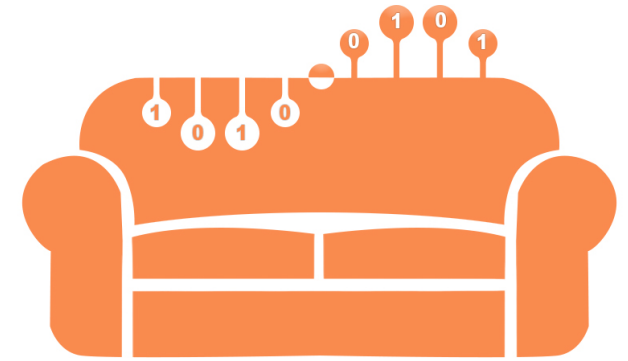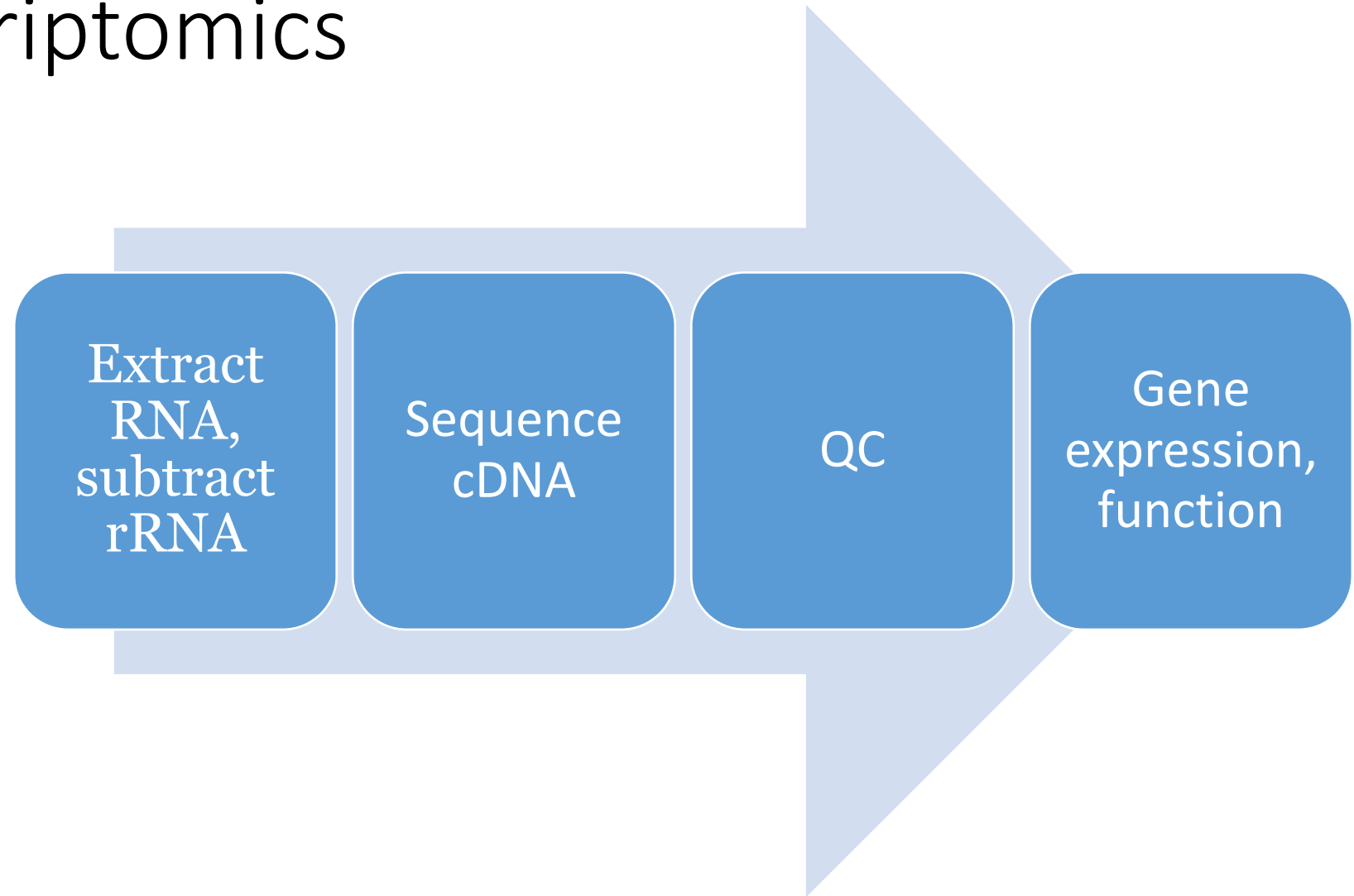# CBC Data Therapy

Metatranscriptomics Discussion
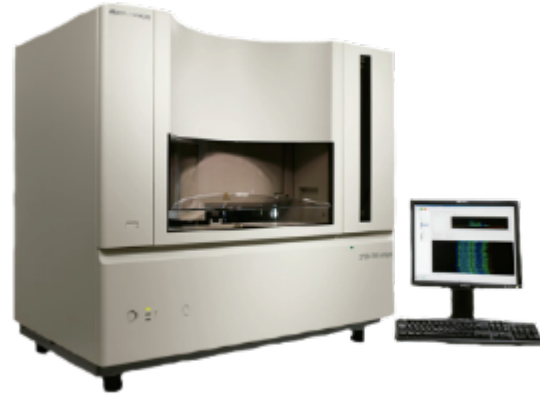
Computational Biology Core

UCONN

UNIVERSITY OF CONNECTICUT

# Metatranscriptomics

# Sequencing

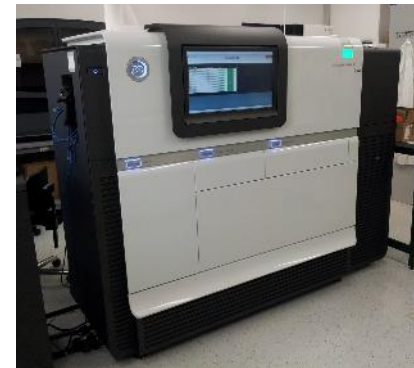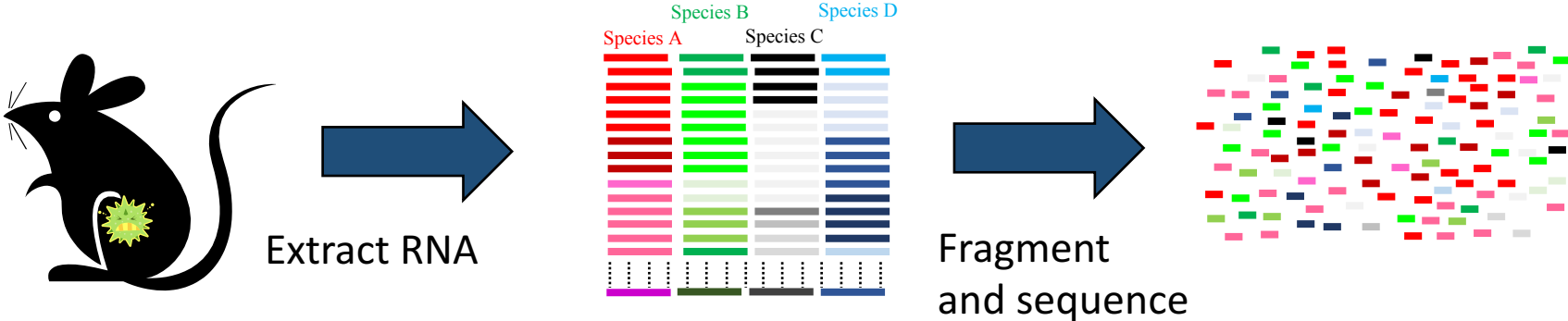
Sanger


Ion Torrent


Roche 454


Illumina *Seq


Pacific Biosciences


Nanopore

# Metatranscriptomics focuses on community activity

Metatranscriptomics exploits RNA-Seq to determine which genes and pathways are being actively expressed within a community

Genes involved in pathways associated with cell wall biogenesis



Relative abundance



Relative abundance and contributing taxa

Metatranscriptomics can reveal active *functions* (knowing the taxa responsible is unimportant)

It can also reveal which taxa are responsible for the active functions

# Metatranscriptomics through RNA Seq



Extract RNA

Species A
Species B
Species C
Species D

Fragment and sequence

Align reads to known transcripts to obtain relative expression

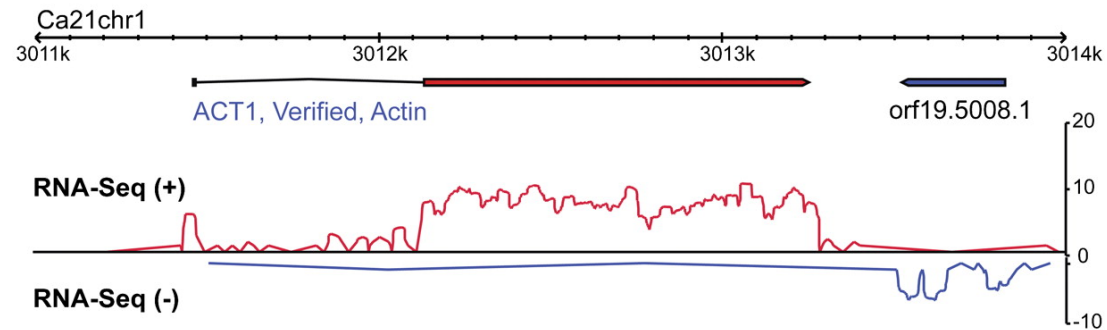| Gene | |
|------|---|
| Gene A1 | 6 |
| Gene A2 | 3 |
| Gene A3 | 3 |
| Gene A4 | 1 |
| Gene B1 | 2 |
| Gene B2 | 2 |
| Gene B3 | 6 |
| Gene C1 | 4 |
| Gene D4 | 3 |

RNA-Seq is the unbiased sequencing of an RNA sample to yield a digital readout of the relative expression of transcripts within a sample

Typically applied to organisms with a reference (sequenced) genome, microbiome applications face a number of challenges

III          IV

# Metatranscriptomics: Challenges

In a typical RNA-Seq experiment applied to a single eukaryotic organism, mRNA is isolated. After fragmentation and sequencing, reads are mapped to a reference genome using standard software such as MAQ and BWA to provide: 1) support that the transcript is expressed; 2) the relative abundance of the transcript; and 3) the presence and abundance of isoforms
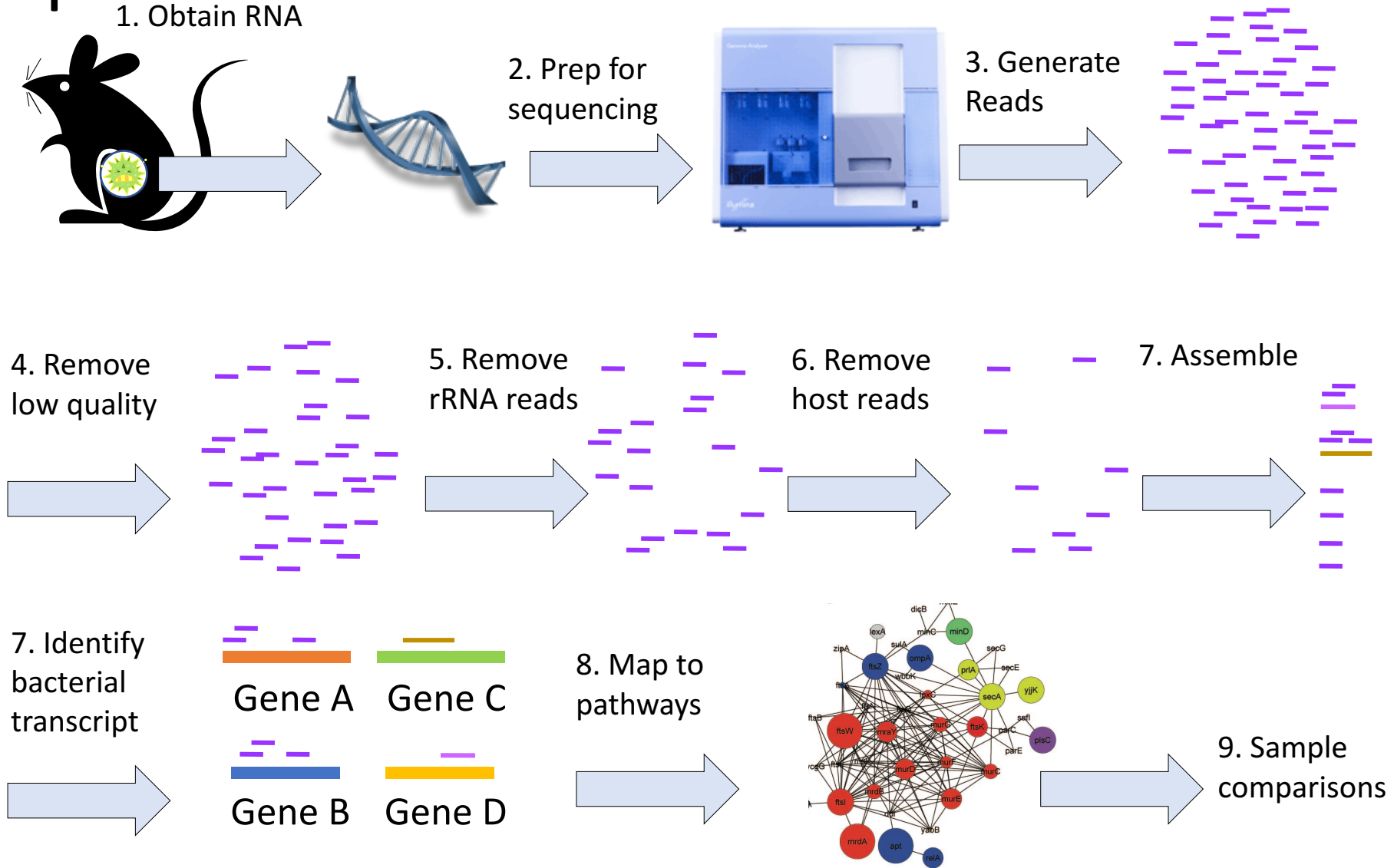
Resource

Comprehensive annotation of the transcriptome of the human fungal pathogen *Candida albicans* using RNA-seq

Vincent M. Bruno,[1] Zhong Wang,[2] Sadie L. Marjani,[3] Ghia M. Euskirchen,[4] Jeffrey Martin,[2] Gavin Sherlock,[4,5] and Michael Snyder[1,4,5]

[1]Department of Molecular, Cellular, and Developmental Biology, Yale University, New Haven, Connecticut 06520, USA; [2]DOE Joint Genome Institute (JGI), Walnut Creek, California 94598, USA; [3]Department of Genetics, Yale University School of Medicine, New Haven, Connecticut 06520, USA; [4]Department of Genetics, Stanford University Medical School, Stanford, California 94305-5120, USA



For microbiome samples we have the following problems:
a) Lack of a polyA signal makes it difficult to isolate bacterial mRNA and resulting in (massive) rRNA contamination
b) Environmental microbiome samples lack reference genomes making it difficult to map reads back to their source transcripts

III       IV

# A typical metatranscriptomic analytical pipeline

1. Obtain RNA

2. Prep for sequencing

3. Generate Reads

4. Remove low quality

5. Remove rRNA reads

6. Remove host reads

7. Assemble

7. Identify bacterial transcript

Gene A    Gene C

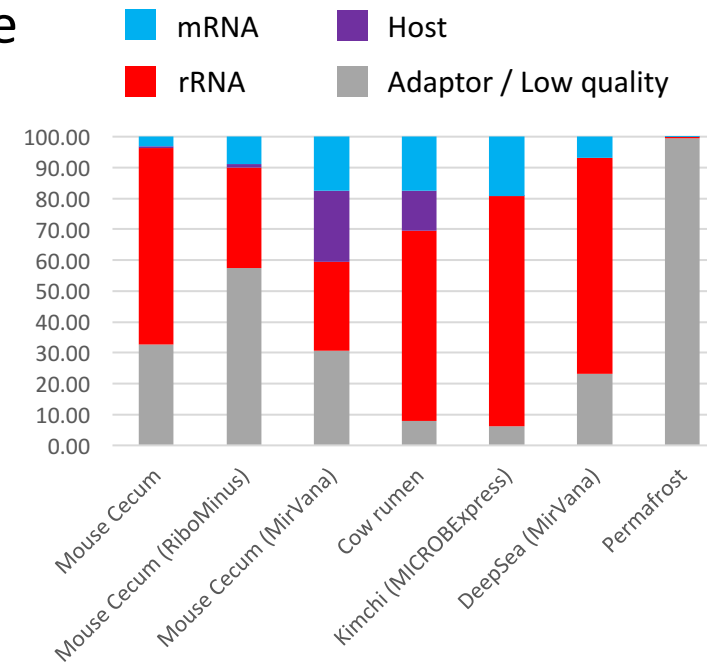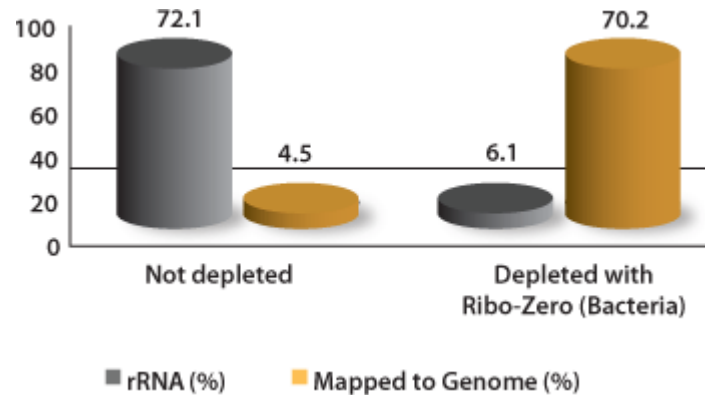Gene B    Gene D

8. Map to pathways

9. Sample comparisons

# Preparing sample for sequencing

Bacterial mRNA's lack a polyA tail so how to remove abundant rRNA species?

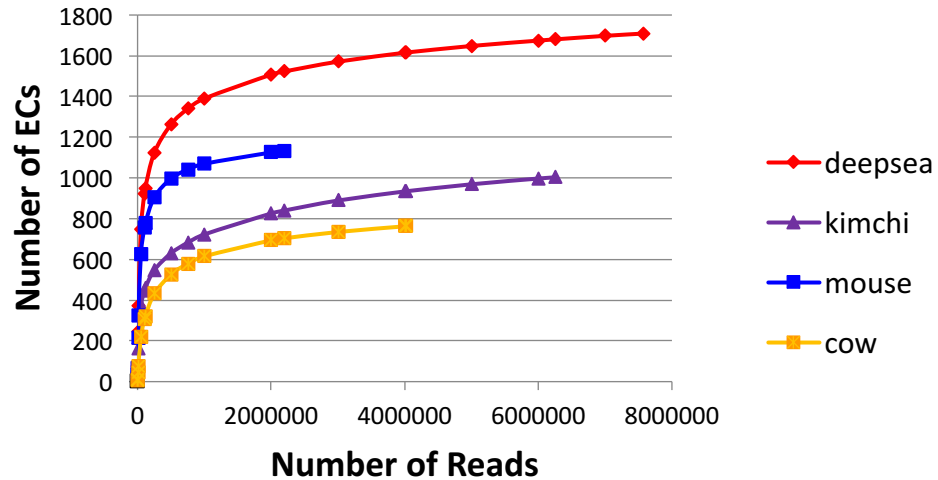Once RNA has been extracted, several kits are available to remove rRNA – need 500ng-2.5ug RNA/sample

Ribo-Zero (Illumina) provides reasonable success



Host mRNAs can also prove challenging – can also be informative!

# Generating reads

How many reads are "enough"?



20 million/sample mRNA

While PB and MiSeq provide long reads useful for annotation HiSeq (or NextSeq) provide sequencing depth and offer possibility of multiplexing

# Read processing - filtering

To identify reads derived from mRNA bioinformatics pipelines need to be in place that remove contaminating reads:

Low quality - *Trimmomatic*
Adaptors – *Trimmomatic* & *Cross_Match*
Host - *BWA*
rRNA – *BLAT / Infernal*

*Trimmomatic* uses a sliding window approach from the 5` end to identify low quality regions which are then trimmed from the 3` end. Reads < 36 bp are discarded

# Read processing - Assembly

contig1

contig2

attagcggcgattttcggcgatcttatcttgatctgggcgcgtatcggtagcgtagcgattcgtagc
attagcggcgattttcggcgatcttatcttgatctgggcgcgtatcggtagcgtagcgattcgtagc
attagcggcgattttcggcgatcttatcttgatctgggcgcgtatcggtagcgtagcgattcgtagc

Assembly improves annotation accuracy



Trinity appears to provide best performance in terms of reads that can be annotated

Chimera's, misassembled contigs, can become a problem due to reads derived from orthologs from different species

# Read processing – functional annotation

One solution is to work in peptide space and use *BLASTX* to search protein databases - this is very time consuming and requires cloud/cluster computing

Other solutions USEARCH/VSEARCH or DIAMOND (issues over quality and cost)



Even with BLAST many reads remain unannotated

Can be improved with longer read length

# Read processing – converting mappings to expression

To normalize expression levels to account for differences in gene length, read counts are converted to **Reads per kilobase of transcript mapped (RPKM)**

Expression is biased for gene length (longer transcripts should have more reads) to normalize, reads are converted to Reads per Kilobase of transcript per million reads mapped

| # | | RPKM |
|---|---|---|
| 8 | | 1 |
| 24 | | 6 |
| 8 | | 12 |

$$RPKM_{geneA} = 10^9 \, C_{geneA} / NL$$

$C_{geneA}$ = number of reads mapped to geneA
N = total number of reads
L = length of transcript in units of Kb

Several software tools available to do mapping and calculate normalized expression measurements across different samples including Bowtie and Cufflinks

# Read processing – taxonomic annotation

Alignment based methods such as BLAST and BWA can fail where we lack suitable reference genomes – particularly for short read datasets where assignments may be ambiguous

Compositional methods (e.g. nt frequency, codon bias) offer alternative strategies



Here a sequences is classified into frequencies of 3-mers



Nearest neighbours methods then try to assign a sequence to the genome with the closest distribution



NBC

MetaCV

RITA

# Visualizing results

Metabolic pathways are among the most highly conserved and best characterized systems
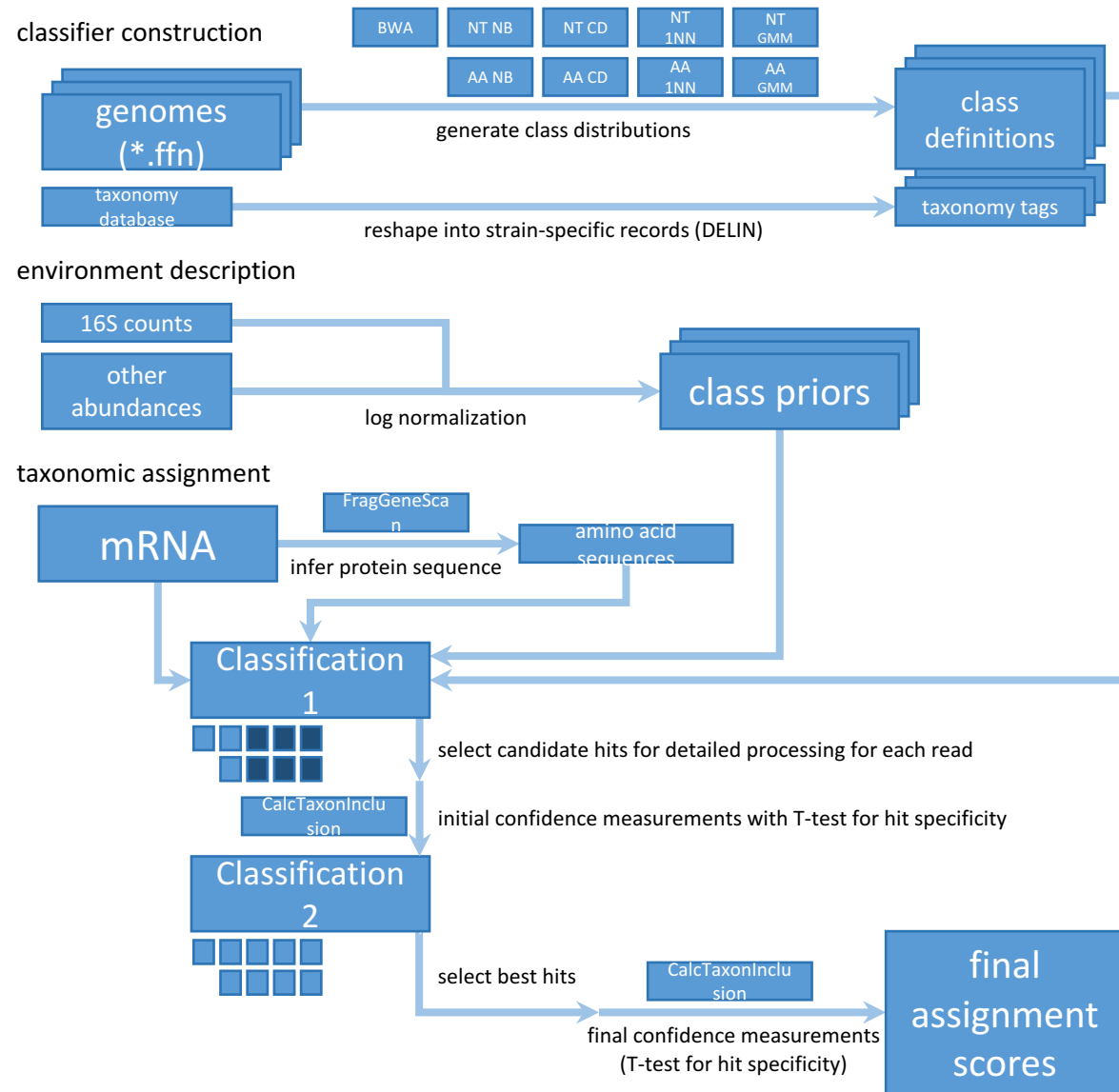


MG-RAST and MEGAN are automated metagenomic annotation tools that rely on KEGG
A major problem with KEGG pathway definitions is that the boundaries of pathways are arbitrarily defined and links between pathways (i.e. functional relationships) can be lost

# Read processing – Gist

Gist is a computational pipeline for accurate assignment of reads to individual species

Integrates several methods, but uniquely assigns different weights to methods for each genome

Can also take in expected sequence distributions (e.g. based on 16S rRNA surveys)

classifier construction

| BWA | NT NB | NT CD | NT 1NN | NT GMM |
| AA NB | AA CD | AA 1NN | AA GMM |

genomes (*.ffn)

generate class distributions

class definitions

taxonomy database

reshape into strain-specific records (DELIN)

taxonomy tags

environment description

16S counts

other abundances

log normalization

class priors

taxonomic assignment

mRNA

FragGeneScan

infer protein sequence

amino acid sequences

Classification 1

select candidate hits for detailed processing for each read

CalcTaxonInclusion

initial confidence measurements with T-test for hit specificity

Classification 2

select best hits

CalcTaxonInclusion

final confidence measurements (T-test for hit specificity)

final assignment scores

# Statistical considerations

There is no dedicated software or statistical tool for statistical comparisons of metatranscriptomic datasets
- Number of biological replicates? (preferably at least three)
- Differential expression of individual genes
- Gene set enrichment analyses

Ultimately metatranscriptomics could be viewed as hypothesis generating requiring subsequent targeted validation
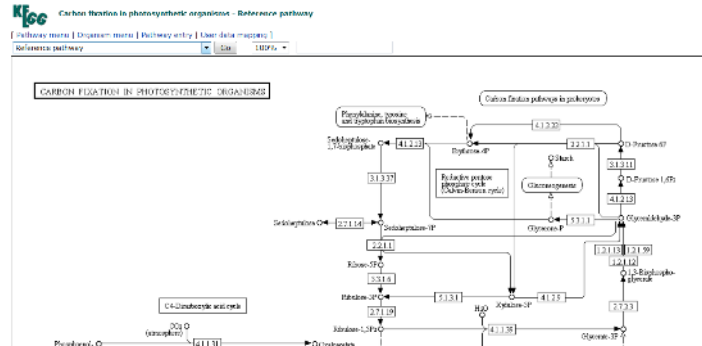
# Statistical considerations

While there are no dedicated tools for metatranscriptomics analyses, tools used for RNA Seq offer potential
- DESeq, EdgeR, ALDEx
- Alternatively simply rely on fold change (Gfold)
- Challenges include which genes to include (minimum RPKM?)

Differentially expressed genes can be subsequently analysed through Gene Set Enrichment Approaches

# Resources (Function)


KEGG


CARD


UniProtKB


Gene Ontology