

COM SCI C121 Week 3

Aidan Jan

April 16, 2024

RNA-seq

- "-seq": probing the molecular biology of the cell

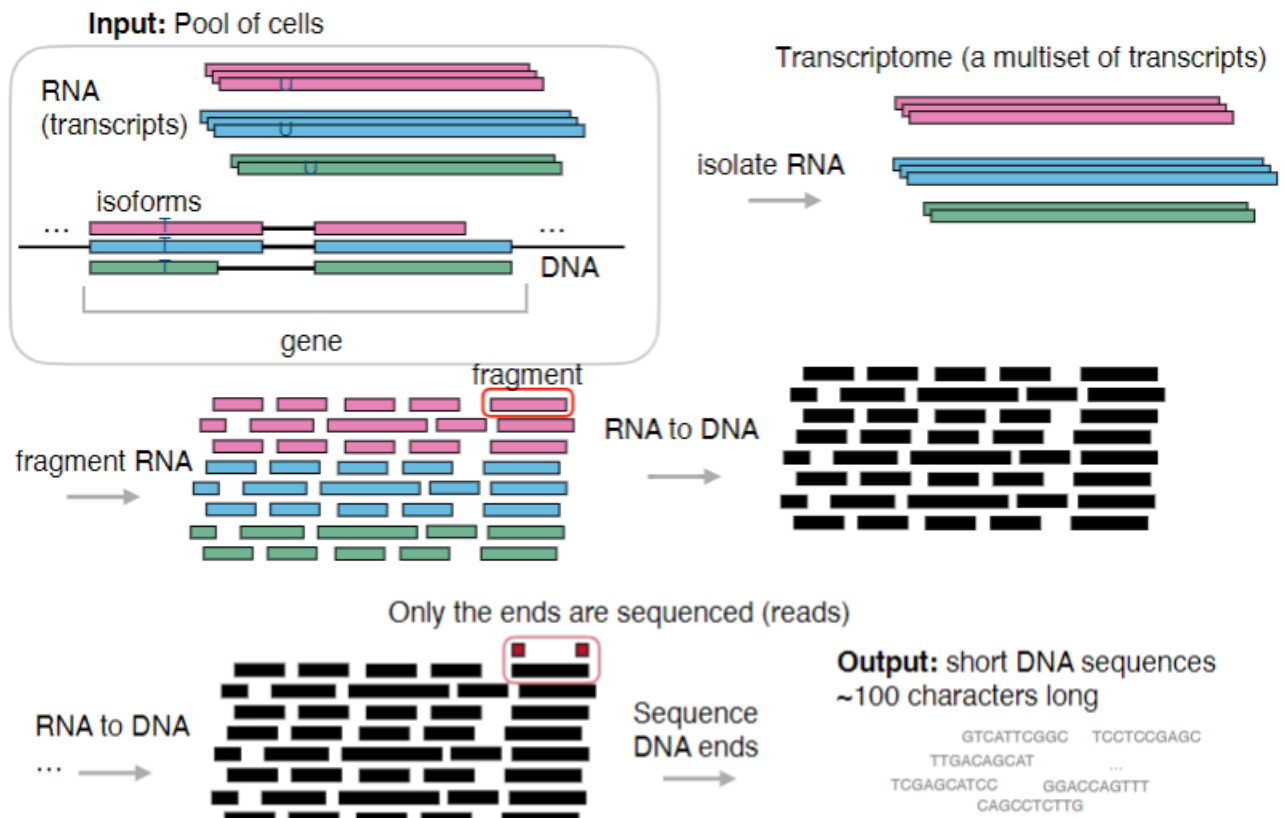
A lot of times, when we want to measure some trait, the easiest way is to reduce it to sequencing, then sequence the DNA, count occurrences, and analyze.

In this case, we want to measure RNA abundance, so we use the following pipeline.

RNA abundance \rightarrow cDNA Library Prep \rightarrow Sequence \rightarrow Estimate Abundances \rightarrow Differential Analysis

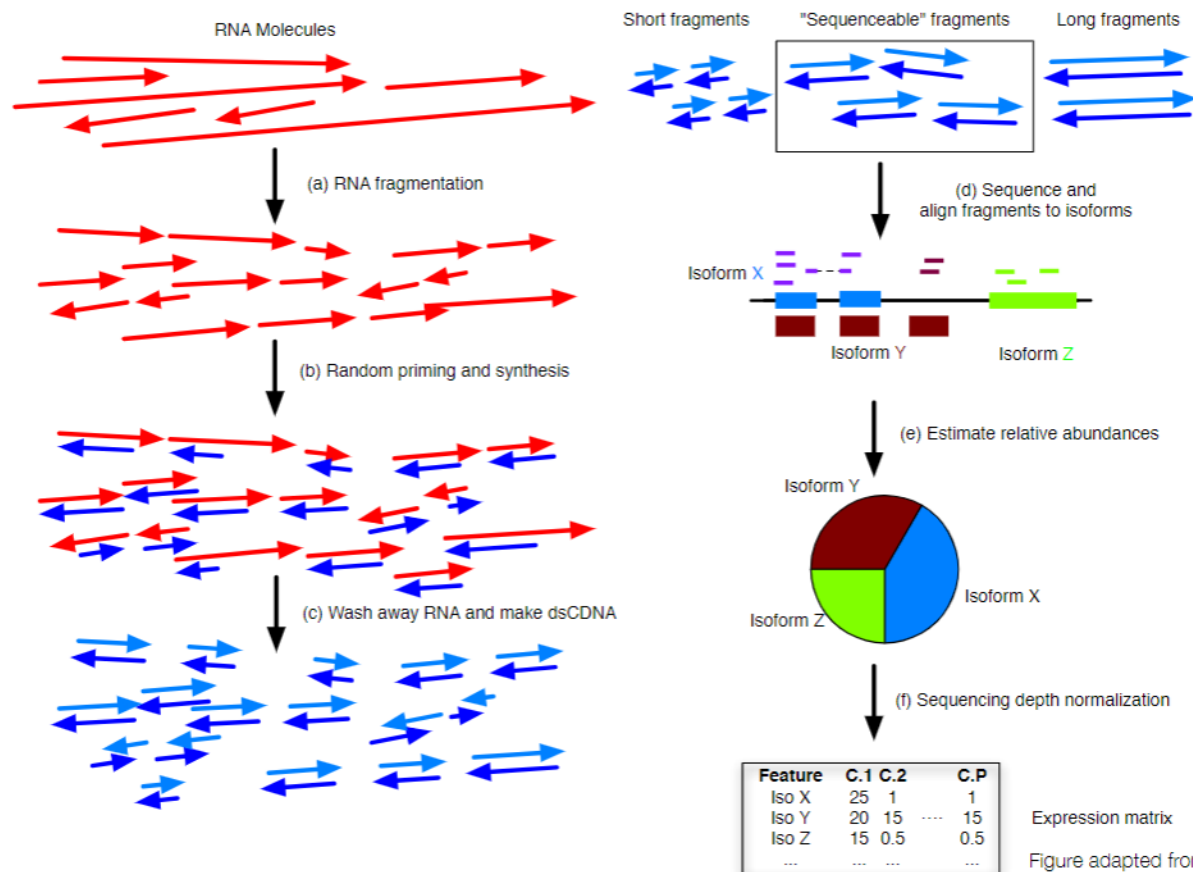
The below image depicts cDNA library prep, where the RNA is copied, isolated, fragmented, then converted to DNA.

- When converting to DNA, the information to where the fragment came from on the original pool is lost. (This is why the DNA is black in the image.)



- Only the ends of each fragment is sequenced because the sequencer does not give good output when sequencing long sections.
- Also, the middle is not necessarily needed since the sequenced ends have enough base pairs for us to figure out which other fragments it connects to.

Image of Converting RNA to DNA



- Notice that Isoform X and Isoform Y share the same DNA coding region in the image. This is referred to as an *ambiguous read* - the limit of RNA-seq. We want to know how much each X and Y there are, but since the two cover the same DNA region, sequencing cannot give you that information.
- What makes this worse is that in real life (where you are not the oracle), you don't know that the isoforms are overlapping.

RNA-seq quantification: a computational problem

Goal: given a known set on isoform targets (genes) and RNA-seq fragments, recover the distribution of RNA molecules.



All we want is that output pie chart that describes how common each isoform is.

Unlike DNA reads, where we can assume that all fragments appear at a relatively constant frequency, this is not true for RNA reads, where some isoforms may be more common than others. This makes solving the probabilities and thus the genome incredibly hard. (This is an unsolved problem.)

What is the "RNA Distribution"?

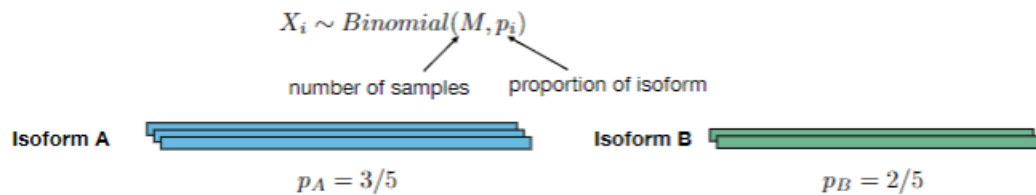
- In reality, there are a *finite* number of RNA molecules in each cell
- By nature of sequencing, we cannot directly sequence every molecule
- Instead, we mix a ton of cells together, isolate their RNA, then get a **relative measurement**.

Use cases for RNA-seq

- Tissue specific gene expression in *D. melanogaster*
- Cancer specific gene expression
- Genetic variation effects on gene expression and their relationships to tissues and complex traits

Binomial Sampling

- I'm going to sample M transcripts at random. Given the proportions below, what is a good model?



Now, suppose we have a very large M and many isoforms, where the proportion of each isoform is close to zero. What is a good model now?

- It turns out that the Poisson distribution is a good model.

$$X_i \sim \text{Poisson}(Mp_i)$$

– where M is the number of reads (samples) and p_i is the original isoform proportion.

- This makes a strong assumption about sampling, that all the isoform lengths are the same. This is not true in reality.
- We have to normalize the number of counts for each isoform based on the length of the isoform.
 - For example, if we have isoform A with length L , and isoform B with length $\frac{L}{2}$, then we would expect half as many reads in B than A . Therefore, to normalize the number of reads, we need to scale the raw count of reads of B by a factor of 2.

Transcript per Million (TPM)

$$\text{TPM}_i = \frac{X_i}{\tilde{l}_i} \cdot \left(\frac{1}{\sum_j \frac{X_j}{\tilde{l}_j}} \right) \cdot 10^6$$

where

- X_i is the number of counts
- \tilde{l}_i is the length
- $\sum_j \frac{X_j}{\tilde{l}_j}$ is the normalization constant
- 10^6 is a big number (the 'Million' part in TPM)

Assuming every site has equal probability of being sampled, what should the expectation of L squiggle be? Remember, not all fragments are of the same length. There's a fragment length probability in the expectation.

Suppose we have a transcript of length 5. Then:

- if length of fragment (F) = 3, then there are three different sites. (e.g., [0, 2], [1, 3], [2, 4])
- if length(F) = 2, then there are 4 sites.
- if length(F) = 1, then there are 5 sites.
- In general, (number of sites) = $l - \text{length}(F) + 1$

A simple model for RNA-seq

Conceptually:

(let n represent the read number)

1. Randomly select an isoform

$$I_n | p \sim \text{Categorical}(p)$$

2. Randomly select a fragment length

$$L_n | I_n = i_n \sim \text{Fragment length distribution}(\text{Length}(I_n))$$

3. Randomly select a position to generate a fragment from

$$R_n | L_n = l_n \sim \text{Uniform}(1, \text{Length}(I_n) - l_n + 1)$$

4. Observe and repeat

What is the probability of a particular arrangement $P(r_n, l_n, i_n)$? Hint: use the Bayes Theorem.

Answer:

$$P(r_n, l_n, i_n) = P(r_n | l_n, i_n) \cdot P(l_n | i_n) \cdot P(i_n)$$