

Alignment and Quantification of Gene Expression with Salmon

March 2023

Differential Gene Expression Analysis Workflow

Alignment and Quantification overview

Traditional Alignment

AIM: Given a reference sequence and a set of short reads, align each read to the reference sequence finding the most likely origin of the read sequence.

Alignment - Splicing aware alignment

Aligners: STAR, HISAT2

Alignment

- ▶ Traditional alignment perform base-by-base alignment
- ▶ It is (relatively) slow and computationally intensive

Alignment

- ▶ Traditional alignment perform base-by-base alignment
- ▶ Traditional alignment is (relatively) slow and computationally intensive

Alignment

- ▶ Traditional alignment perform base-by-base alignment
- ▶ Traditional alignment is (relatively) slow and computationally intensive

Switch to *quasi-mapping* (Salmon) or *pseudo-alignment* (Kallisto)

Why are Pseudo-alignment methods faster?

- ▶ These tools avoid base-to-base alignment of the reads
- ▶ ~ 20 times faster than the traditional alignment tools like STAR, HISAT2 etc
- ▶ Unlike alignment based methods, pseudo-alignment methods focus on transcriptome (~2% of genome in human)
- ▶ Use exact kmer matching rather than aligning whole reads with mismatches and indels

Quantification tools

- ▶ Broadly classified into two types ...
 - ▶ Alignment based:
 - ▶ Takes bam file as input, therefore reads must be mapped prior to quantification
 - ▶ quantifies using simple counting procedure
 - ▶ Pros: Intuitive
 - ▶ Cons: Slow and can not correct biases in RNAseq data
 - ▶ Tools: HTseq, SubRead etc.
 - ▶ Alignment-free:
 - ▶ Also called quasi-mapping or pseudoalignment
 - ▶ Starts from fastq files and base-to-base alignment of the reads is avoided
 - ▶ Pros: Very fast and removes biases
 - ▶ Cons: Not intuitive
 - ▶ Tools: Kallisto, Sailfish, **Salmon** etc

What is read quantification?

- ▶ **Quantification:** How many reads have come from a genomic feature?
 - ▶ genomic feature can be gene or transcript or exon, but usually gene

If we had mapped our reads to the genome (rather than the transcript sequences), our mapping would look like this:

We also know the locations of exons of genes on the genome, from an annotation file (e.g. GFF or GTF)

So the simplest approach is to count how many reads overlap each gene.

What is read quantification?

However, Salmon does not work this way. We have mapped to the transcript sequences, not the genome. Quantification is performed as part of the quasi-mapping process.

Salmon also takes account of biases:

- ▶ **Multimapping:** Reads which map equally well to multiple locations
- ▶ **GC bias:** Higher GC content sequences are less likely to be observed as PCR is not efficient with high GC content sequences.
- ▶ **Positional bias:** for most sequencing methods, the 3 prime end of transcripts are more likely to be observed.
- ▶ **Complexity bias:** some sequences are easier to be bound and amplified than others.
- ▶ **Sequence-based bias:** Bias in read start positions arising from the differential binding efficiency of random hexamer primers

Salmon workflow

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

Practical

1. Create and index to the transcriptome with Salmon
2. Quantify transcript expression using Salmon