

Understanding the Microbiome: *Metatranscriptomics*

Marcus Claesson

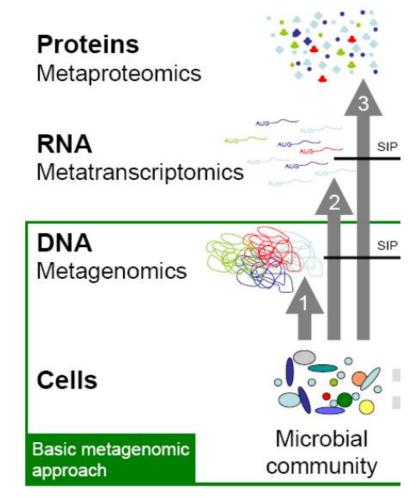
APC Microbiome Symposium 2015

Metatranscriptomics

Definition

(genetics, ecology)
A branch of
transcriptomics that
studies and correlates,
the transcriptomes of a
group of interacting
organisms or species

--Wiktionary--



Warnecke & Hugenholtz (2007) Genome Biology

Sequence-based 'omics' technologies

Metatranscriptome:

- Protein-coding RNA (mRNA)
- Non-coding RNA (rRNA, tRNA, regulatory RNA, etc)

Metatranscriptomics studies:

- Community functions
- Response to different environments / treatments; differential gene expression
- Regulation of gene expression

Metagenomics shotgun sequencing

- Encoded potential functions of the microbiota
- What CAN they do?

Microbiota compositional analysis

- Quantification of the ubiquitous 16S rRNA gene
- WHAT organisms are there?



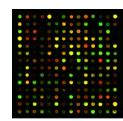
Metatranscriptomics cDNA/mRNA sequencing

- Microbial gene expression at certain times and/or locations
- What ARE they doing?

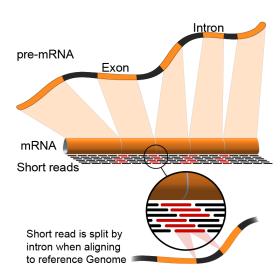
Measuring the transcriptome

CDNA clone libraries + Sanger sequencing

Microarrays: hybridizing mRNA onto cDNA/oligonucleotide proves on glass slide

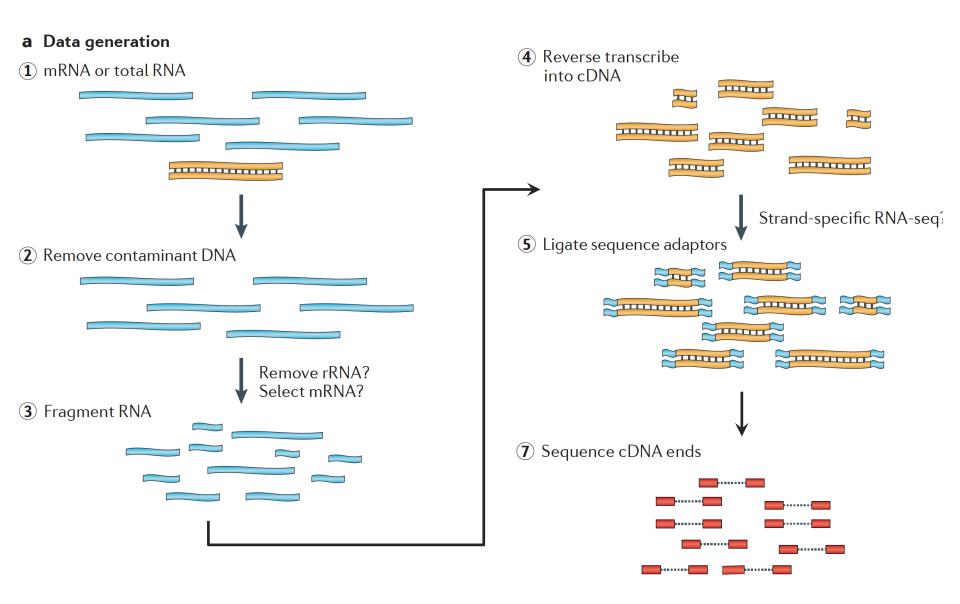


 RNA-seq enabled by nextgeneration sequencing technologies
 usually Illumina HiSeq



RNA-seq superior to microarrays for gene expression analysis of microbial communities

From RNA to sequence data



Challenges and considerations

Wet lab

- Deplete host RNA or not? Dual RNA-Seq an option
- Instability of RNA (half-lives of minutes)
- High rRNA content in total RNA (mRNA<5% of total RNA)</p>
- Single-end or Paired-end?
- Stranded cDNA transcription or not?
- How much sequencing per sample?
- Technical replicates not necessary any more, but biological are
- Avoid batch effect -> randomize samples across runs!



Bioinformatics

- General challenges with short reads and large data size
- Lack of metagenome/genome reference
- Statistical considerations
- Assemble or map reads, or both?



http://www.nwfsc.noaa.gov/index.cfm

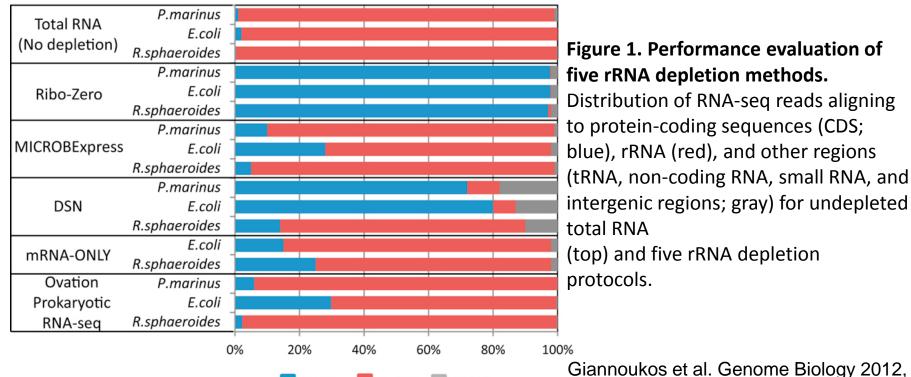
rRNA removal methods

Majority of bacterial mRNA is not polyadenylated, can't be isolated using oligo-dT selection

Subtractive Hybridization

Exonuclease Digestion





% Other

% CDS

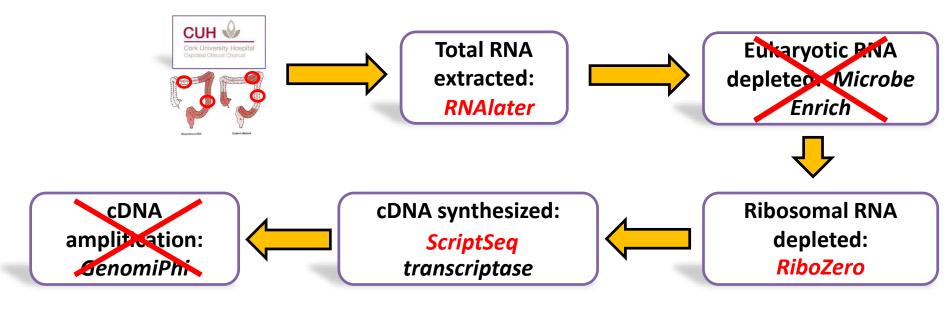
% rRNA

Experimental pipeline

- Inflamed/uninflamed colonic biopsies
 - Pilot project: 12 CD & 6 UC
 - Main project: 60 CD, 86 UC & 30 HC



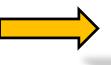
PostDoc Fmilio Laserna



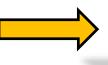


RNA-Seq library prep:

Illumina TruSeq

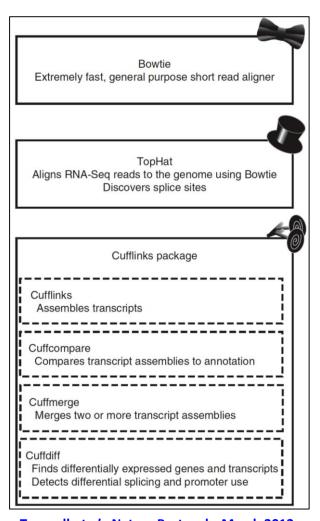


RNA-Seq:
Illumina HiSeq
2500



15 million
125 bp PE reads per sample

Differential Gene Expression



Trapnell et al., Nature Protocols, March 2012

<u>Bowtie2</u> and Bowtie use Burrows-Wheeler indexing for aligning reads. With bowtie2 there is no upper limit on the read length

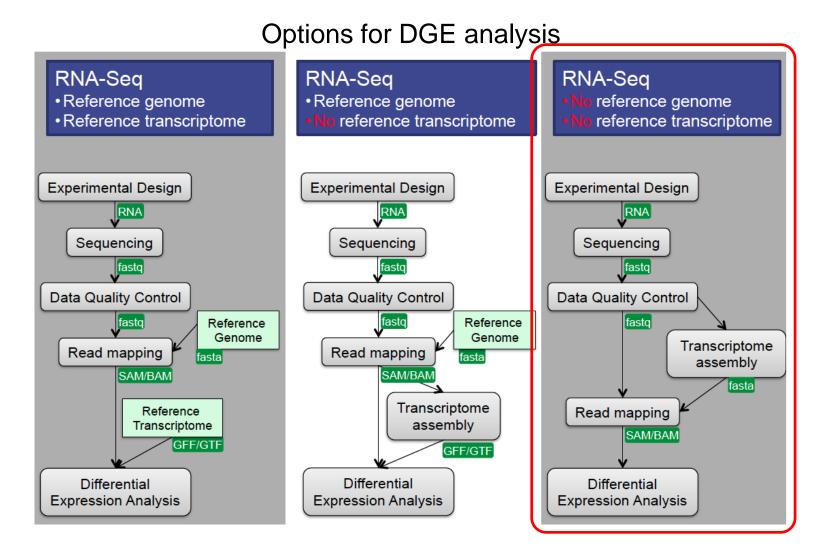
<u>Tophat2</u> uses Bowtie2 to align reads in a splice-aware manner and aids the discovery of new splice junctions

The <u>Cufflinks2 package</u> has 4 components, the 2 major ones are listed below -

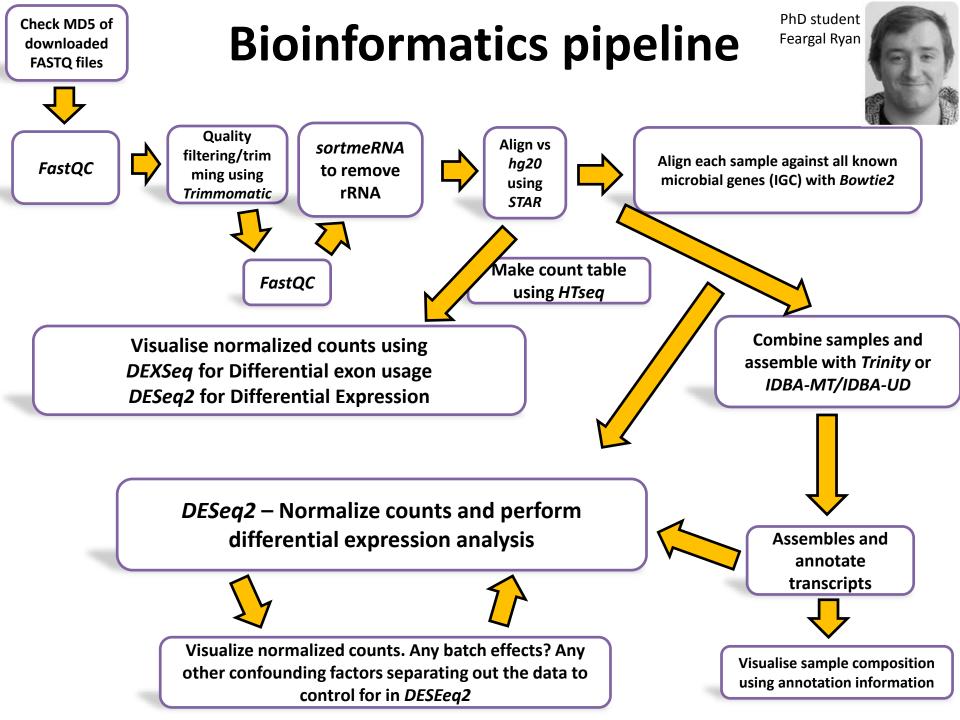
<u>Cufflinks2</u> does both *de novo* and reference-based transcriptome assembly

<u>Cuffdiff2</u> does statistical analysis and identifies differentially expressed transcripts in a simple pairwise comparison, and a series of pairwise comparisons in a time-course experiment

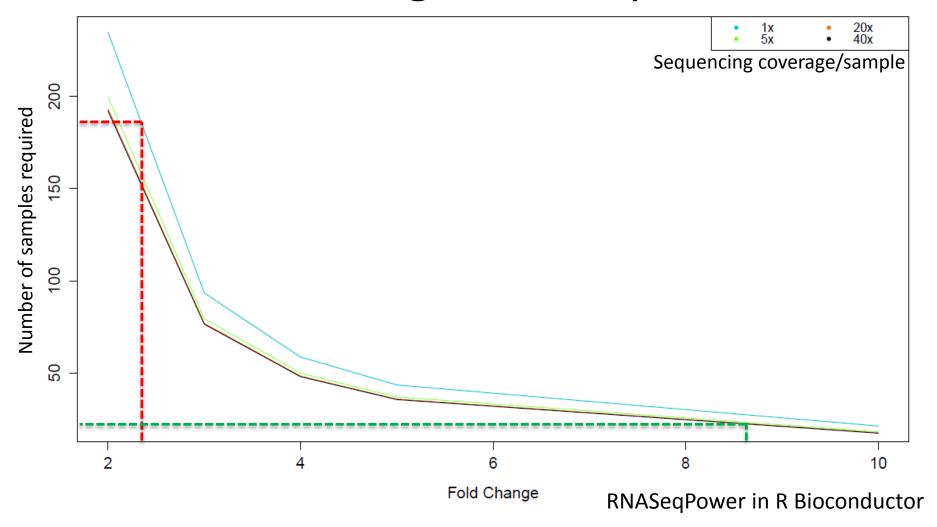
Differential Gene Expression



Want to learn more about the formats?
https://genome.ucsc.edu/FAQ/FAQformat.html



Power calculation: Read coverage vs. sample size



What the Pilot Study taught us

- Differential gene expression for inflamed and non-inflamed in <u>both UC and CD</u>
 higher *n* required to confirm
- Use a more effective rRNA depletion method
- Don't deplete host mRNA if also interested in it
- No extra cDNA amplification needed
- High n more important than sample coverage
- Experimental & bioinformatics pipelines
- SIRG & Second Genome study:
 - 200 subjects
 - Genotype, disease activity, diet, medication