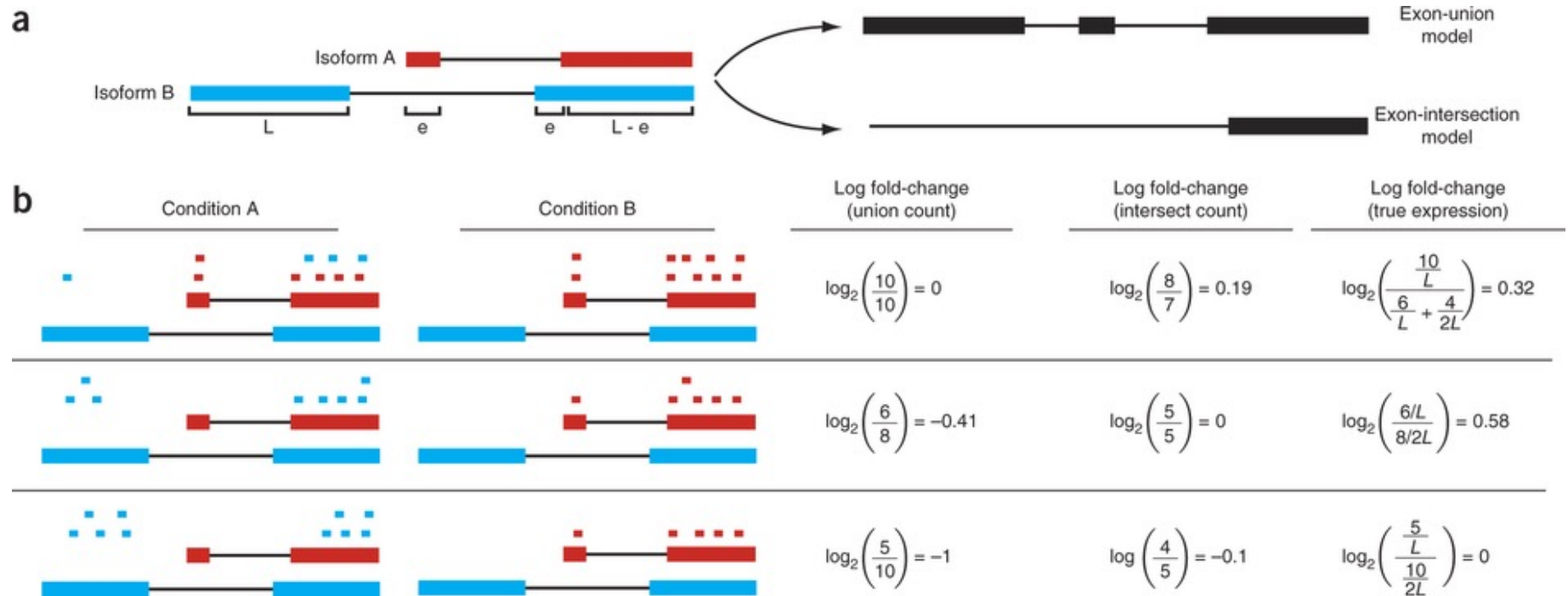
3. TPM: Transcripts Per Million (Li et al, 2011)

=> 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?



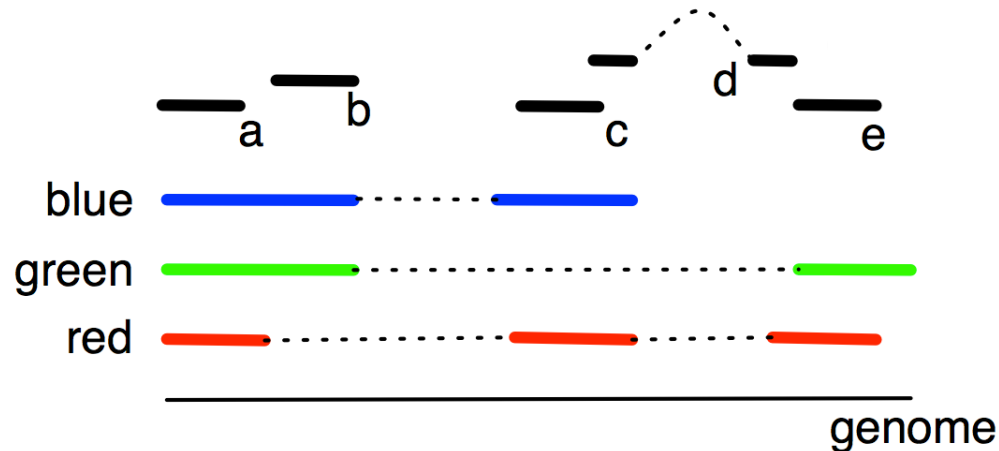
Key point : The length of the actual molecule from which the fragments derive is crucially important to obtaining accurate abundance estimates.

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]