

Estimating Gene Expression

Lecture 27
Nov 4, 2016

GO VOTE!

If you're a U.S. citizen and 18 or older, you have the right to vote. Registering to vote is simpler and faster than you may think!

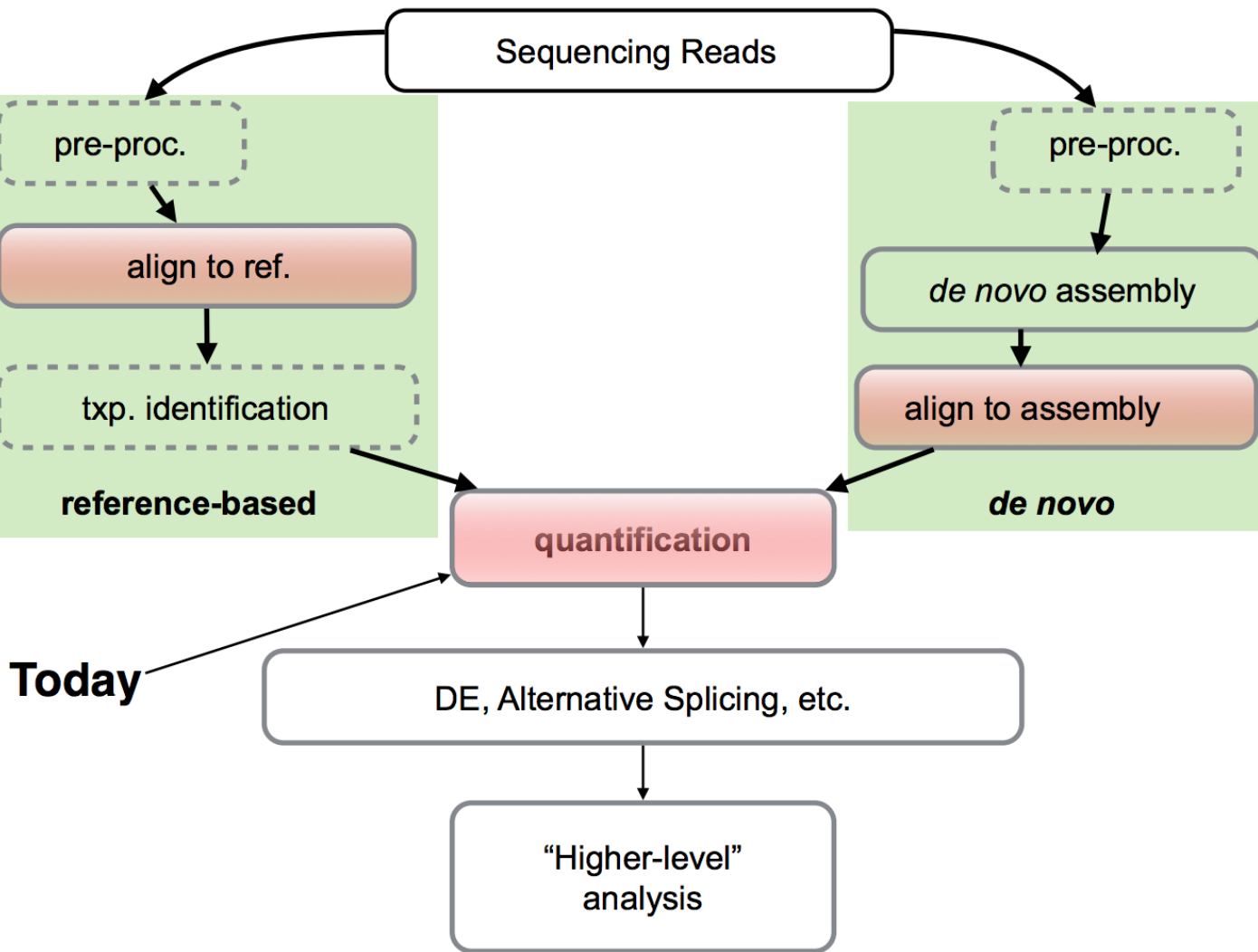
Important Dates

General Election: Tuesday, November 8, 2016

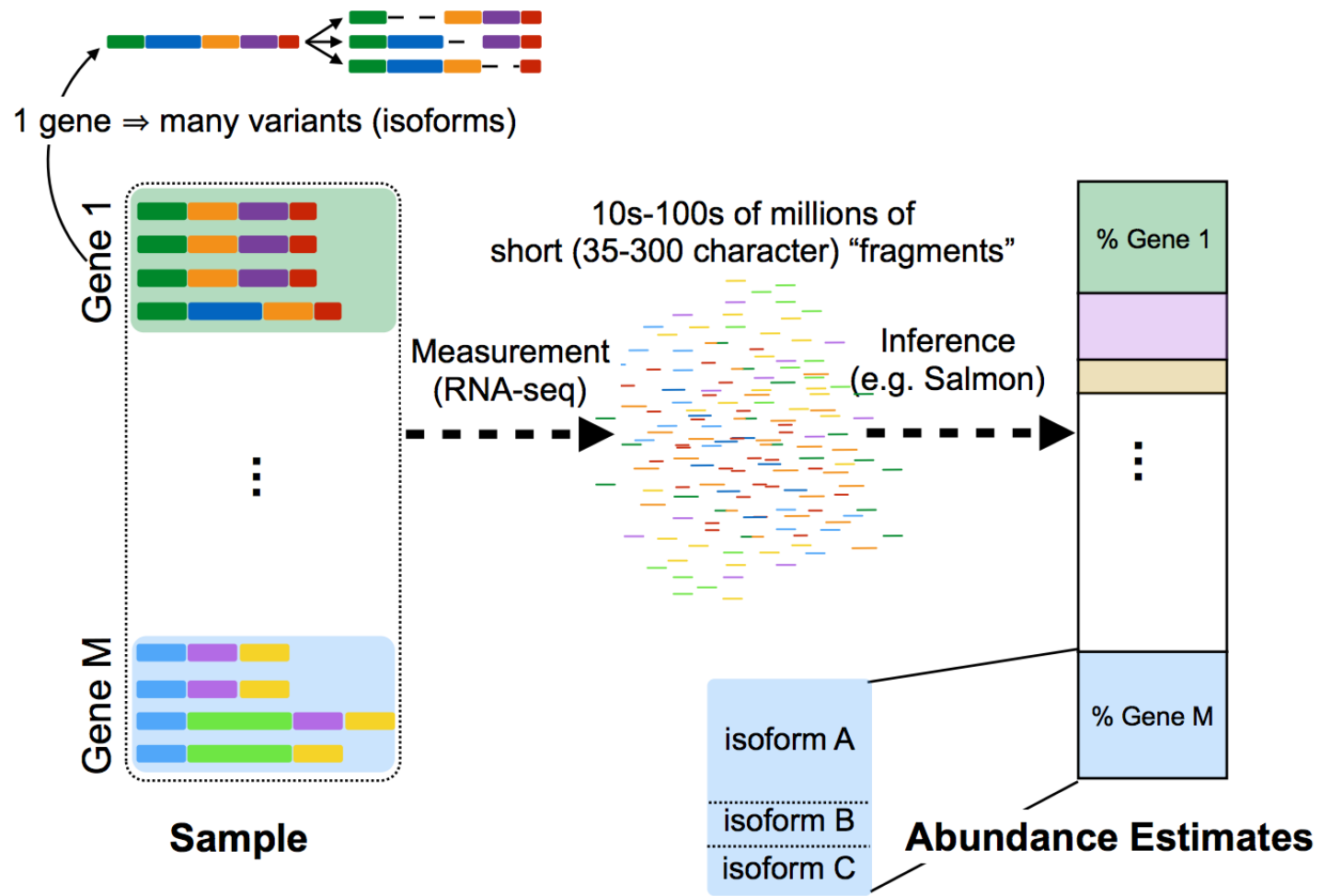
For Durham students, shuttle service will be running from the Holloway Commons bus stop on Main St. to the polling location at Oyster River High School from 7 a.m. to 7 p.m. on Election Day. Please contact the Dean of Students office with any questions.

Voting Info

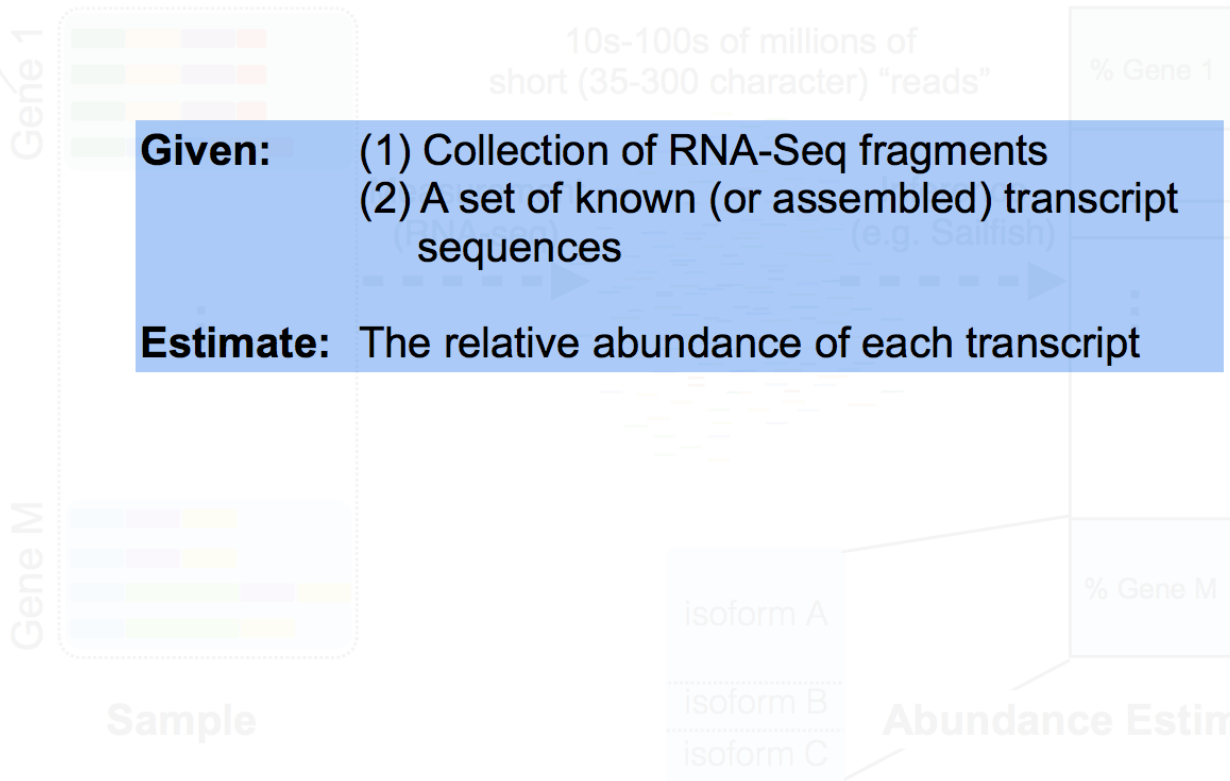
Voting in New Hampshire is quick and easy, especially when you bring appropriate identification. A valid state driver's license or non-driver ID proves your identity, your age, and your domicile (if it shows the address you are claiming as your voting domicile). A lease or piece of mail will also show your domicile. A birth certificate (or copy) or passport will prove your citizenship. If you lack proper identification, you may fill out an affidavit (a legal document attesting to your identity and domicile). You may register in person at the town or city clerk's office up to 10 days prior to an election, or on Election Day at the polling place.



Transcript Quantification: An Overview

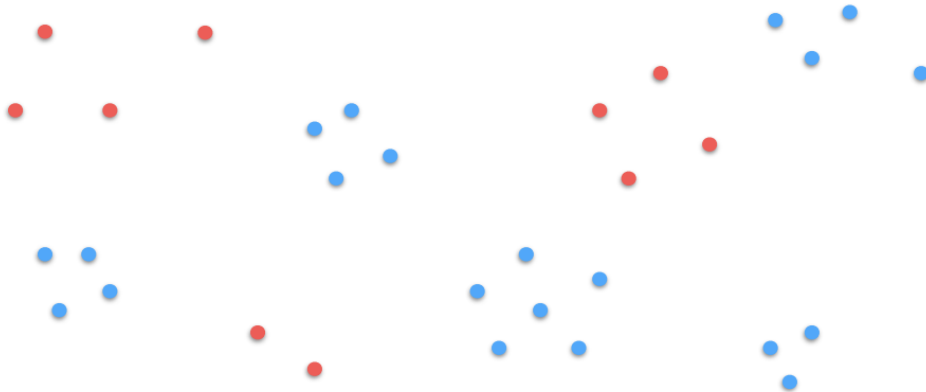


1 gene \Rightarrow many variants (isoforms)



First, consider this non-Biological example

Imagine I have two colors of circle, **red** and **blue**. I want to estimate the **fraction of circles** that are **red** and **blue**. I'll *sample* from them by tossing down darts.



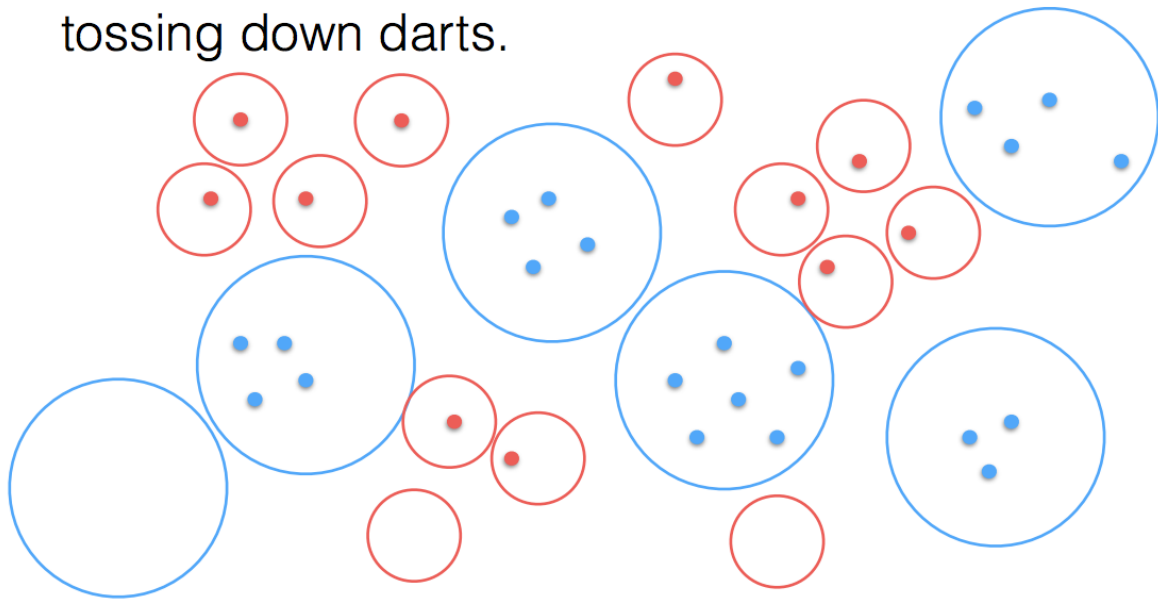
Here, a dot of a color means I hit a circle of that color.

What type of circle is more prevalent?

What is the fraction of red / blue circles?

First, consider this non-Biological example

Imagine I have two colors of circle, **red** and **blue**. I want to estimate the **fraction of circles** that are **red** and **blue**. I'll *sample* from them by tossing down darts.

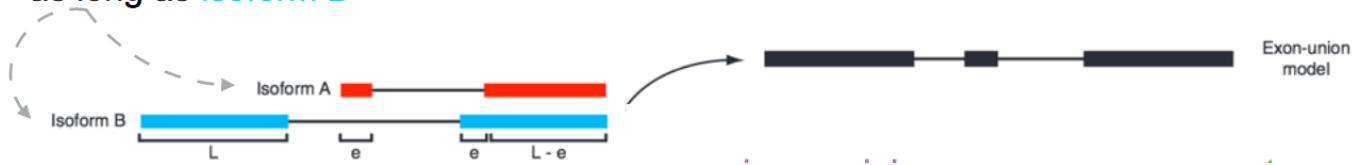


You're missing a **crucial piece of information!**

The areas!

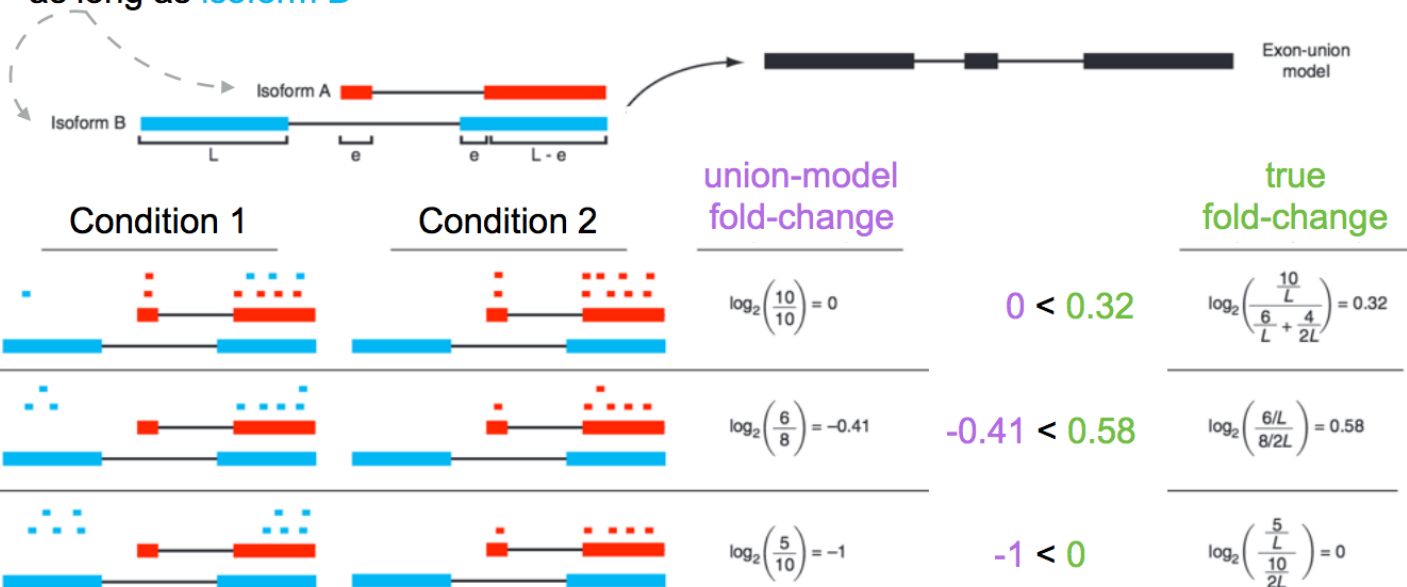
Resolving multi-mapping is fundamental to quantification

Isoform A is half
as long as isoform B



Resolving multi-mapping is fundamental to quantification

Isoform A is half
as long as isoform B



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

How to assess “abundance”

RPKM — Reads per kilobase per million mapped reads

FPKM — Fragments per kilobase per million mapped reads

↖
Don't use these measures, TPM measures the
“same thing”, but in a better way.

TPM — Transcripts per million

↖
Useful for visualization / assessment etc.

(Estimated) Number of Reads

↖
These are what are used (after normalization)
for differential expression. Why can't we use TPM?

GENE EXPRESSION

<http://haroldpimentel.wordpress.com/2014/05/08/what-the-fpkm-a-review-rna-seq-expression-units/>

A probabilistic view of RNA-Seq quantification

nucleotide fractions true read origins assumes independence of fragments

$$\Pr\{\mathcal{F} \mid \boldsymbol{\eta}, \mathbf{Z}, \mathcal{T}\} = \prod_{j=1}^N \Pr\{f_j \mid \boldsymbol{\eta}, \mathbf{Z}, \mathcal{T}\}$$

observed fragments (reads)

We can safely truncate $\Pr\{t_i \mid \boldsymbol{\eta}\}$ to 0 for transcripts where a fragment doesn't map.

$$= \prod_{j=1}^N \sum_{i=1}^M \Pr\{t_i \mid \boldsymbol{\eta}\} \cdot \Pr\{f_j \mid t_i, z_{ji} = 1\}$$

Prob. of selecting t_i given $\boldsymbol{\eta}$ Prob. of generating fragment f_j given t_i

Depends on abundance estimate Independent of abundance estimate

We want to find the values of $\boldsymbol{\eta}$ that **maximize** this probability.
 We can do this (at least locally) using the EM algorithm.

Biases abound in RNA-seq data

Biases in prep & sequencing can have a significant effect on the fragments we see.

Fragment gc-bias¹—

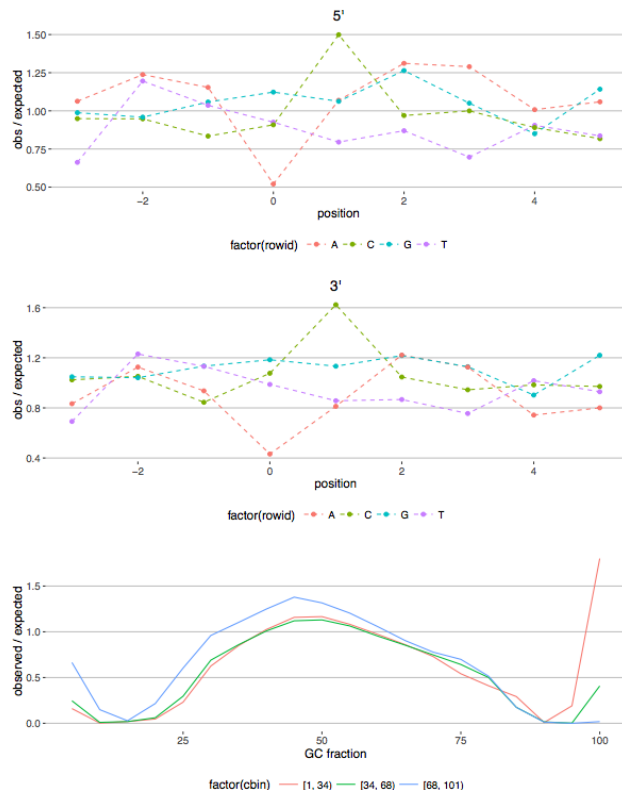
The GC-content of the fragment affects the likelihood of sequencing

Sequence-specific bias²—

sequences surrounding fragment affect the likelihood of sequencing

Positional bias²—

fragments sequenced non-uniformly across the body of a transcript



1: Love, Michael I., John B. Hogenesch, and Rafael A. Irizarry. "Modeling of RNA-seq fragment sequence bias reduces systematic errors in transcript abundance estimation." *bioRxiv* (2015): 025767.

2: Roberts, Adam, et al. "Improving RNA-Seq expression estimates by correcting for fragment bias." *Genome biology* 12.3 (2011): 1.