RNA Sequencing Analysis

Sol Katzman BME 235 May 6, 2015

Today's main topics

- RNA sequencing basics (Illumina)
- Aligning (with or without gene model)
- Gene by gene coverage
- Comparing conditions or tissues
- Alternative splicing
- Comparing gene to gene
- Alternative technologies (PacBio)
- Gene assembly (linking exons)
- Small rna (miRNA)
- CLIPseq (RNA binding proteins), Riboprof, ...

DNAseq issues

- Data
 - How many reads do I need?
 - Quality of reads
 - PairedEnd vs SingleEnd sequencing
 - Length of reads
- Information
 - Mapping of reads to reference genome
 - Coverage of target region by reads
 - PCR duplicates

RNAseq issues

- All DNAseq issues, plus...
- Samples and library prep
 - Ribosomal RNA elimination
 - Strand specificity
 - Biological replicates
 - Infinite set of conditions, tissues to choose from
- Information
 - Mapping across splice junctions
 - RNA levels are not protein levels
 - 5' or 3' bias
 - Normalization among samples

RNAseq applications

- Annotated (model organism) genome
 - Relative gene expression
 - Alternative isoform expression
 - Differential gene or isoform expression
 - Novel (non-coding) gene expression
- Genome annotation
- Small rna (miRNA, etc.) expression
- Allele specific expression
- RNA editing



Quality is encoded:

- C = Q34
- D = Q35

. . .

Q = -10 log₁₀ (prob (error)) Q20: 1/100 prob error Q30: 1/1000 prob error Quality distribution by position in read huESC_CBP2_68_2_frac1000



DNAseq mapping

- Mapping issues
 - Errors in the reads
 - Errors in the reference genome
 - Polymorphisms
 - Repetitive regions in the genome
- Mapping algorithms
 - Allow for mismatches, multiple matches, low quality positions in the read
 - Must be fast, parallelizable

RNAseq mapping

- All DNAseq mapping issues, plus
 - Splice junction mappings
 - High copy genes (rRNA, tRNA)
 - Paralogs
 - Pseudogenes
 - Degradation products (introns)
- Mapping algorithms
 - Map to genome, pre-split the reads (tophat1,2)
 - Map to transcriptome (gene model) (tophat2)
 - On-the-fly split reads during mapping (star)

Tophat2 (bowtie2)



tophat2

- Optional genome model for pre-mapping
- Reads split into thirds
- Need long enough reads to be split for splice junction finding (75bp)
- Where is the junction? Canonical splice sites (GU/AG) favored
- Coverage pileup method for finding splices too
- Iterative, more reads per run better (less parallelizable)
- Somewhat slow

STAR

Schematic representation of the Maximum Mappable Prefix search in the STAR algorithm for detecting (a) splice junctions, (b) mismatches and (c) tails.



STAR

- Optional gene model included in mapping target (splice junction database)
- On-the-fly splicing
- Soft clipped mappings
- Post-filter based on splice sites (GU/AG...) (and reverse complement?)
- Can be run iteratively (internally or externally)
- Very fast (20X to 50X tophat)
- Separate output of "chimeric" mappings

Other mappers

- BWA
- MapSplice (TCGA)
- SpliceMap
- Bowtie "genomic" target only
 - for QC, for filtering
 - peLength estimation
 - mapping to transcriptome only (RSEM)
 - mapping to isoform alternatives (SpliceTrap)
 - it is the mapping engine for Tophat

RNAseq split reads



RNAseq exons coverage



RNAseq issues (1)

- Ribosomal RNA elimination from library
 - Over 90% of RNA in the cell
 - polyA selection
 - 3' bias
 - Pol II transcripts only
 - riboMinus, riboZero, (polyT priming)
 - 5' bias?
 - pre-mRNA
 - QC: Filter reads against repetitive genomic elements before mapping
 - QC: 5' to 3' bias after mapping

RiboZero positional bias



relative position in transcript from 5' to 3' end (18240 transcripts)

Ribo Prof Exons positional bias



relative position in transcript from 5' to 3' end (12668 transcripts)

Ribo Prof CDS positional bias



RNAseq issues (2)

- Strand specificity
 - Detect antisense transcription of known genes
 - Detect transcription of unknown genes, which strand?
 - Splice site bias (vs. revComp) for mapper

RNAseq issues (3)

- PCR duplicates
 - "PCR is not your friend"
 - Wastes sequencing \$\$
 - Wastes computational time, storage
 - 10-12 cycles max! But extracting RNA is hard.
 Extracting non-rRNA is harder still.
 - Detection and bioinformatic removal after mapping. Mitigate removal with:
 - Paired-end reads for an extra degree of freedom (if randomly digested)
 - Internal random barcodes in adaptors before PCR

RNAseq issues (other)

Mitochondrial genes

Muscle cells

- "Uniquely" mapping reads only?
 - Determined by mapping parameters
 - Can lose common genes with paralogs (mouse GAPDH)
 - One best mapping (if scored) is reasonable
- Lengths of fragments vs. read length
 - Longer reads for more splice junctions (mapper must support multiple junctions per read)
 - Longer fragments for more mapping information (no overlap read1/read2)

Today's main topics

- RNA sequencing basics (Illumina)
- Aligning (with or without gene model)
- Gene by gene coverage
- Comparing conditions or tissues
- Alternative splicing
- Comparing genes
- Alternative technologies (PacBio)
- Gene assembly
- Small rna (miRNA)
- CLIPseq (RNA binding proteins), Riboprof, ...

What is a Gene?

- Genes have many isoforms
- Canonical transcript?
- Any exon in any transcript?
 - Count reads? Read starts? Intron/exon reads?
 - Count total coverage, divide by read length to estimate number of reads in a "gene"

Today's main topics

- RNA sequencing basics (Illumina)
- Aligning (with or without gene model)
- Gene by gene coverage
- Comparing conditions or tissues
- Alternative splicing
- Comparing genes
- Alternative technologies (PacBio)
- Gene assembly
- Small rna (miRNA)
- CLIPseq (RNA binding proteins), Riboprof, ...

Comparing expression

• FPKM: fragments per kilobase of gene length per million mapped reads

from Dillies et al. Brief Bioinform. 2013:

- Normalization of RNA-seq data in the context of differential analysis is essential in order to account for the presence of systematic variation between samples as well as differences in library composition.
- The Total Count and RPKM normalization methods, both of which are still widely in use, are ineffective and should be definitively abandoned in the context of differential analysis.
- Only the DESeq and TMM normalization methods are robust to the presence of different library sizes and widely different library compositions, both of which are typical of real RNA-seq data.

DESeq

- Uses replicates per "condition"
- Estimates dispersion as a function of mean expression for a gene
- Uses median relative expression to normalize sample gene counts
- Uses dispersion and normalized counts to determine significance of differential expression between conditions

DESeq fitting dispersion estimates



per Gene dispersion estimates

normalized means

DESeq normalization



genes ranked by expression ratio to all samples (N = 6119)

RNAseq diffex with replicates



Today's main topics

- RNA sequencing basics (Illumina)
- Aligning (with or without gene model)
- Gene by gene coverage
- Comparing conditions or tissues
- Alternative splicing
- Comparing genes
- Alternative technologies (PacBio)
- Gene assembly
- Small rna (miRNA)
- CLIPseq (RNA binding proteins), Riboprof, ...

Alt Splicing tools

- MISO
- Splicetrap
- CuffDiff
- rMATs
- •

MISO

- Uses mapped reads as input
- Bayesian posterior estimate of "Percent Spliced In" (PSI) for a skipped exon
- Nice plotting tool

SE21101





SE00574



RNAseq SUMMARY

- PairedEnd, strand-specific library prep is preferred
- Biological replicates are a must
- Avoid PCR duplicates

