

ENSEMBL
PhytoPathDB
SolGenomics
TAIR

AraPort

Data • Tools Provided • USPs

ENSEMBL

url: ensembl.org



Family of DBs
(Human is the flagship)

Assemblies • Genes • Annotations • Expression • Variants •
Comparative Genomics

ENSEMBL PLANTS

url: plants.ensembl.org

Login/Re

 ▾ Search Ensembl Plants.

39 Plant Species

Ensembl Plants • HMMER | BLAST | BioMart | Tools | Downloads | Documentation | Website help

Arabidopsis thaliana (TAIR10) ▾ Location: 1:8,001-18,000

Location-based displays

- Whole genome
- Chromosome summary
- Region overview
- Region in detail**
- Comparative Genomics
 - Alignments (image)
 - Alignments (text)
 - Region Comparison
 - Synteny
- Genetic Variation
- Resequencing
- Linkage Data
- Markers
- Other genome browsers
 - TAIR®
 - Phytozome

Configure this page

Add your data

Export data

Share this page

Bookmark this page

Ensembl Plants is produced in collaboration with Gramene

Chromosome 1: 8,001-18,000

Region in detail

Location: 1:8,001-18,000

Gene:

ENSEMBL PLANTS

url: plants.ensembl.org

39 Plant Species

[Login/Register](#) Search Ensembl Plants...

Arabidopsis thaliana (TAIR10) ▾

Location-based displays

- Whole genome
- Chromosome summary
- Region overview
- Region in detail**
- Comparative Genomics
 - Alignments (image)
 - Alignments (text)
 - Region Comparison
 - Synteny
- Genetic Variation
 - Resequencing
 - Linkage Data
 - Markers
- Other genome browsers
 - TAIR
 - Phytozome

Configure this page

Add your data

Export data

Share this page

Bookmark this page

Ensembl Plants is produced in collaboration with Gramene

Location: 1:8,001-18,000

Chromosome 1: 8,001-18,000

Region in detail

Contigs TAIR gene

Gene Legend

protein coding RNA gene

Location: 1:8001-18000

Gene:

Drag>Select: ↕

EST_Cluster_r.(Dicot)

EST_Cluster_r.(Arabidopsis)

Contigs

TAIR gene

EST_Cluster_r.(Arabidopsis)

EST_Cluster_r.(Dicot)

Variant - All sources

SV - Smaller variants

%GC

Reverse strand

Variant Legend

- Misense variant
- Synonymous variant
- 3 prime UTR variant
- Upstream gene variant
- Intron variant
- CNV

Structural Variant Legend

Gene Legend

Protein Coding

protein coding

There are currently 85 tracks turned off.

Ensembl Plants Arabidopsis thaliana version 84.10 (TAIR10) Chromosome 1: 8,001 - 18,000

ENSEMBL PLANTS

url: plants.ensembl.org

Gene

- trees – gene families, across plant kingdom, wider!
- secondary structure
- orthologs
- GO terms
- PO terms
- genomic alignments
- gene expression (outlinks)
- variants

Transcript

- exons
- splice variants
- cDNA

Protein

- domain
 - prints, PFAM, PROSite, Panther, Gene3D
- variants

ENSEMBL PLANTS

url: plants.ensembl.org

Tools

HMMER



This website uses cookies. By continuing to browse this site, you are agreeing to the use of our site cookies. To find out more, see our Terms of Use.

OK

Score Taxonomy Domain Download

PHMMER Results

Sequence Matches and Features

Plam SH2 163

hit coverage
hit similarity

disorder coiled-coil tm & signal peptide

Show hit details

BLAST

ENSEMBL PLANTS

url: plants.ensembl.org

Tools

BIOMART

New Count Results

Dataset
Brachypodium distachyon genes
(v1.0 (2010-02-Brachy1.2))

Filters

- GO term accession [e.g. GO:0050789]: [ID-list specified]
- GO term name [e.g. regulation of biological process]: regulation of biological process
- Chromosome/scaffold: 5

Attributes

- Gene stable ID
- Transcript stable ID

Dataset
[None Selected]

Please restrict your query using criteria
(If filter values are truncated in any lists, hover over the list item)

- REGION:
- GENE:
- GENE ONTOLOGY:
- PLANT ONTOLOGY:
- ENVIRONMENT ONTOLOGY:
- GRAMENE TAXONOMIC ONTOLOGY:
- GROWTH STAGE ONTOLOGY:
- MULTI-SPECIES COMPARISONS:
- PROTEIN DOMAINS:
- VARIATION:

ENSEMBL FUNGI, BACTERIA, PROTISTS

url: fungi.ensembl.org, bacteria.ensembl.org, protists.ensembl.org

Gene

- ~~trees — gene families, across plant kingdom, wider!~~
- secondary structure
- orthologs
- ~~GO terms~~
- ~~PO terms~~
- ~~genomic alignments~~
- ~~gene expression (outlinks)~~
- variants

Transcript

- exons
- splice variants
- cDNA

Protein

- domain
 - prints, PFAM, PROSite, Panther, Gene3D
- variants

PhytoPath

url: phytopath.org

Based on ENSEMBL Portal for Plant Pathogens

The screenshot shows the PhytoPath website homepage. At the top is a green header bar with the PhytoPath logo, navigation links (Home, News, Pathogens, Contact, About, Advanced Search), and a search bar. Below the header is a large banner image of plant roots. The main content area has a yellow background. It features a "Welcome to PhytoPath" section with a brief introduction, a "New feature: Advanced Search" section with a screenshot of the search interface, a "Training Course in Bioinformatics of Plants and Plant Pathogens" section with details about a May 2016 course, and a "This release" section with information about the April 2016 release.

Welcome to PhytoPath

PhytoPath is a new bioinformatics resource that integrates genome-scale data from important plant pathogen species with literature-curated information about the phenotypes of host infection. Using the [Ensembl Genomes](#) browser, it provides access to complete genome assembly and gene models of priority crop and model-fungal, oomycete and bacterial phytopathogens. PhytoPath also links genes to disease progression using data from the curated [PHI-base](#) resource.

New feature: Advanced Search

The Advanced Search lets you build complex queries to find genes across all the PhytoPath species. You can search for genes with specific pathogenic phenotypes, gene ontology annotation or protein domains, and download your results for further analysis.

Training Course in Bioinformatics of Plants and Plant Pathogens

A training course on the Bioinformatics of Plants and Plant Pathogens will be held at EBI from 23rd-25th May, 2016. Please click [here](#) for more details and to register.

This release

The current release of PhytoPath is built from the 31th release (April 2016) of Ensembl Genomes and version 4.1 of [PHI-base](#) and was released on April 2016. For all species, available data includes genome sequence and gene models, functional annotation, and protein-based comparative analysis with other fungal/oomycete species (provided by Ensembl Genomes),

Literature Information,
disease progression info from PHI-BASE

SolGenomics

url: solgenomics.net

Community Site for Solanaceae

Sol Genomics Network Search Maps Genomes Tools About

Potato Genome
Genome sequence and analysis of the tuber crop potato

Browse the potato genome
Find sequences by similarity
Download the annotations
Download the Genome Sequence

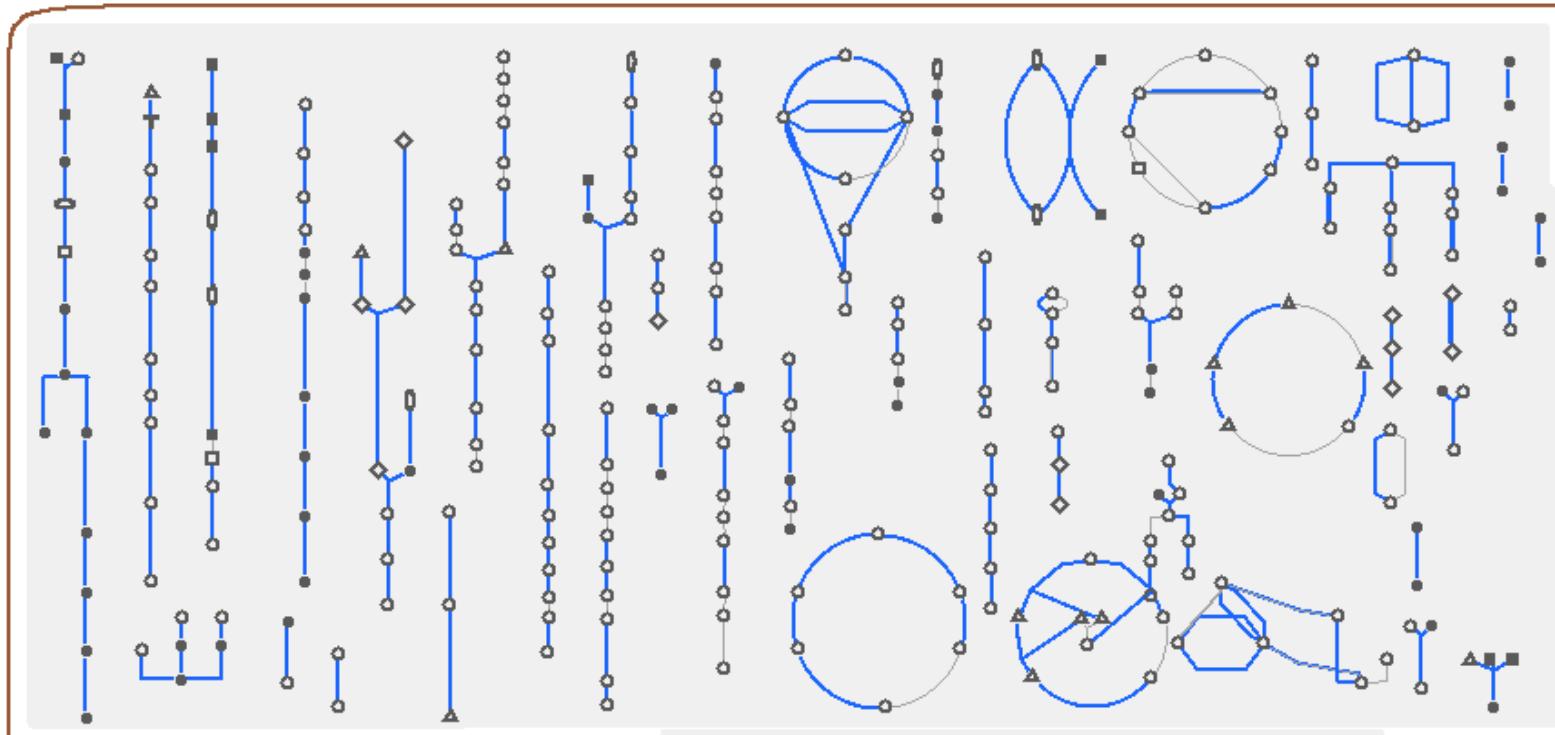
Tomato Potato Pepper *N. benthamiana* Petunia Eggplant SOL Meeting

SolGenomics

url: solgenomics.net

Tools

Biochemical pathways



SolGenomics

url: solgenomics.net

Tools

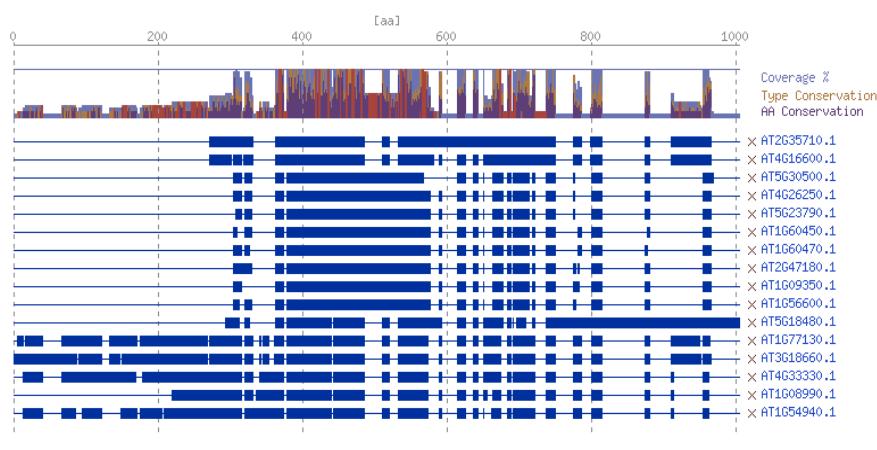
CAPS experiment designer

VIGs tool

Alignment Analyser

Tree Browser

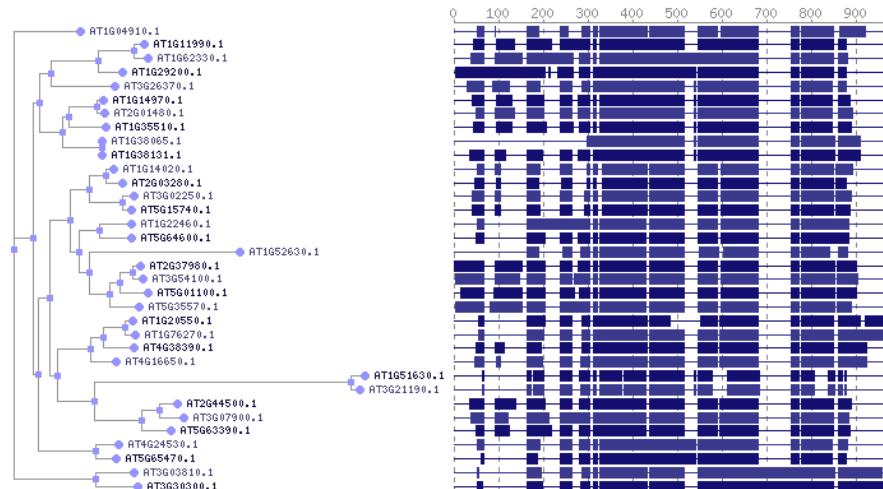
Alignment Image



[Hide Domain Information](#)

Click on a region of the image to zoom in.

You must submit the form below to change alignment parameters or remove members from analysis.





Gold Standard DB for Arabidopsis

Data

Genes - Transcripts – Proteins - Genomes

GO / PO

Locus history

Markers

Array elements, expression, labs

Protocols, Publications

Seeds, Germplasm

Tools

Gbrowse

Synteny browser

Protein-Protein Interaction Browser

AraCyc Metabolic Pathways

BLAST & WU-BLAST

Pattern/Motif Finders

Expression / Array Clustering

Portals

Clone stocks

Education

Gene Expression

Gene Annotation

Metabolome

Mutant (Genetics)

AraPort

url: araport.org



Next-Gen Portal for Arabidopsis

Data

TAIR 10 Genes - Transcripts – Proteins - Genomes

Araport 11

Live community annotation

Science Apps Catalog

This is the current set of public apps deployed at Araport. Once you are logged in, you can preview their functions one-by-one by clicking each icon, or you can install and configure them permanently into [your own private workspace](#).

Tools

JBrowse

App Store

20+ apps

sequence tree views

search

phosphorylation

interactions

BLAST+

50 Years of
Arabidopsis
Research
v0.0.9

[Preview This App](#)

Araport File Browser
v0.2.3

[Preview This App](#)

AT Expression
Profiling
v1.0.3

[Preview This App](#)

Atted Science
Application
v2.1.4

[Preview This App](#)

BAR Expression
Viewer
v1.0.7

[Preview This App](#)

BAR Gene Slider
v0.0.28

The TAIR10 Col-0 Arabidopsis thaliana sequence is decorated with JASPAR transcription factor binding

[Preview This App](#)

BAR Interactions
Viewer
v0.1.8

[Preview This App](#)

BLAST+
v0.1.8

[Preview This App](#)

Genome Assembly and Annotation

Best Practices, Best Tools

Structure

AIM:

Help you develop a strategy for performing a genome assembly and assessing its quality

Theory

Sequence Types

Algorithm Basics

Assessing an Assembly

Tools

For Assembly

For Assessment

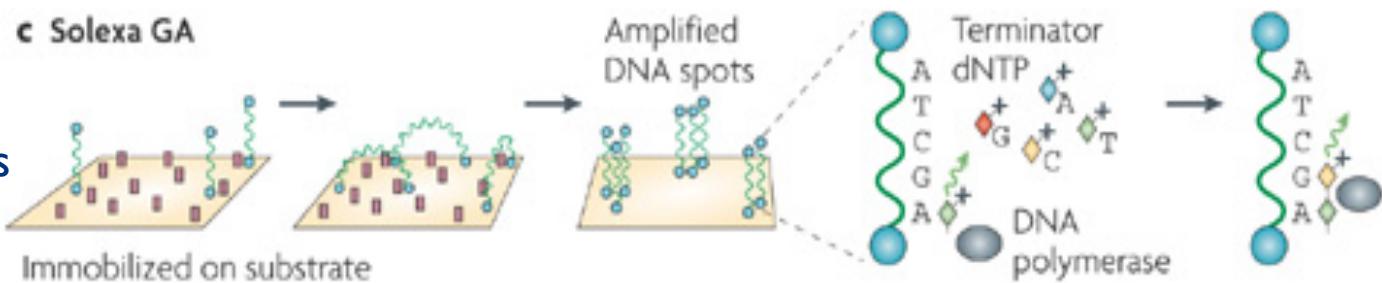
For Annotation

Sequencing

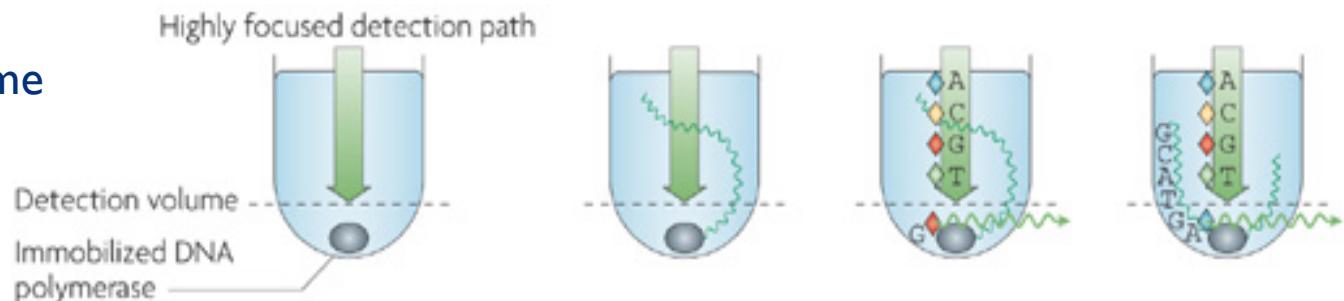
Two major kinds

- Illumina
- Pacific BioSciences

Illumina
Sequence by Synthesis



PacBio
Single Molecule Real Time



Paired-end and Mate Paired Reads

Paired End

Read both ends of the molecule

Separated by ~ 500 bp

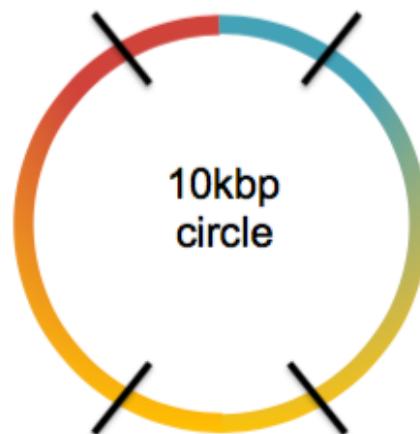


Mate Pair

Circularize long molecules (1-10kbp)

Shear into fragments then sequence

10kbp



2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)



Brief sequence comparison

	Illumina	PacBio
Length	<300bp	<10kb
Reads per run	25 – 400 million	~50 thousand
accuracy	>99%	~86%

The Assembly Problem

Rebuilding a shredded text

Rick Astley accidentally shreds the lyrics to his first album...



Never gonna give you up, never gonna let you down, Never gonna run around and desert you, Never gonna make you cry, never gonna say goodbye, Never gonna tell a lie and hurt you

Never g~~onna~~ g~~onna~~ give you up, never gonna ~~et~~ ~~you~~ down, Never gonna run

Never gonna ~~g~~ive you up, never gonna ~~l~~et you ~~d~~own,, Ne~~w~~er gonna

How can he put them back together?

The fragments from all copies are mixed together

Some are identical

Greedy Algorithm

Never gonna

Never gonna give

gonna give you up, never

never gonna

Never gonna give

never gonna let

Repeated sections make the proper assembly ambiguous

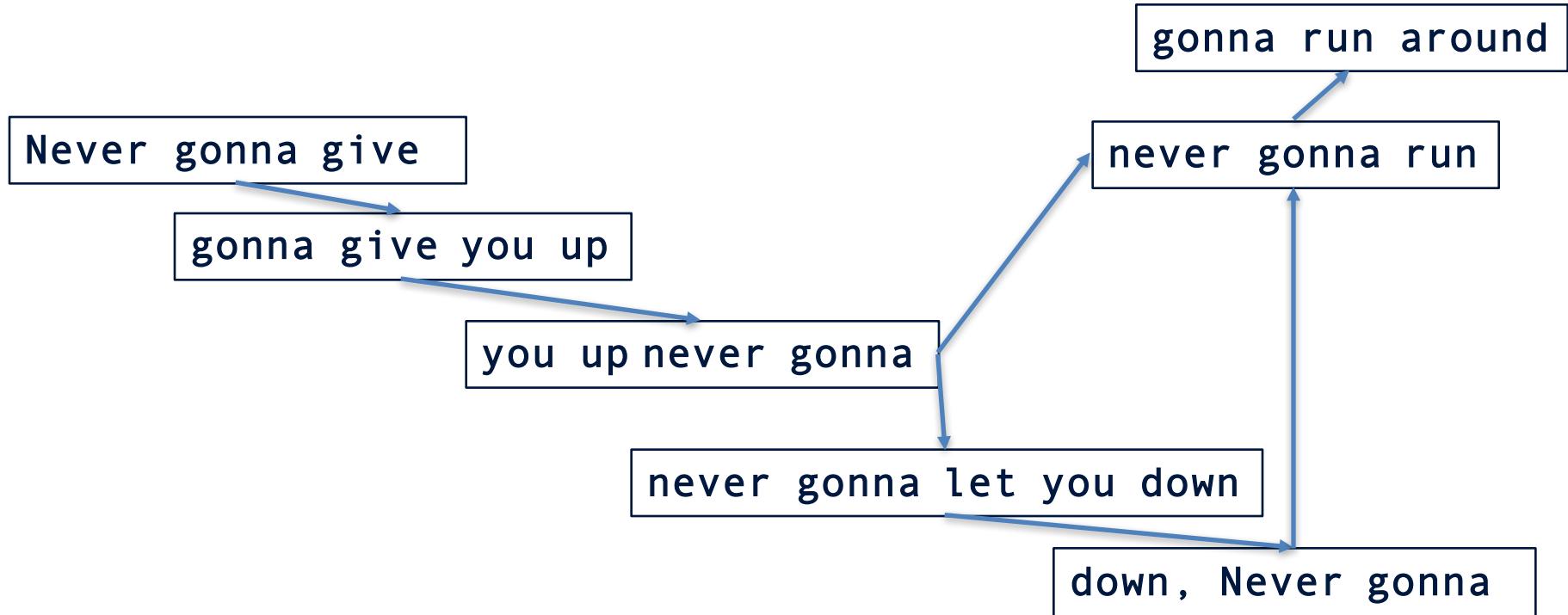
de Bruijn Graph

Make all 'k' length fragments

Network of possible directional links – from overlaps

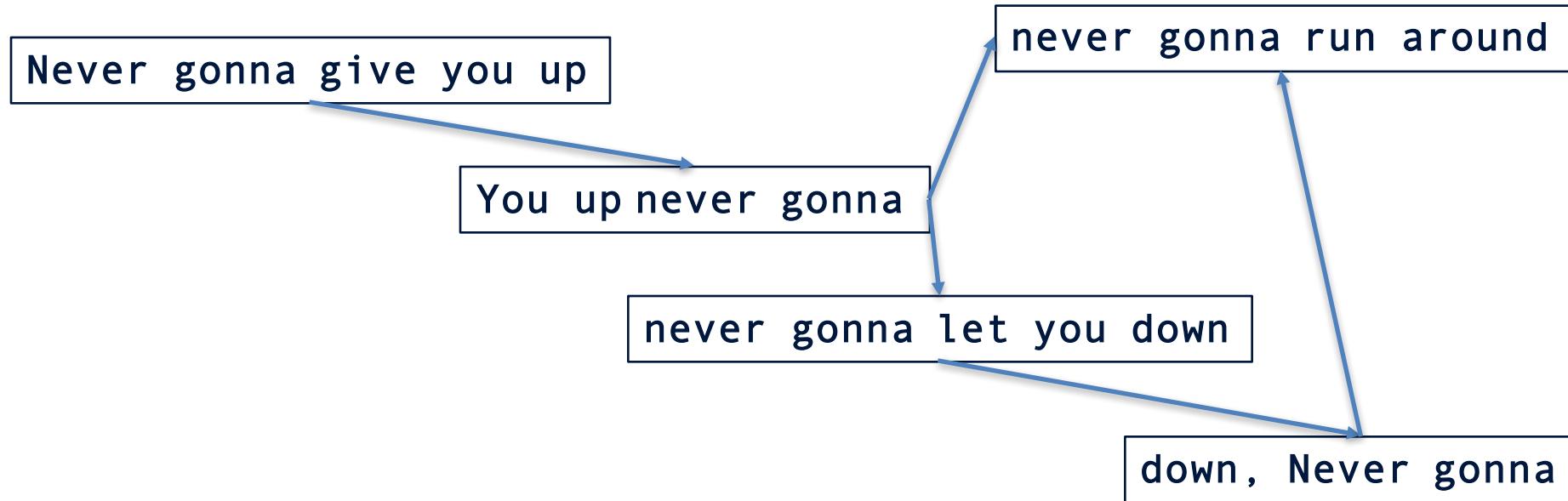


de Bruijn Graph Assembly



Build Graph and
simplify

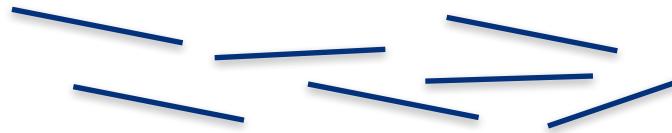
de Bruijn Graph Assembly



Build graph and
simplify

Assembling a genome

I. Shear and sequence DNA



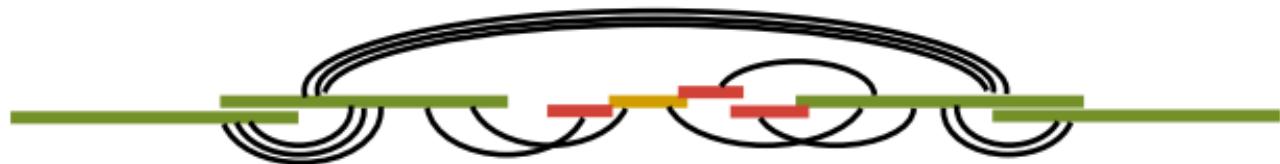
2. Make assembly graph from overlapped reads

AGCCTTAGACCTACAGGATGCGCGCACAGT
GATGCGCGCACAGTAGCCTTAGACT

3. Simplify assembly graph



4. Detangle graph with long reads, mate pairs and other links



Why are genomes hard to assemble?

I. Biological Problems

Polyplody, heterozygosity

Repeats, repeats

2. Sequencing problems

Large genomes genome >>> reads
sequence read errors

3. Computational

Very large genomes

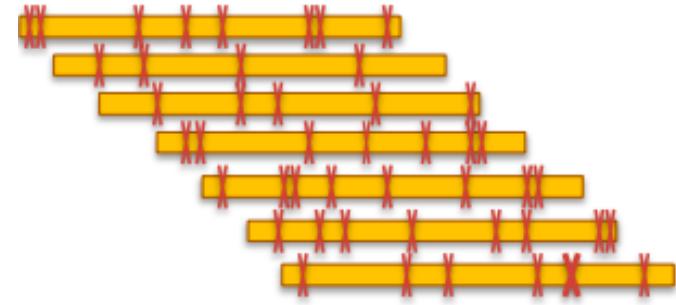
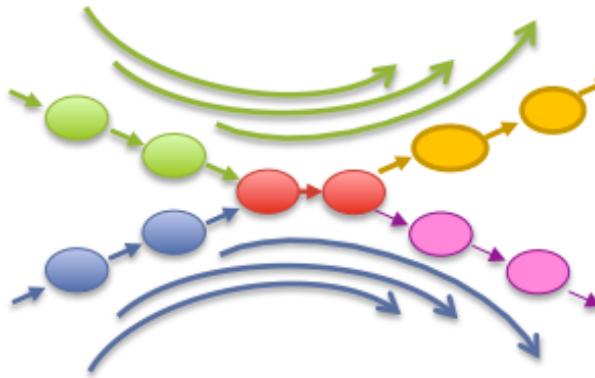
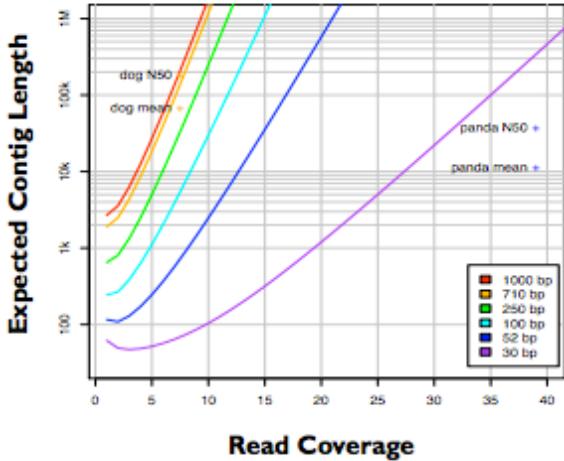
Complex repeat structure

4. Accuracy

What is correct assembly?

How do we tell?

What you need for a good assembly



Coverage

- Oversample

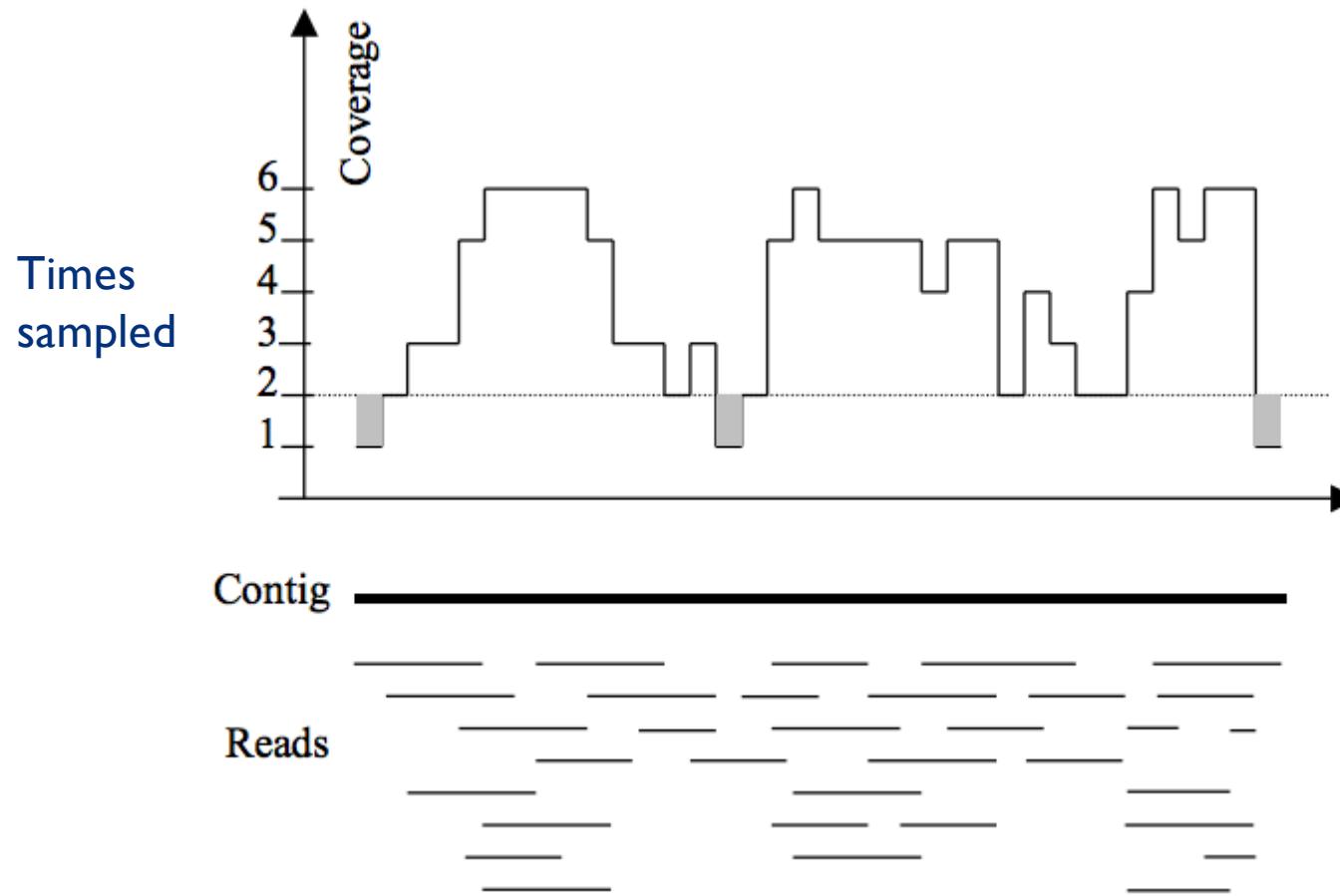
Read Length

- Longer than repeats

Quality

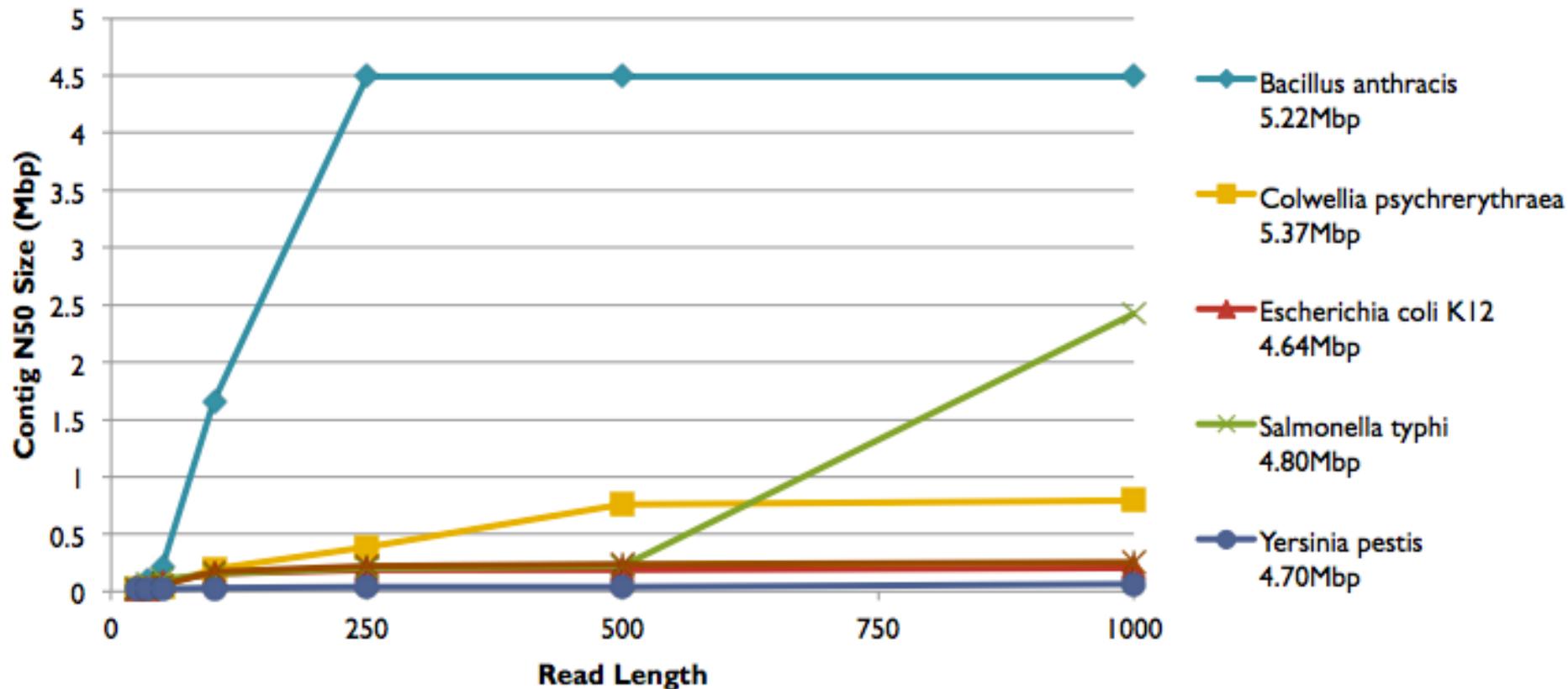
- Errors mess up overlaps

Coverage



Some segments may not even be sequenced!

Read Length and Repeats



Assembly length increases with read length,
especially when it passes the repeat length

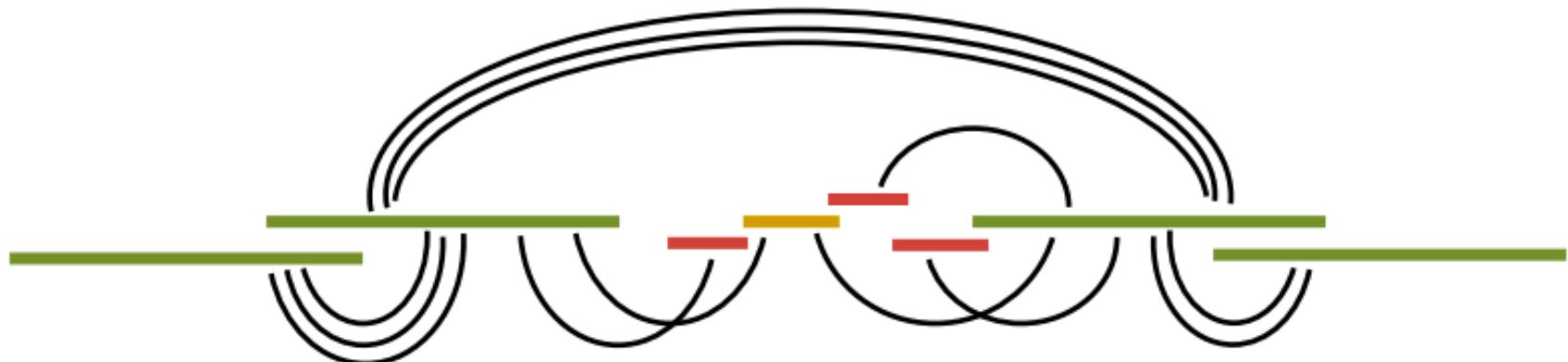
Scaffolding

Contigs end at

- coverage gaps
- errors
- repeat boundaries

Use paired end/mate information to place contigs relative to each other

- has gaps



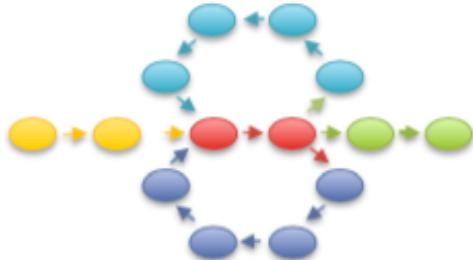
Assembly Tools

Loads! Including but not limited to ..

SOAPdenovo, HyDA, ABySS, ALLPATHS-LG,
Velvet, Celera, PacBio Corrected Reads, Mira,
meraculous, Newbler, Ray, SGA, CLC Assembly
Cell ...

Assembly Tools

SOAPdenovo



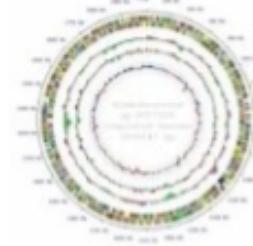
BGI's assembler
(Li et al. 2010)

De bruijn graph
Short reads

Most flexible, but requires a
lot of tuning

[http://soap.genomics.org.cn/
soapdenovo.html](http://soap.genomics.org.cn/soapdenovo.html)

Celera Assembler



JCVI's assembler
(Miller et al. 2008)

Overlap graph
Medium + Long reads

Supports Illumina/454/PacBio
Hybrid assemblies

<http://wgs-assembler.sf.net>

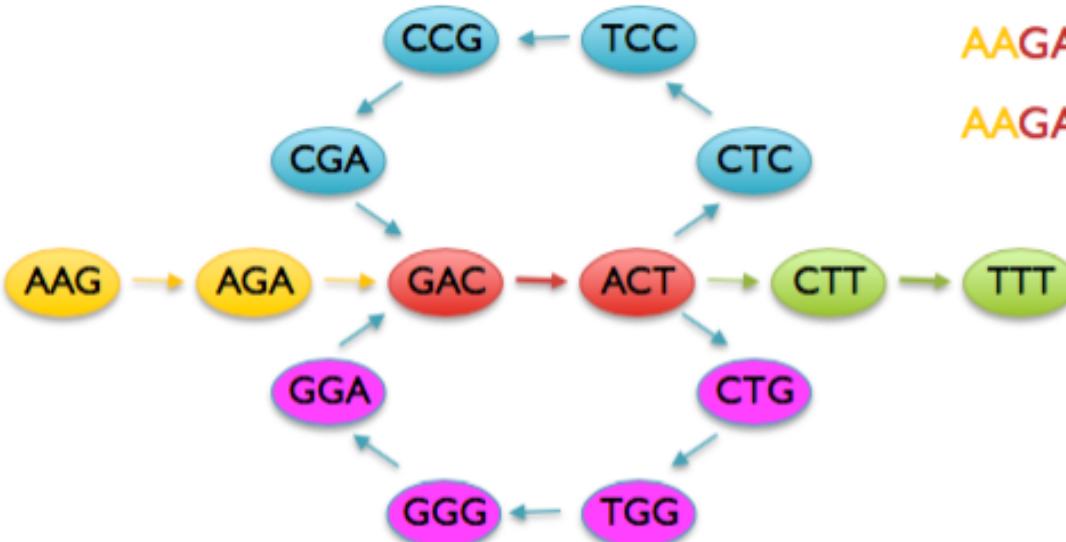
Short read assembly

SOAPdenovo

Reads

AAGA
ACTT
ACTC
ACTG
AGAG
CCGA
CGAC
CTCC
CTGG
CTTT
...

de Bruijn Graph



Potential Genomes

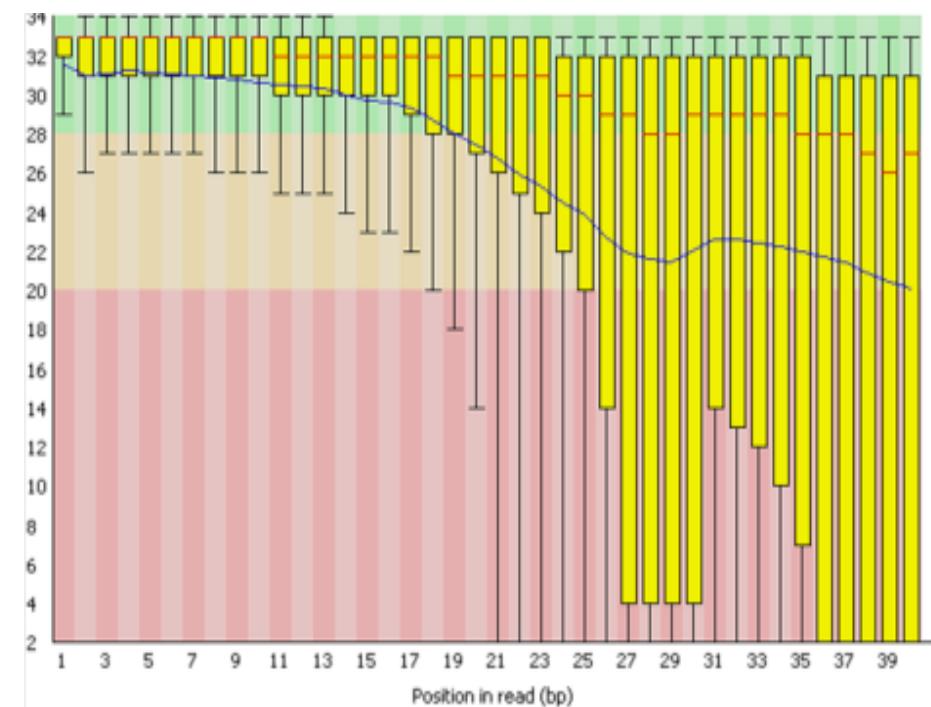
AAGACTCCGACTGGGACTTT
AAGACTGGGACTCCGACTTT

I. Get a big computer

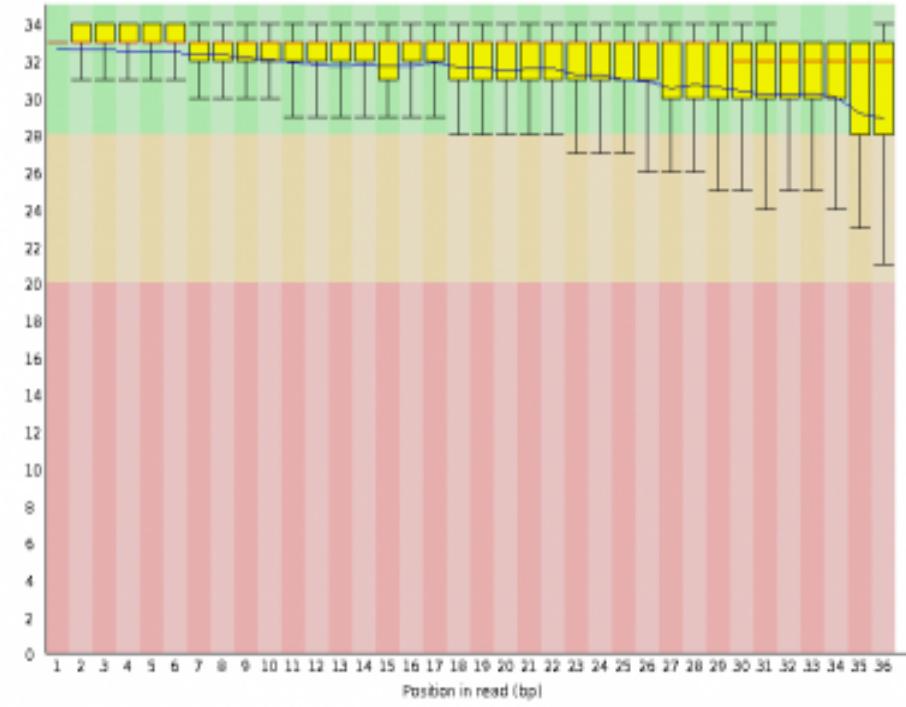
- I. Different specs for every genome and sequence data.
2. > 128Gb RAM for Pathogen genomes
3. > 1Tb for Plant genomes

PreProcessing

FASTQC plots quality scores and sequence metrics and allows development of filter thresholds for whole sequence sets



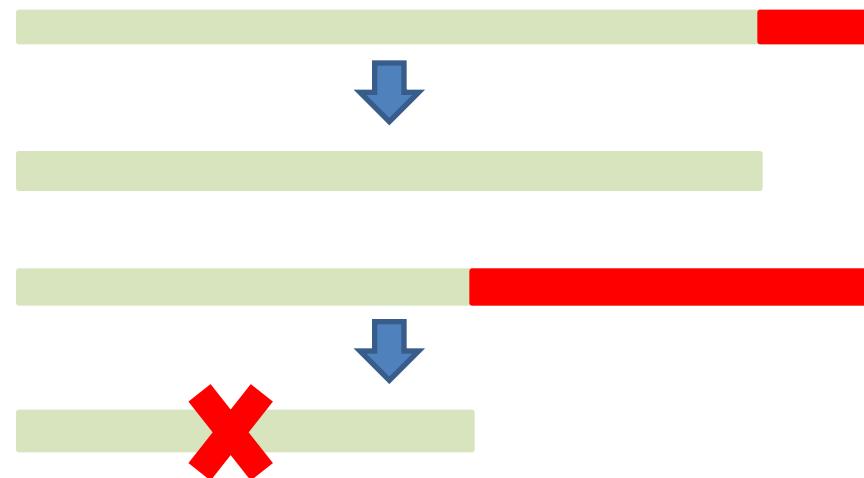
before



after

PreProcessing

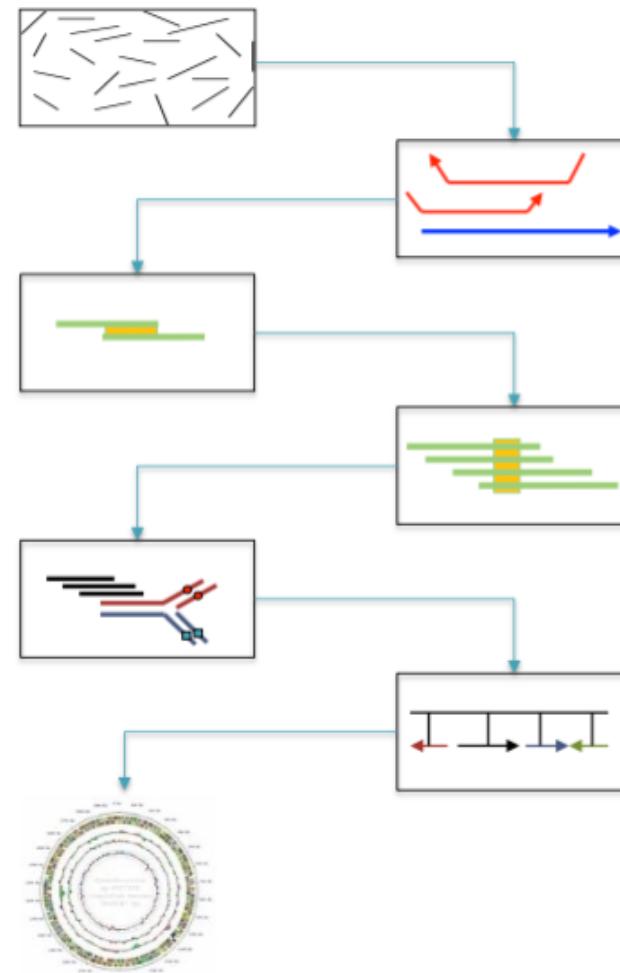
Trimmomatic uses a sliding window approach from the 5` end to identify low quality regions which are then trimmed from the 3` end. Reads < 36 bp are discarded



Long read assembly

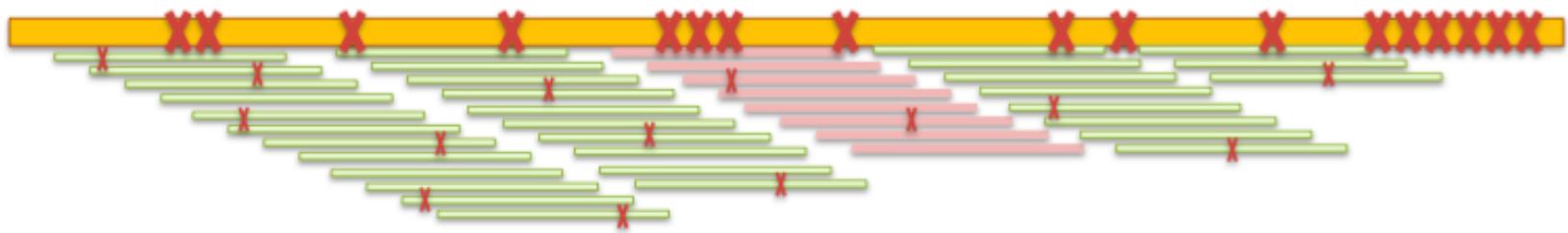
Celera Assembler

1. Pre-overlap
 - Consistency checks
2. Trimming
 - Quality trimming & partial overlaps
3. Compute Overlaps
 - Find high quality overlaps
4. Error Correction
 - Evaluate difference in context of overlapping reads
5. Unitigging
 - Merge consistent reads
6. Scaffolding
 - Bundle mates, Order & Orient
7. Finalize Data
 - Build final consensus sequences



Hybrid assembly

1. Map Illumina reads to PacBio reads to correct errors
2. Assemble



Hybrid error correction and de novo assembly of single-molecule sequencing reads.
Koren, S, Schatz, MC, et al. (2012) *Nature Biotechnology*. doi:10.1038/nbt.2280

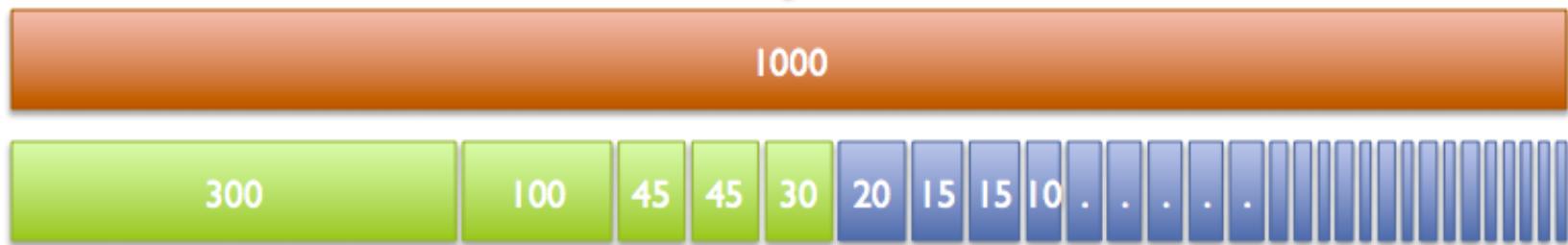
Assessing an assembly

N50

50 % of genome in contigs longer than the N50 contig

Example: 1 Mbp genome

50%
↓



N50 size = 30 kbp

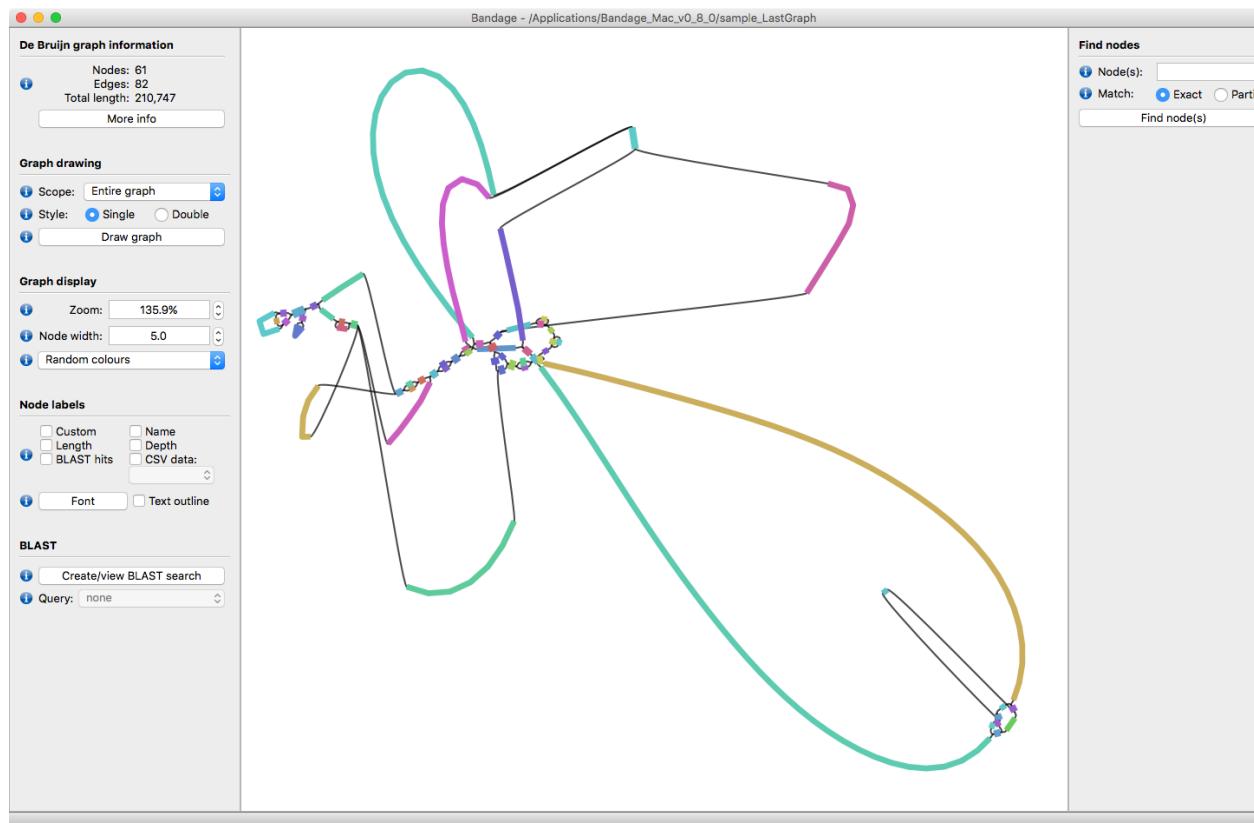
($300k+100k+45k+45k+30k = 520k \geq 500\text{kb}$)

Massively Abused!!

Assessing an assembly

Visual Graph Inspection - Bandage

Allows easy visualisation of some graphs. Assess structure (and sanity) of the assembly



Assessing an assembly

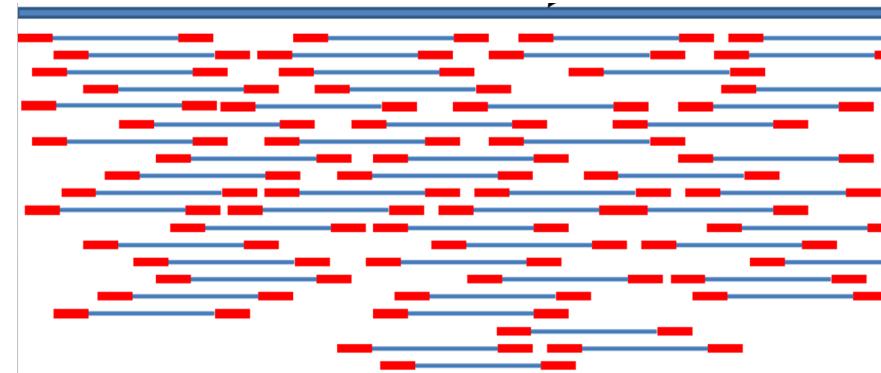
Internal contig consistency

Align original reads to contigs.

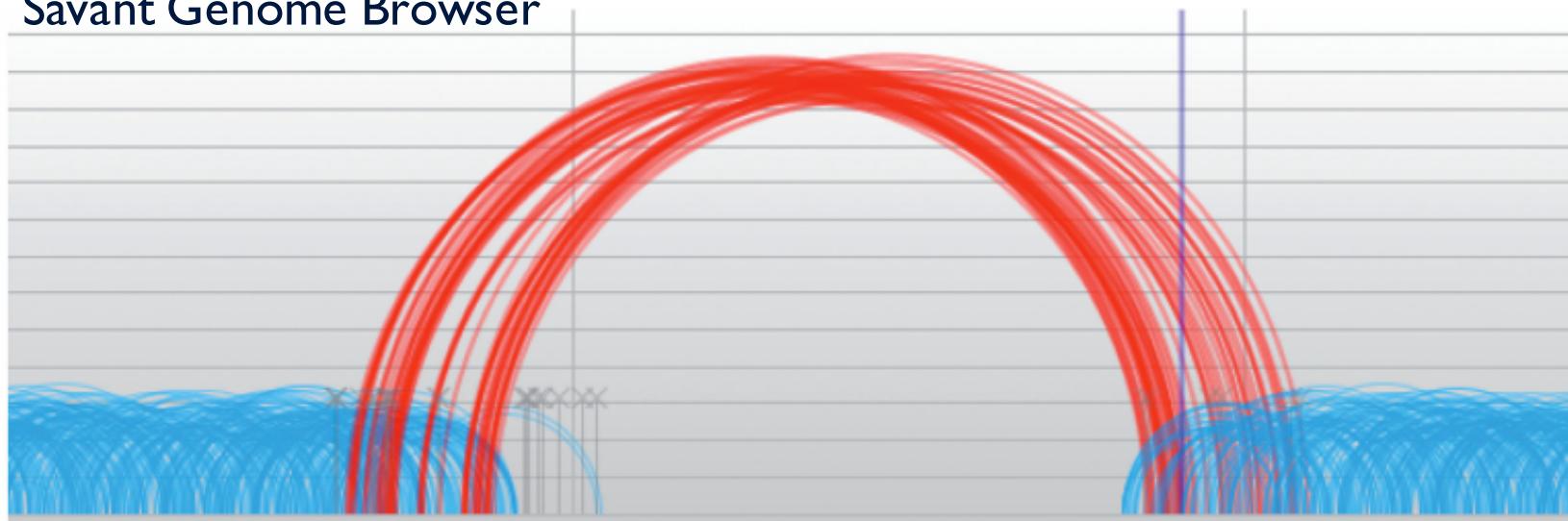
Do they all align without breaks?

Is the distance between pairs consistent?

Is coverage as expected?



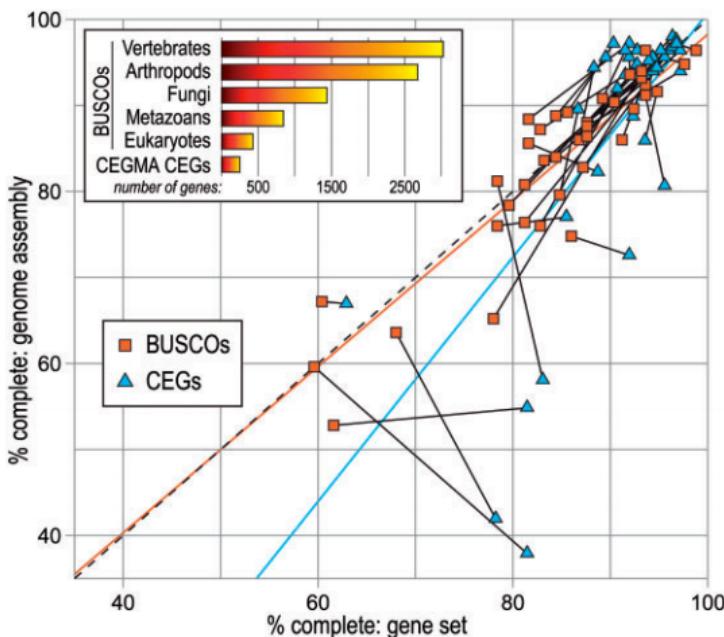
Savant Genome Browser



Assessing an assembly

Gene Content

Core Orthologues – are they all there?



BUSCO – identify and quantify completeness of core orthologues per clade in an assembly

A photograph of a man from the chest up, wearing a dark blue suit jacket, a white shirt, and a dark bow tie. He is standing behind a whiteboard that has a dark blue header bar. The text on the board is overlaid on this header bar.

RNA Seq

Best Practices, Best Tools

Structure

AIM:

Help you develop a strategy for performing RNAseq and assembly and assessing analysis quality

Theory

Sequence Types

Algorithm Basics/ Tool Types

Assessing an Alignment/Assembly

Tools

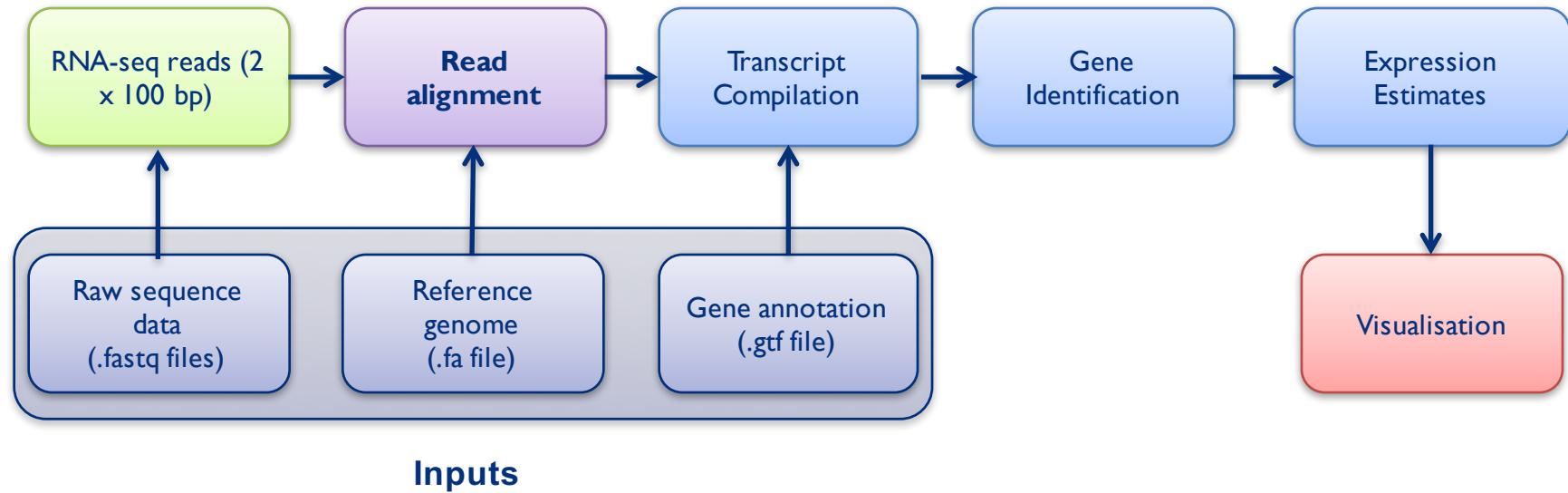
For Alignment/Assembly

For Differential Expression

Why RNA Seq?

- Gene expression and differential expression
- Alternative expression analysis
- Transcript discovery and annotation
- Allele specific expression
- Relating to SNPs or mutations
- Mutation discovery
- Fusion detection
- RNA editing

The RNAseq Pipeline



Alignment with Millions of Sequences

- Computational cost
 - 100's of millions of reads
 - Oldstyle (Smith-Waterman) – Newstyle (Burrows Wheeler Transform, Suffix Trees)
- Introns!
 - Spliced vs. unspliced alignments
- Can I just align my data once using one approach and be done with it?
 - Unfortunately probably not

Alignment (mapping) strategies

De novo assembly

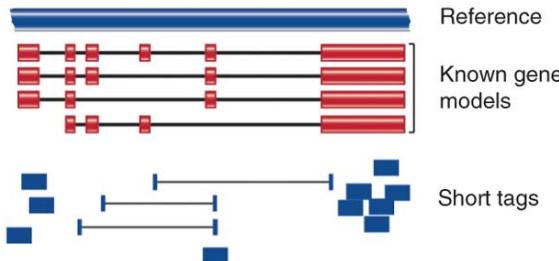


Assemble transcripts from overlapping tags



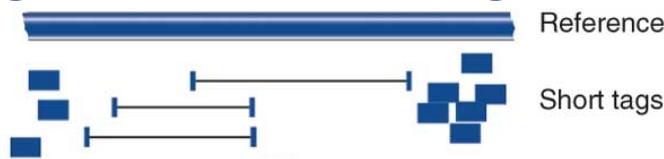
Optional: align to genome to get exon structure

Align to transcriptome



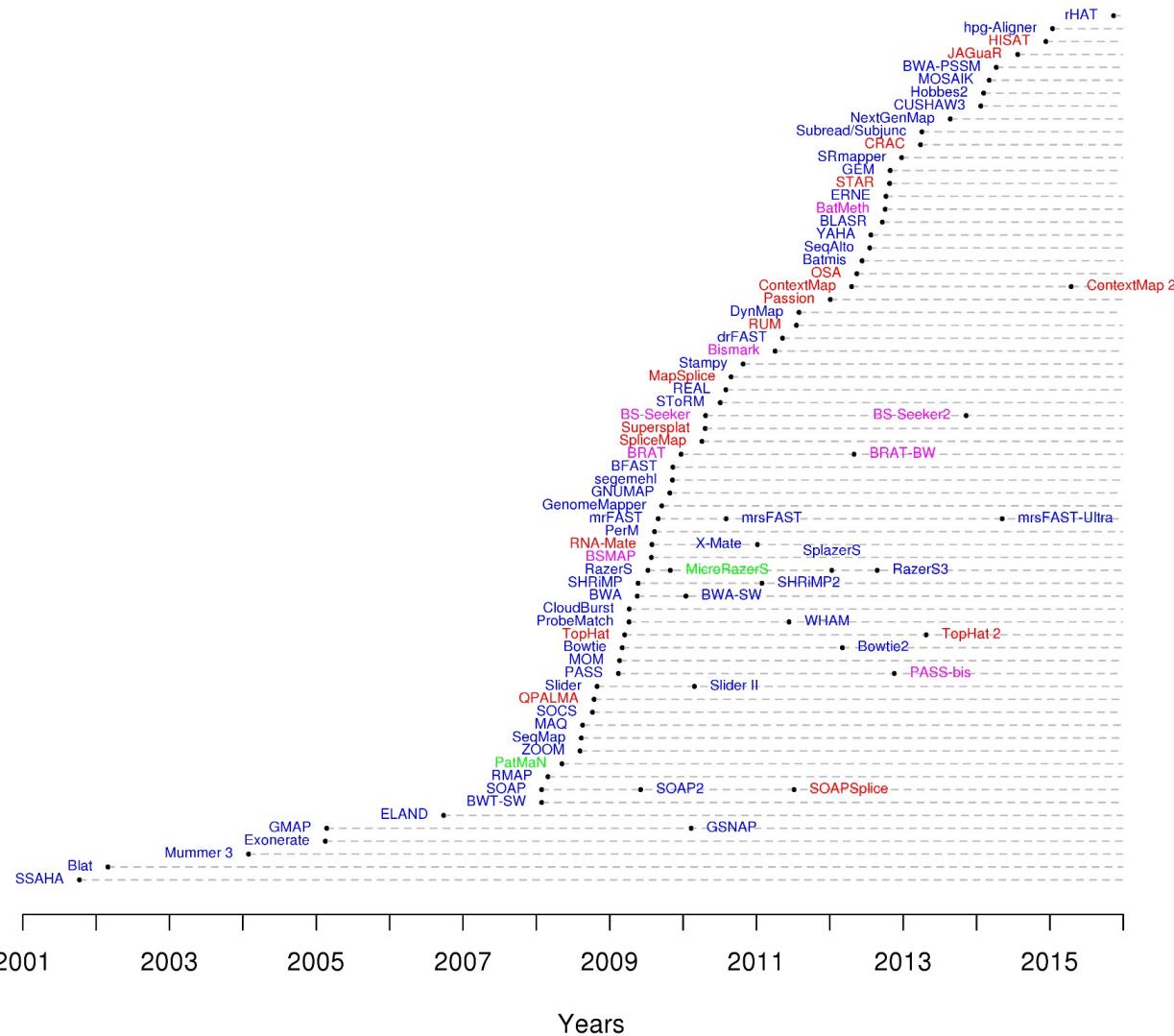
Use known and/or predicted gene models to examine individual features

Align to reference genome



Infer possible transcripts and abundance

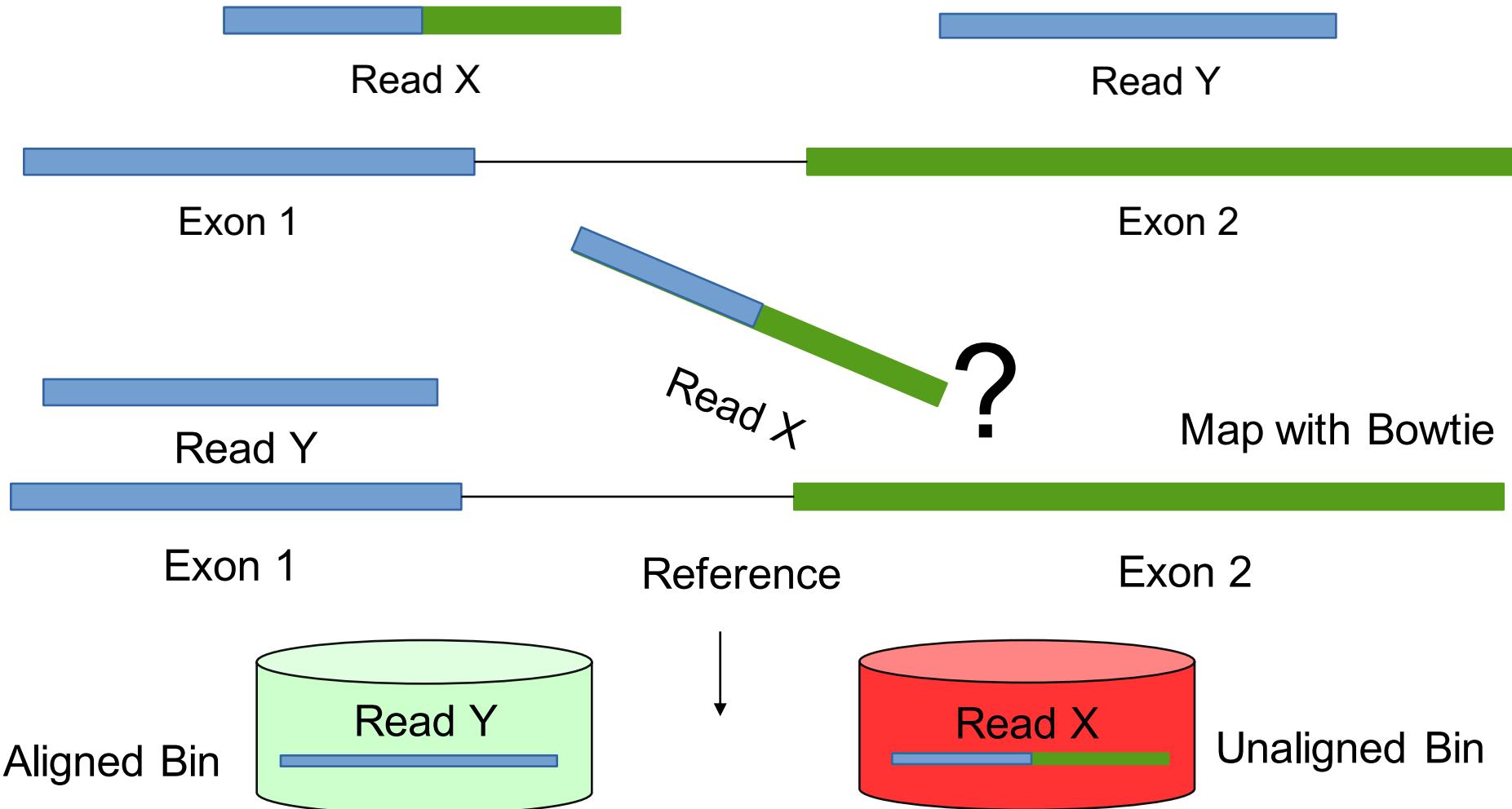
Alignment Tools



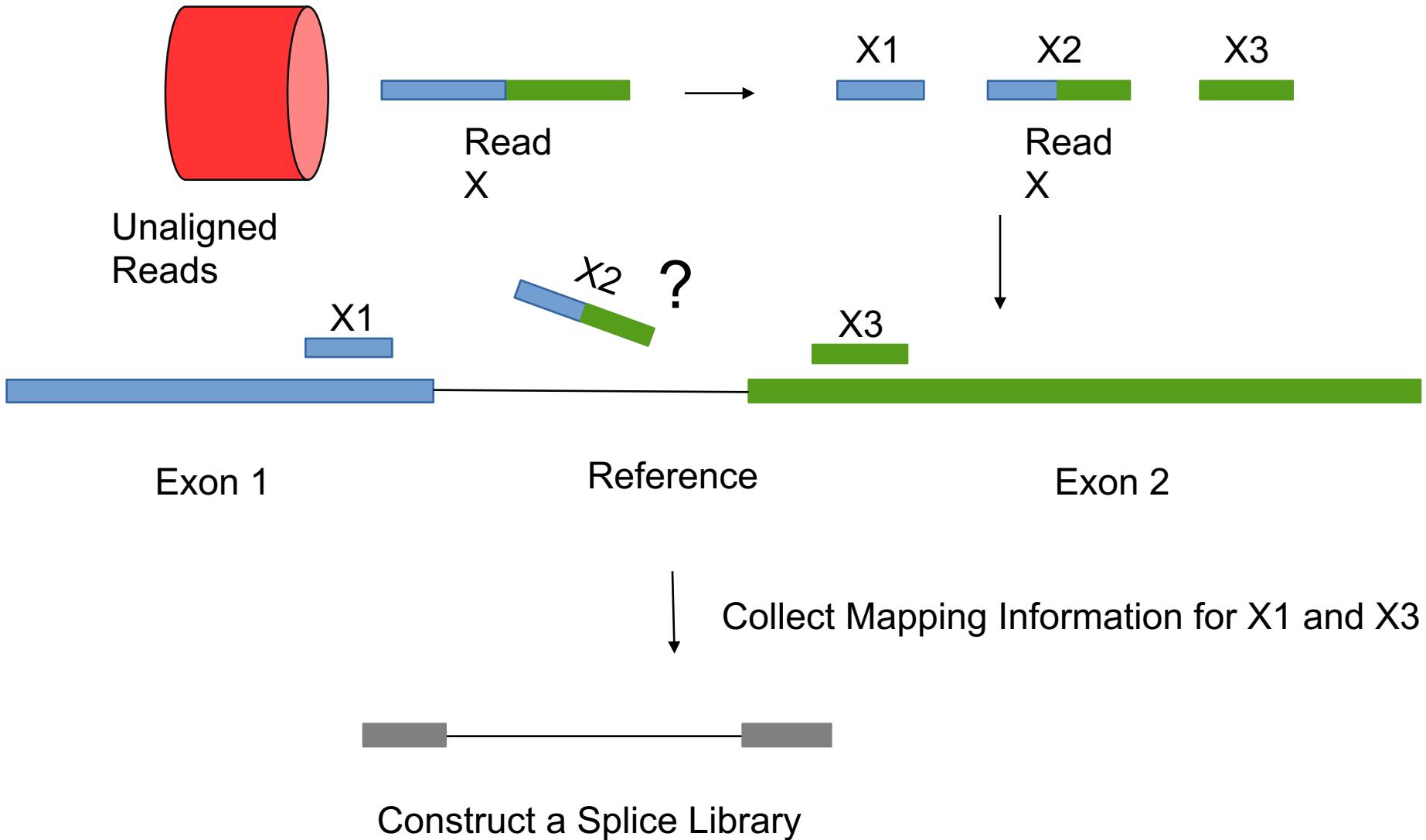
RNA
Bisulfite
DNA
microRNA

Splice Aware Aligners – Tophat/Bowtie

RNA-seq reads may span introns



Splice Aware Aligners - Tophat



SAM/BAM

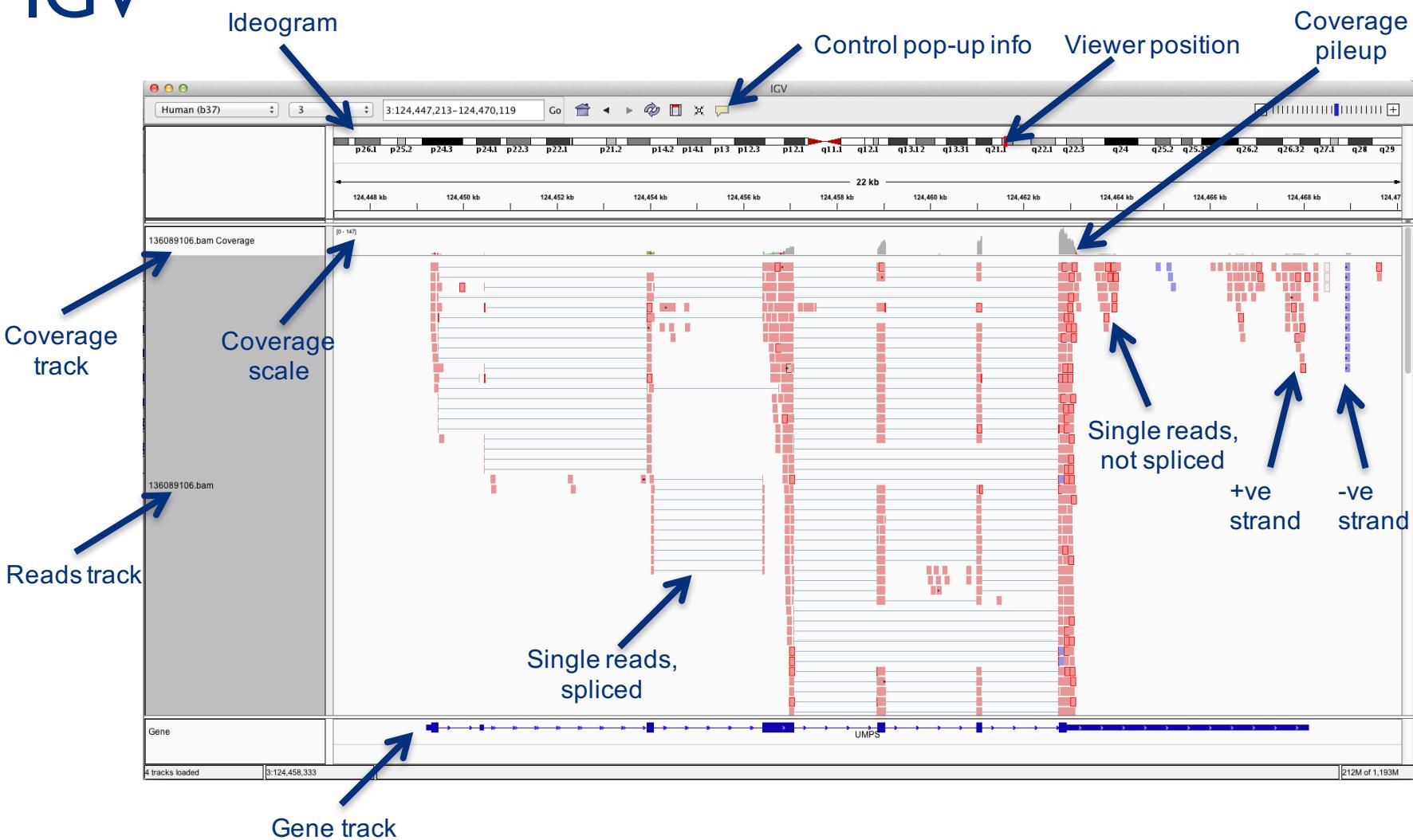
Example SAM/BAM header section (abbreviated)

```
mgriffit@linus270 ~ samtools view -H /gscmnt/gc13001/info/model_data/2891632684/build136494552/alignments/136080019.bam | grep -P "SN\|22|HD|RG|PG"
@HD VN:1.4 SO:coordinate
@SQ SN:22 LN:51034566 UR:ftp://ftp.ncbi.nih.gov/genbank/genomes/Eukaryotes/vertebrates_mammals/Homo_sapiens/GRCh37/special_requests/GRCh37-lite.fa.gz AS:GRCh37-lite M5:a718aca6135fdca8357d5bfe9
4211d SP:Homo sapiens
@RG ID:2888721359 PL:illumina PU:D1BA4ACXX.3 LB:H_KA-452198-0817007-cDNA-3-lib1 PI:365 DS:paired end DT:2012-10-03T19:00:00-0500 SM:H_KA-452198-0817007 CN:WUGSC
@PG ID:2888721359 VN:2.0.8 CL:tophat --library-type fr-secondstrand --bowtie-version=2.1.0
@PG ID:MarkDuplicates PN:MarkDuplicates PP:2888721359 VN:1.85(exported) Cl:net.sf.picard.sam.MarkDuplicates INPUT=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/scratch-Ilg6Y/H_KA-452198-0817007-cDNA-3-lib1-2888360300.bam] OUTPUT=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/scratch-Ilg6Y/H_KA-452198-0817007-cDNA-3-lib1-2888360300-post_dup.bam METRICS_FILE=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/staging-liuJS/H_KA-452198-0817007-cDNA-3-lib1-2888360300.metrics REMOVE_DUPLICATES=false ASSUME_SORTED=true MAX_FILE_HANDLES_FOR_READ_ENDS_MAP=9500 TMP_DIR=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/scratch-Ilg6Y] VALIDATION_STRINGENCY=SILENT MAX_RECORDS_IN_RAM=5000000 PROGRAM_RECORD_ID=MarkDuplicates PROGRAM_GROUP=Name=MarkDuplicates MAX_SEQUENCES_FOR_DISK_READ_ENDS_MAP=50000 SORTING_COLLECTION_SIZE_RATIO=0.25 READ_NAME_REGEX=[a-zA-Z0-9]+:[0-9]+:[0-9]+.* OPTICAL_DUPLICATE_PIXEL_DISTANCE=100 VERBOSITY=INFO QUIET=false COMPRESSION_LEVEL=5 CREATE_INDEX=false CREATE_MDS_FILE=false
mgriffit@linus270 ~
```

Example SAM/BAM alignment section (only 10 alignments shown)

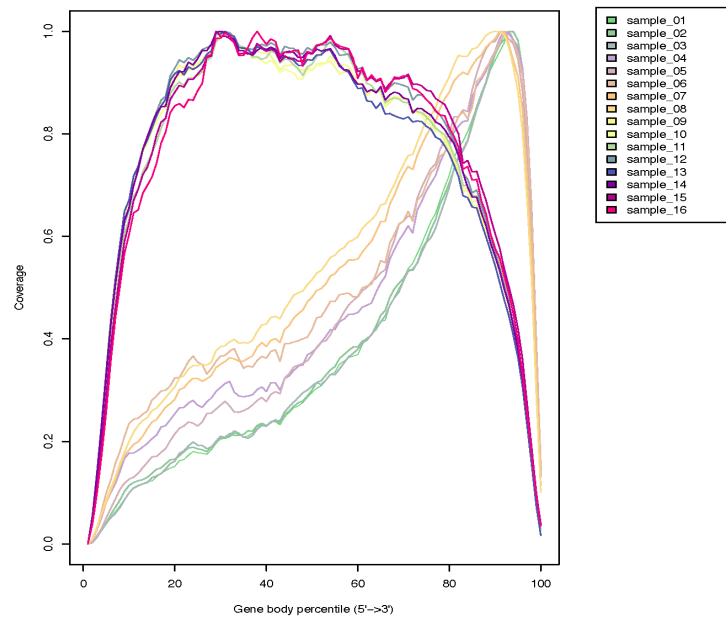
Visualizing Alignments

IGV

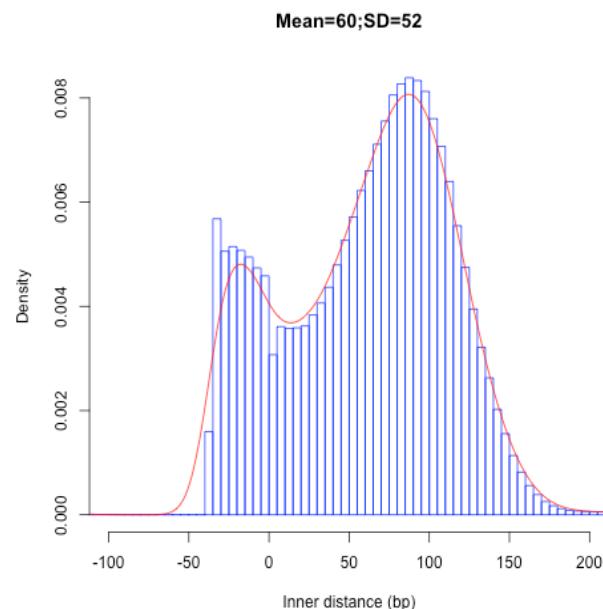


Alignment QC

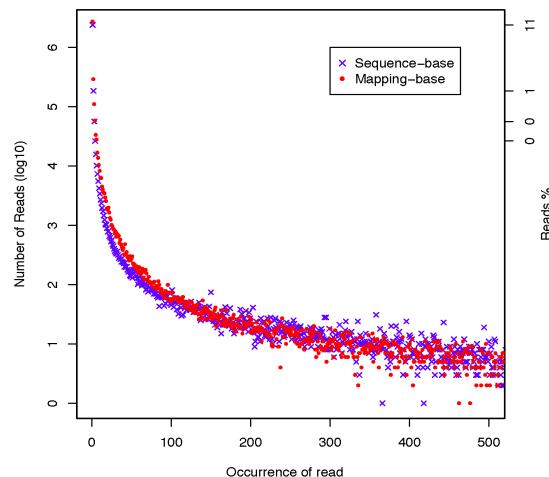
Transcript coverage 3' Bias



Insert Size

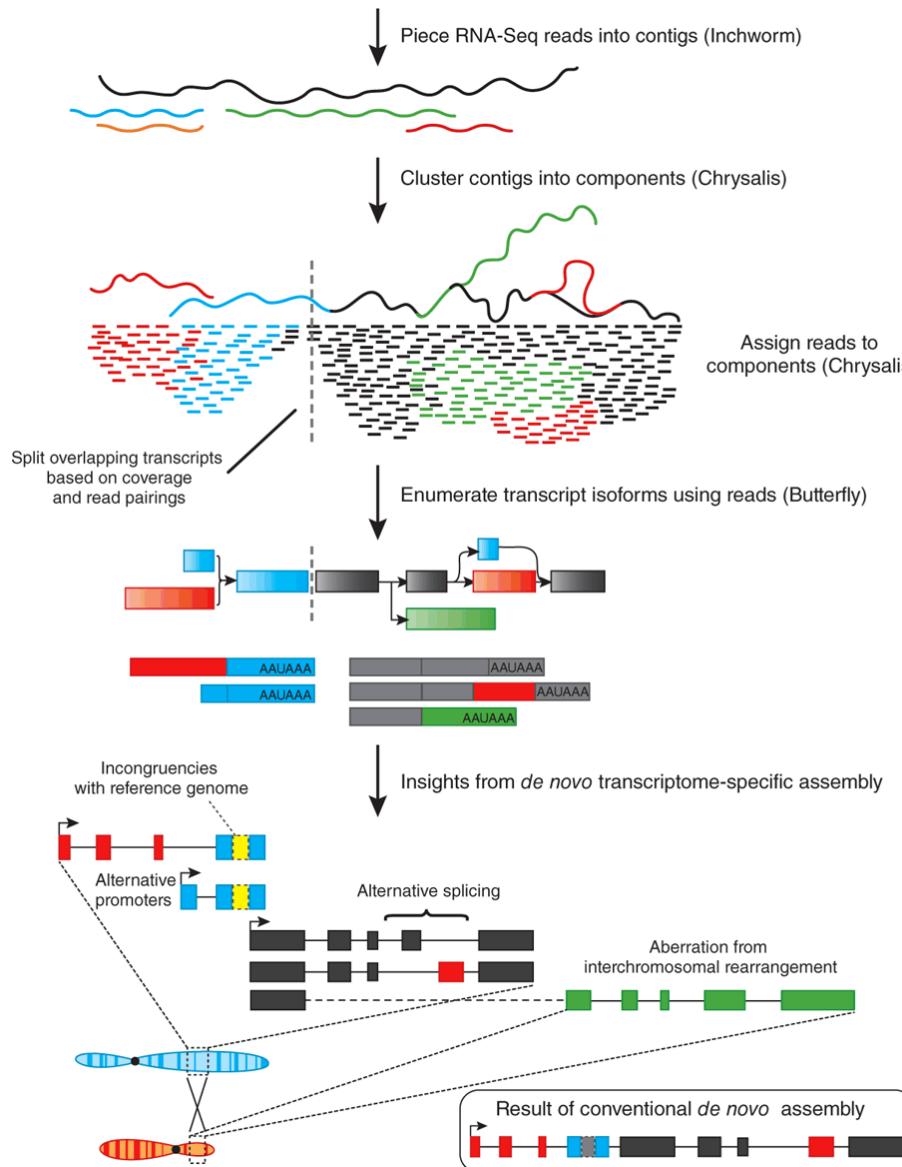


PCR Duplicates



<http://rseqc.sourceforge.net>

Complete de novo – no reference

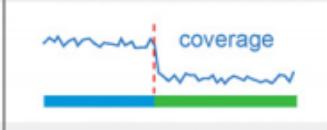
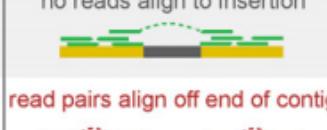
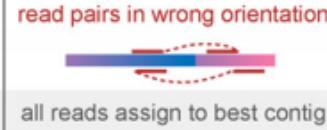


Trinity

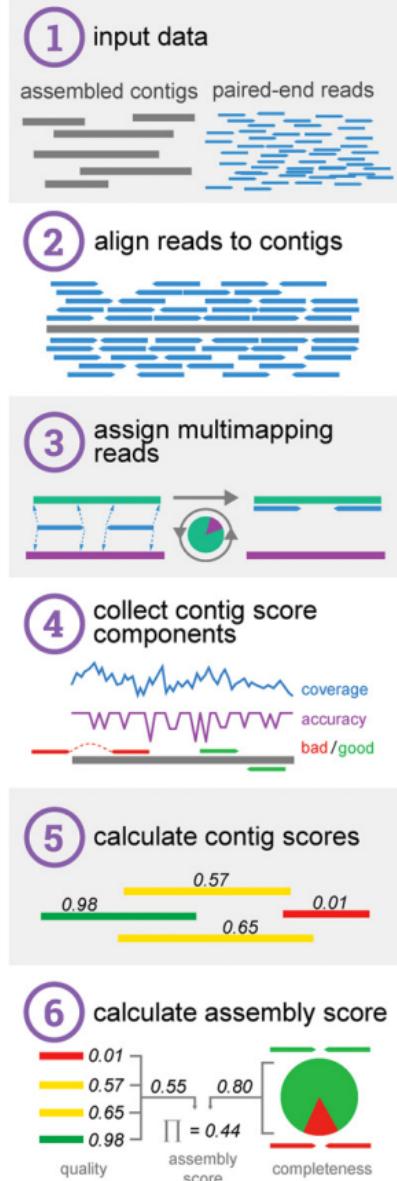
<https://github.com/trinityrnaseq>

Assembly QC

TransRate

Error type	Transcripts	Assembly	Read evidence
Family collapse	geneAA geneAB geneAC n=3	n=1	
Chimerism	geneC geneB n=2	n=1	
Unsupported insertion	n=1	n=1	
Incompleteness	n=1	n=1	
Fragmentation	n=1	n=4	
Local misassembly	n=1	n=1	
Redundancy	n=1	n=3	

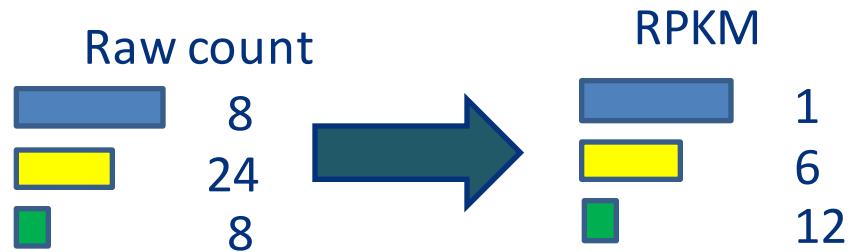
 TransRate



Alignments Mappings to expression

Normalize counts for expression levels –

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



$$\text{RPKM}_{\text{geneA}} = 10^9 \frac{\text{C}_{\text{geneA}}}{\text{NL}}$$

C_{geneA} = number of reads mapped to geneA

N = total number of reads

L = length of transcript in units of Kb

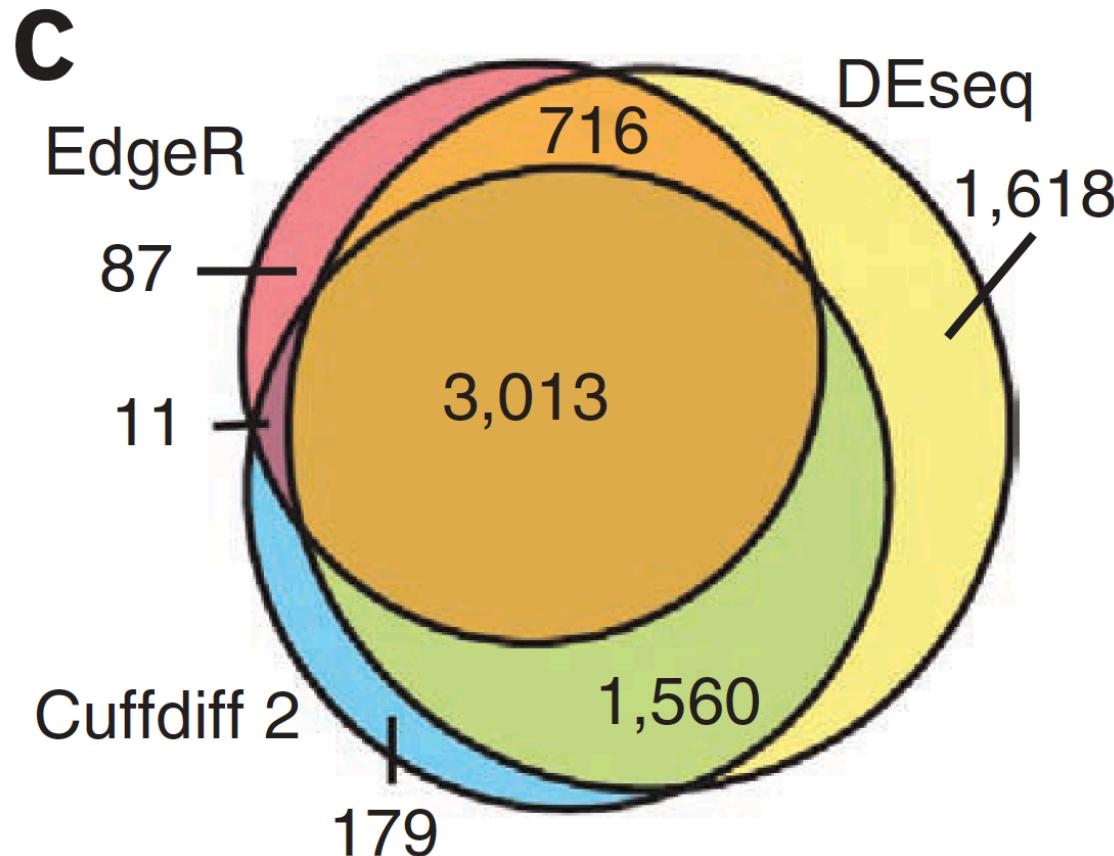
Cufflinks/CuffDiff
Trinity

Alternative differential expression methods

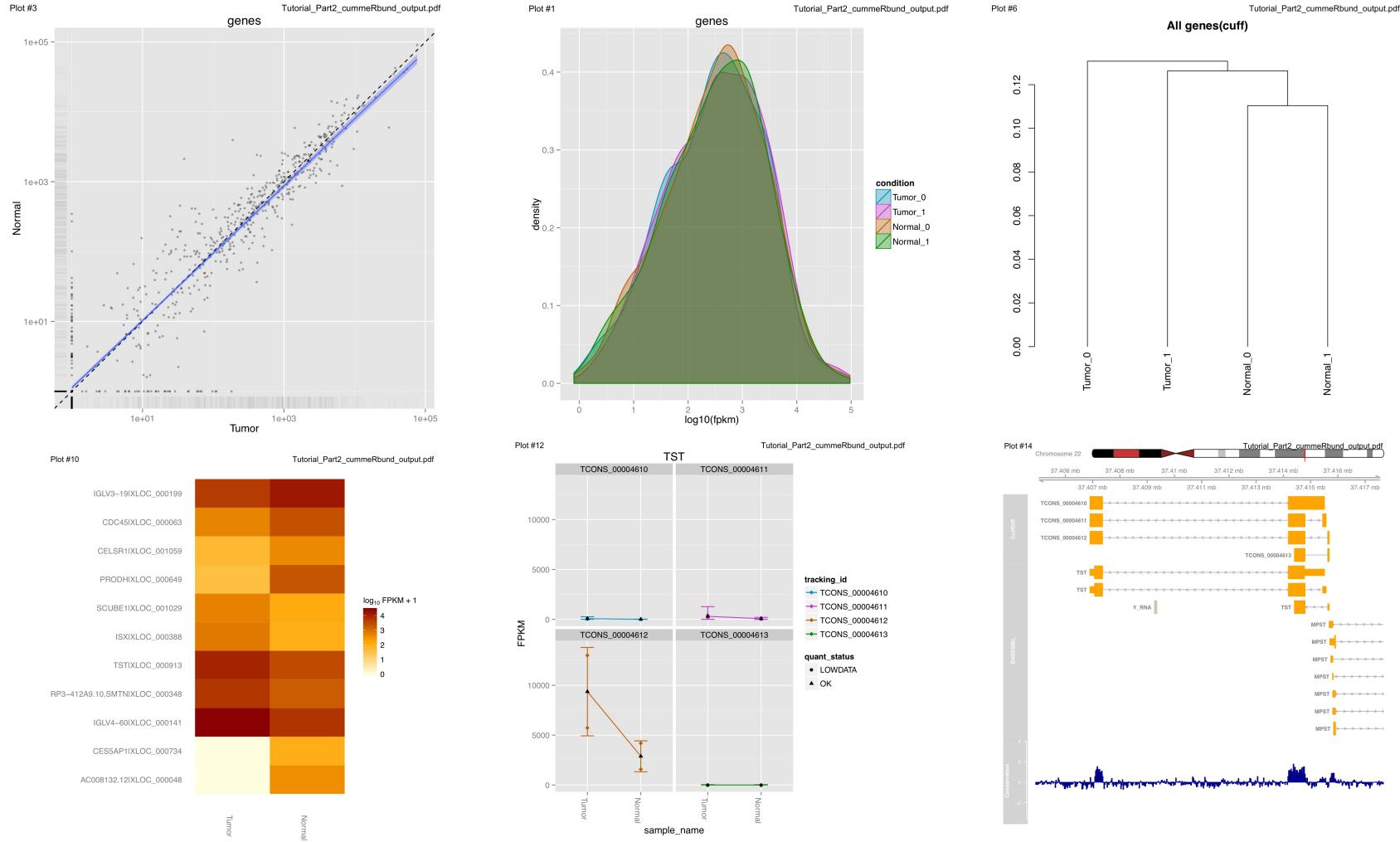
Raw count approaches

DESeq - <http://www-huber.embl.de/users/anders/DESeq/>

edgeR - <http://www.bioconductor.org/packages/release/bioc/html/edgeR.html>



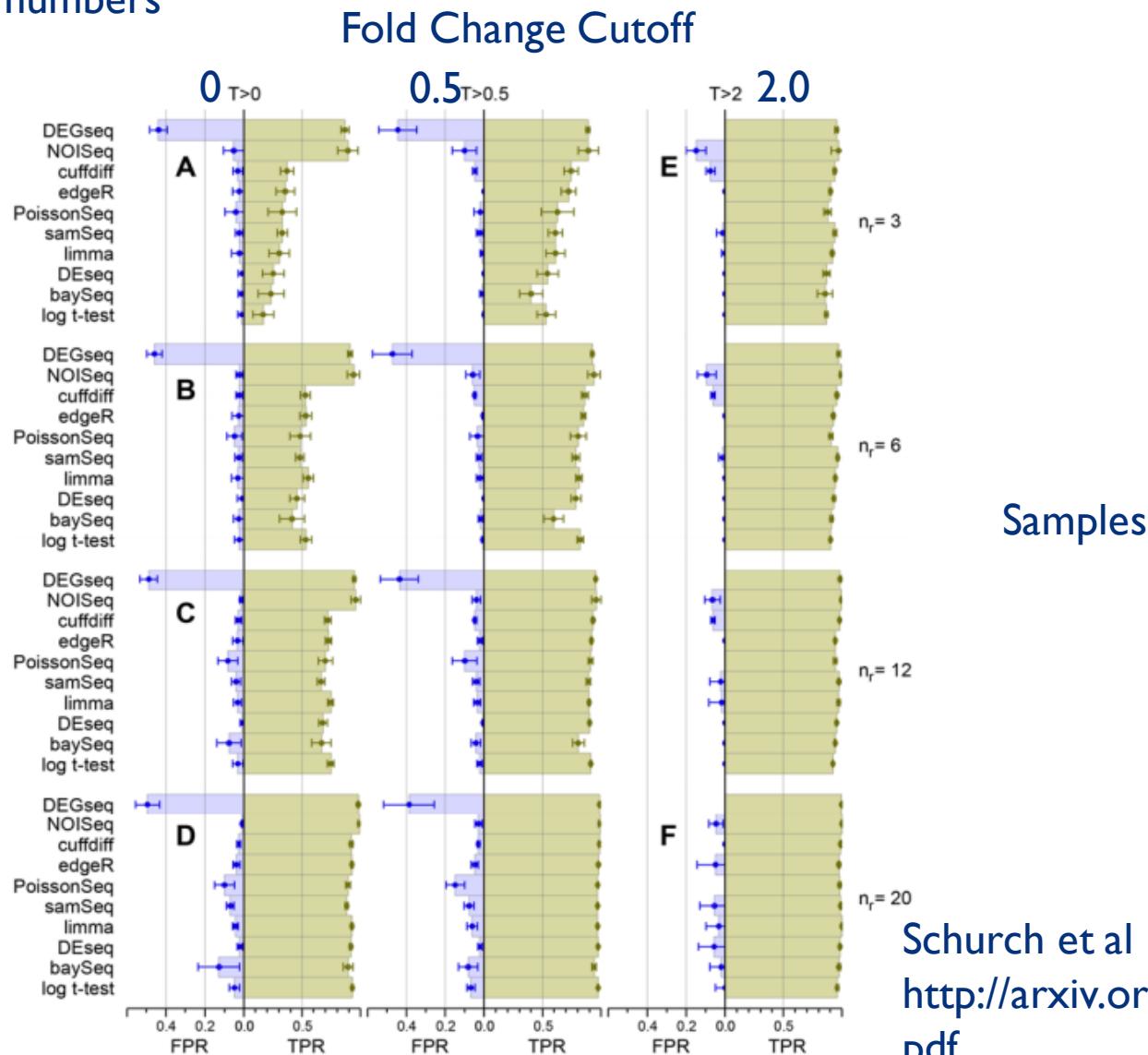
Analysis of expression in R



Which tool? Which experimental design?

False/True Positive Rate at different fold changes
and replicate numbers

Tool



Which tool? Which experimental design?

- At least six replicates per condition for all experiments.
- At least 12 replicates per condition for experiments where identifying the majority of all DE genes is important.
- For experiments with <12 replicates per condition; use edgeR(exact).
- For experiments with >12 replicates per condition; use DESeq.
- Apply a fold-change threshold appropriate to the number of replicates per condition between $0.1 \leq T \leq 0.5$.

Schurch et al

<http://arxiv.org/pdf/1505.02017v2.pdf>



Metagenomics

Best Practices, Best Tools

Structure

AIM:

Help you understand strategies used in metagenomics experiment

Experimental Aims
Pipelines and Tools

Metagenome Taxonomy
Metagenome Function
Metatranscriptomes

What is metagenomics?

Environmental genomics

Samples of mixed genomes, from some place

Microbiomes

e.g. genomes of microbes in soil samples

Plant Pathogen Infection

e.g. field pathogenomics

Not so much marker gene surveys

e.g. 16S

Effective way to profile structure and function of communities

Metagenomics Sequencing



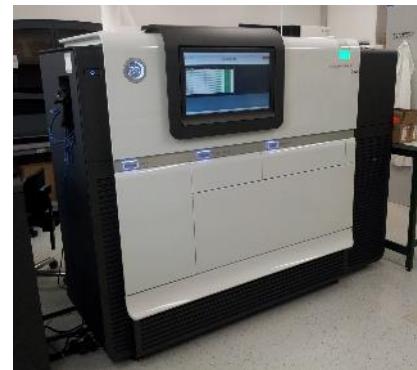
Sanger



Roche 454



Illumina *Seq



Pacific Biosciences



Ion Torrent

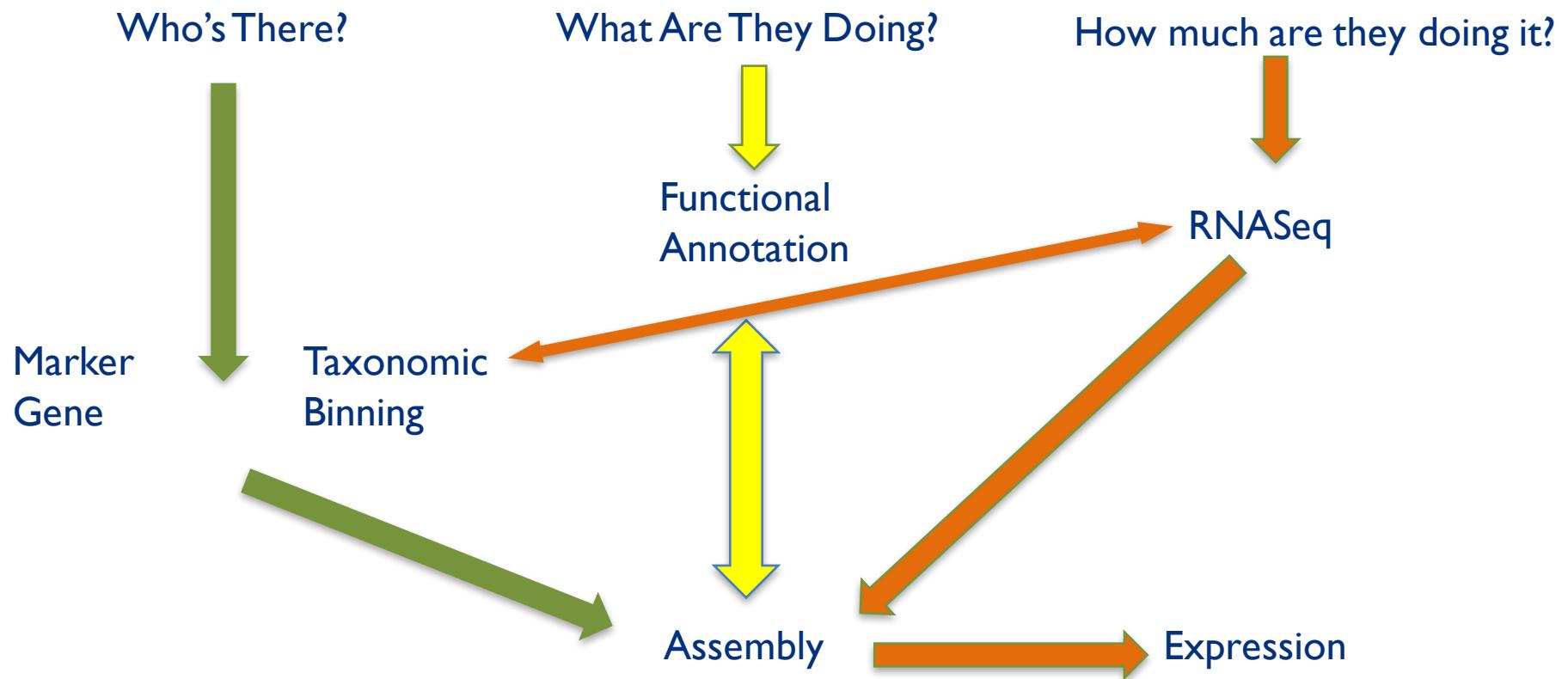


Nanopore

Data Quality Issues

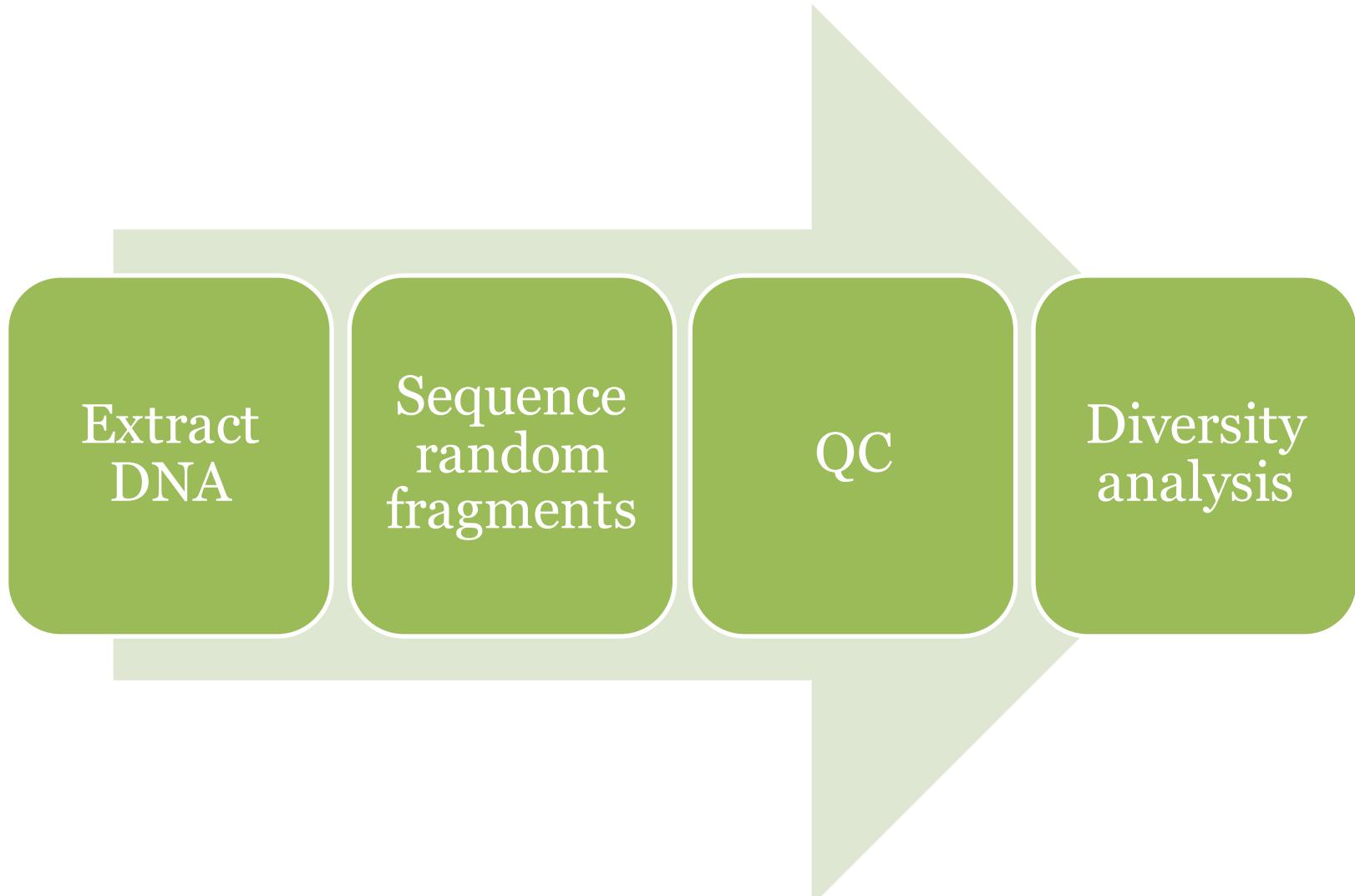
- Sequencing errors
 - Error rates, error type (PacBio: 10% random, Illumina – 0.1% substitution)
- Chimeras
 - Amplification artifacts, cloning of restriction fragments

Metagenomics approaches - Summary



Metagenome Taxonomy

Who's there?



Metagenome Taxonomy

Who's there?

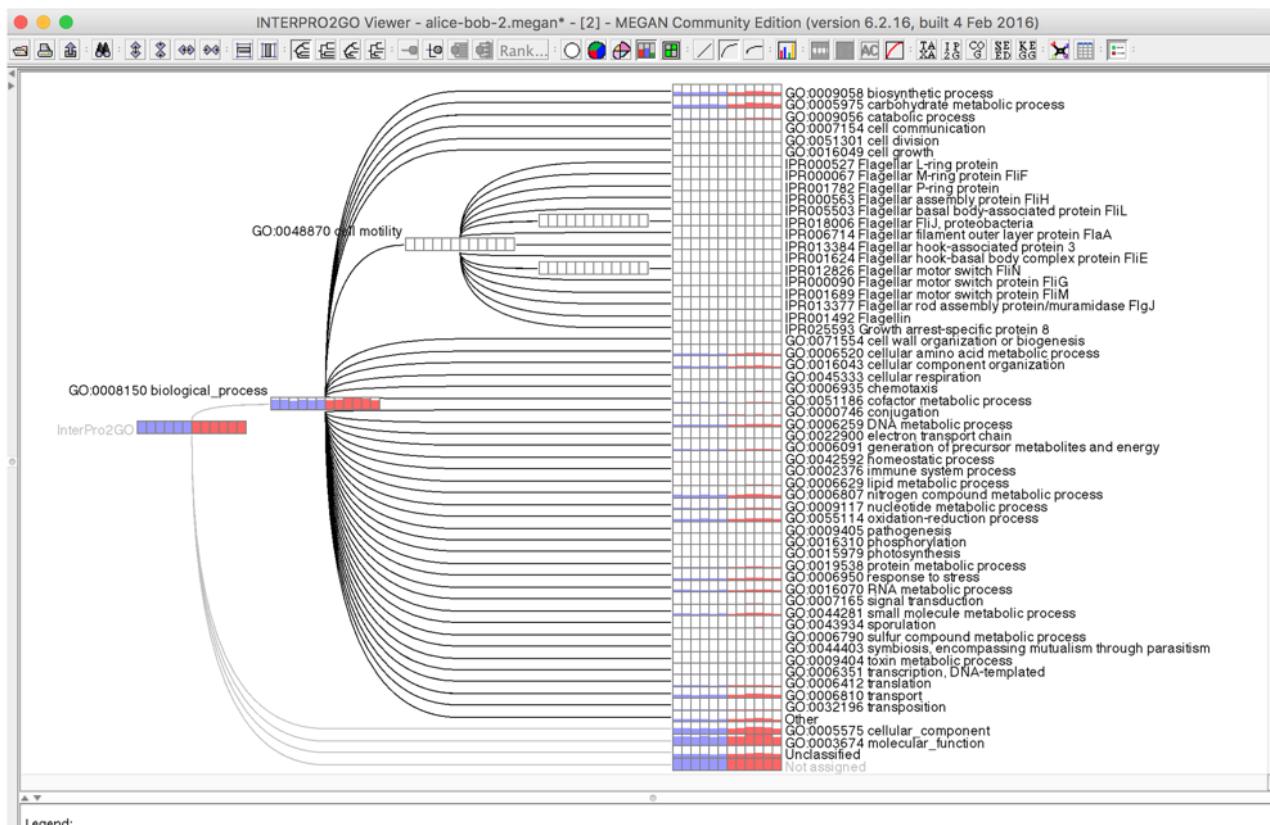
- Goal: Identify the relative abundance of different microbes in a sample
- Problems:
 - Reads are all mixed together
 - Reads can be short (~100bp)
- Two broad approaches
 1. Binning Based
 2. Marker Based

Binning Methods

- Attempts to “bin” reads into the genome from which they originated
- Composition-based
 - Uses GC composition or k-mers (e.g. Naïve Bayes Classifier)
 - Generally not very precise and not recommended
- Sequence-based
 - Compare reads to large reference database using BLAST (or some other similarity search method)
 - Reads are assigned based on “Best-hit” or “Lowest Common Ancestor” approach

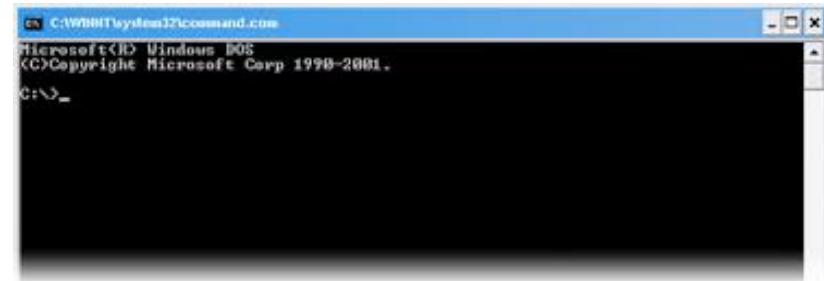
Binning Tools

- Notable Examples:
 - MEGAN: <http://ab.inf.uni-tuebingen.de/software/megan6/>
 - One of the first metagenomic tools
 - Does functional profiling



Binning Tools

- Notable Examples:
 - Kraken: <https://ccb.jhu.edu/software/kraken/>
 - Fastest binning approach to date and very accurate.
 - Large computing requirements (e.g. >128GB RAM)



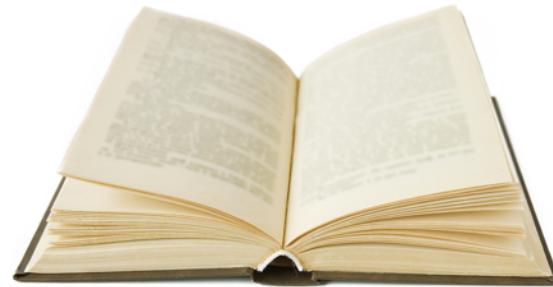
Marker Methods

- Single Gene
 - Identify and extract reads hitting a single marker gene (e.g. 16S, cpn60, or other “universal” genes)
 - Use existing bioinformatics pipeline (e.g. QIIME.)
- Multiple Gene
 - Several universal genes
 - PhyloSift Uses 37 universal single-copy genes

‘left-over’ Reads

- What about the reads that can’t be classified?
 - Errors?
 - Novel Sequence?
 - Just not in database?
 - Try assembling the left-overs

Metagenomic assembly



Metagenomic assembly

Specific assemblers

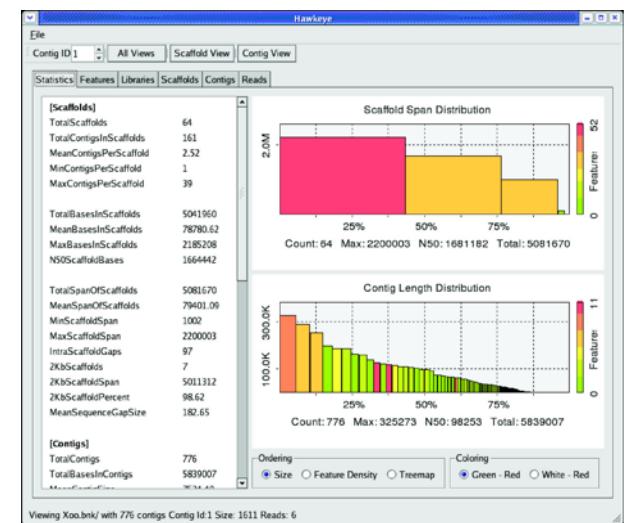
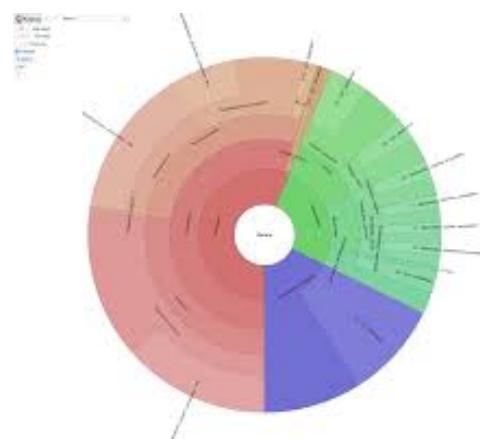
IDBA-UD,
Meta-IDBA,
Meta-Velvet,
Ray Meta,
Meta-Cortex

Validation

LAP,
ALE,
QUAST,
FRCbam,
REAPR

Analysis Pipelines

MetAMOS



Metagenome Function

What are they doing?

What do we mean by function?

- General categories
 - Photosynthesis
 - Nitrogen metabolism
 - Glycolysis
- Specific groups of orthologs
 - EC: I.I.I.I (alcohol dehydrogenase)
 - K00929 (butyrate kinase)

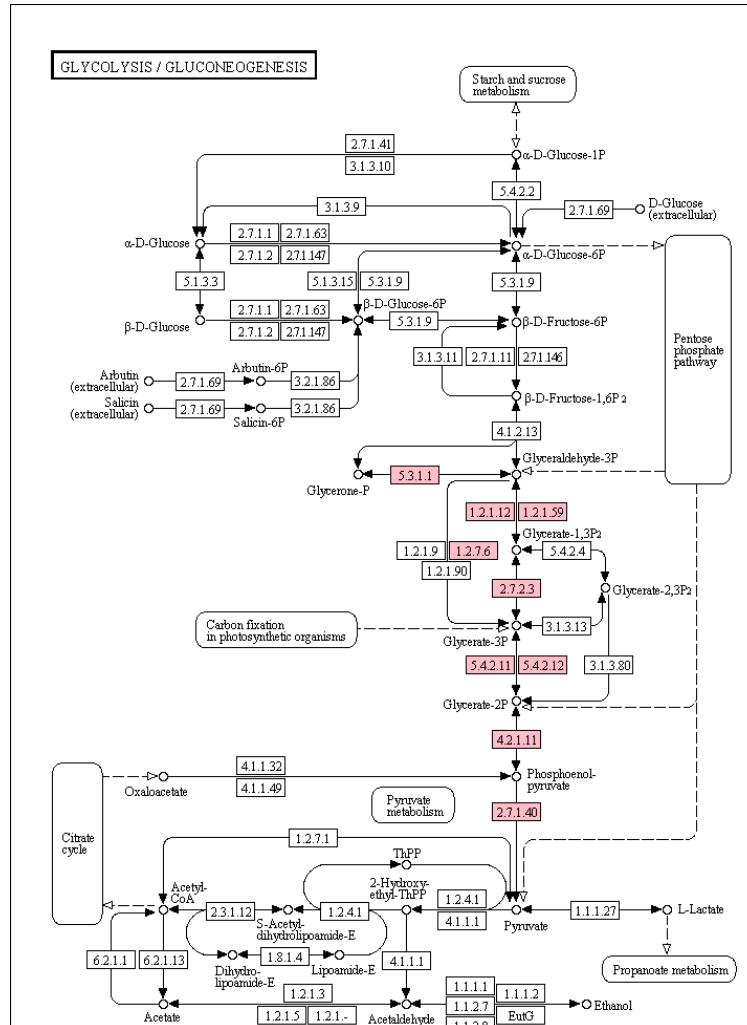
KEGG for functional annotation

- KEGG Orthologs (KOs) <http://www.genome.jp/kegg/ko.html>
 - Most specific. Thought to be homologs and doing the same exact “function”
 - ~12,000 KOs in the database
 - These can be linked into KEGG Modules and KEGG Pathways,
 - Identifiers: K01803, K00231, etc.

Entry	K01803	KO
Name	TPI, tpiA	
Definition	triosephosphate isomerase (TIM) [EC:5.3.1.1]	
Pathway	ko00010 Glycolysis / Gluconeogenesis ko00051 Fructose and mannose metabolism ko00562 Inositol phosphate metabolism ko00710 Carbon fixation in photosynthetic organisms ko01200 Carbon metabolism ko01230 Biosynthesis of amino acids	
Module	M00001 Glycolysis (Embden-Meyerhof pathway), glucose => pyruvate M00002 Glycolysis, core module involving three-carbon compounds M00003 Gluconeogenesis, oxaloacetate => fructose-6P	
Disease	H00664 Anemia due to disorders of glycolytic enzymes	
Brite	KEGG Orthology (KO) [BR: ko00001] Metabolism Overview 01200 Carbon metabolism K01803 TPI, tpiA; triosephosphate isomerase (TIM) 01230 Biosynthesis of amino acids K01803 TPI, tpiA; triosephosphate isomerase (TIM) Carbohydrate metabolism 00010 Glycolysis / Gluconeogenesis	

KEGG for functional annotation

- KEGG Pathways
 - Groups KOs into large pathways (~230)
 - Each pathway has a graphical map
 - Individual KOs or Modules can be highlighted within these maps
 - Pathways can be collapsed into very general functional terms (e.g. Amino Acid Metabolism, Carbohydrate Metabolism, etc.)



Metagenomic Annotation tools

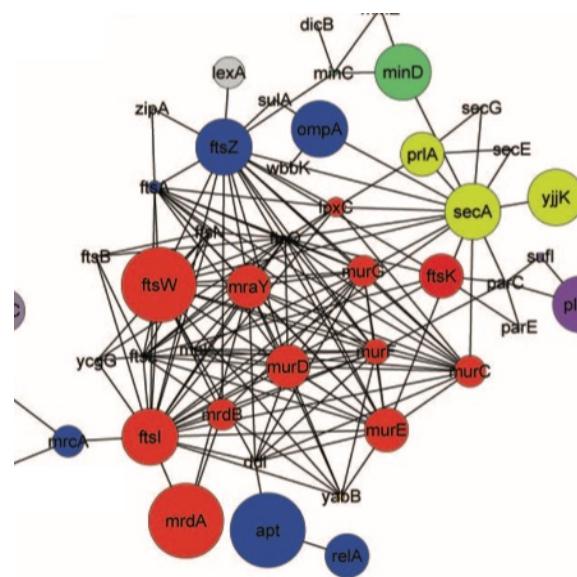
- Web-based
 - EBI Metagenomics Server
 - MG-RAST
 - IMG/M
 - Koala
- GUI-Based
 - MEGAN
 - Allows connection between taxonomy and function
 - ClovR
 - Virtual Machine based - microbial
- Local-based
 - Prokka
 - MetAMOS
 - Built in assembly, highly customizable,
 - DIY
 - Set up your own in-house custom computational pipeline

Metatranscriptomics for community activity

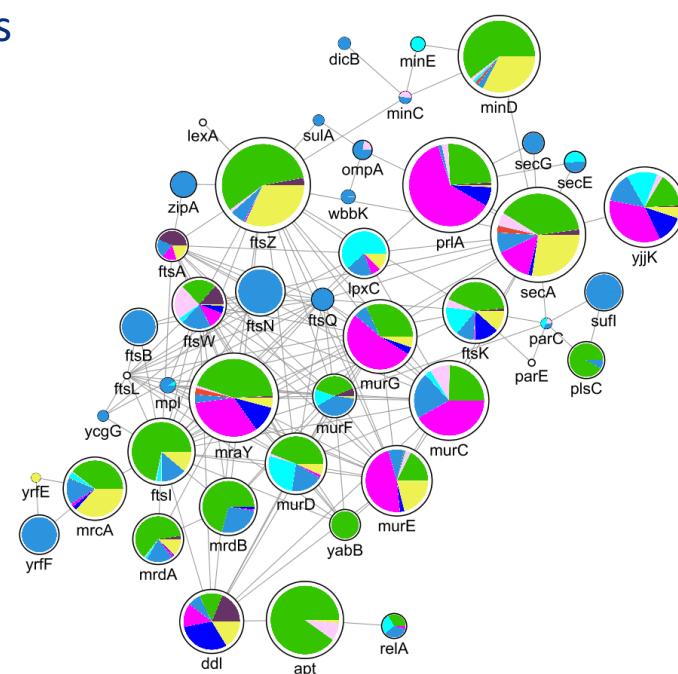
Functional Quantification

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

Genes involved in pathways associated with cell wall biogenesis

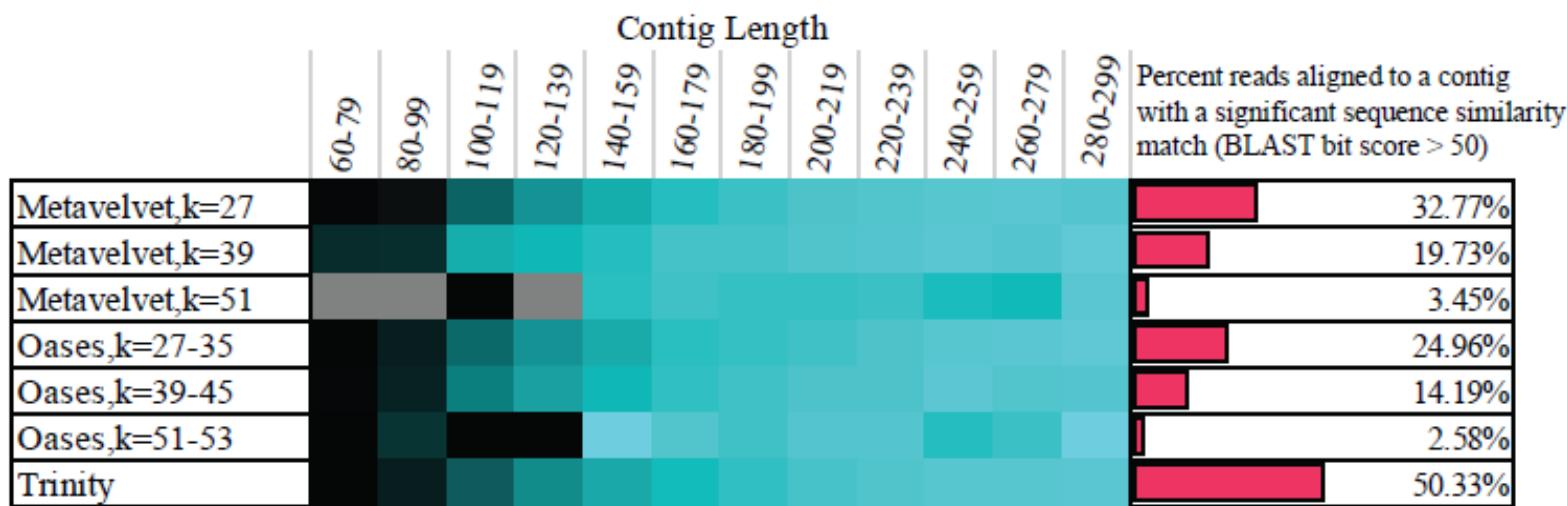
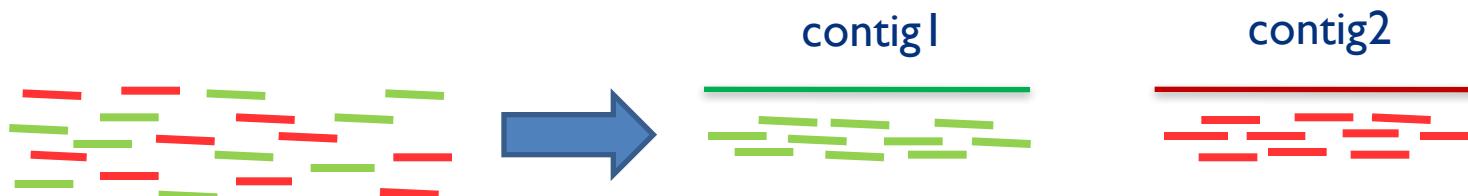


Relative
abundance of
transcripts



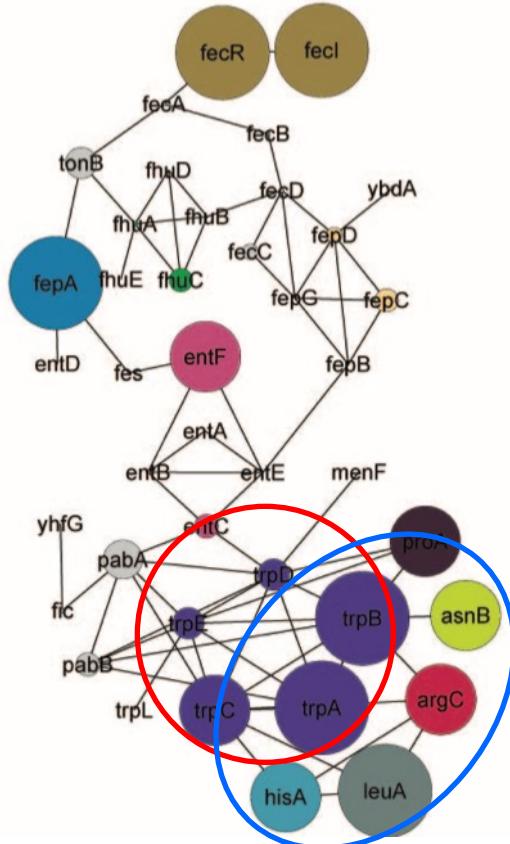
Relative
abundance and
contributing taxa

Metatranscriptomics: Assembling to annotate



Metatranscriptomics: Understanding results

How do we identify those systems which are differentially regulated across samples (e.g. systems up-regulated in disease vrs. health)?



Gene Set Enrichment Analysis

- Node Size = number genes with an annotation

Network analysis

– genes with related function and similar changes in expression

Metatranscriptomics: Understanding results

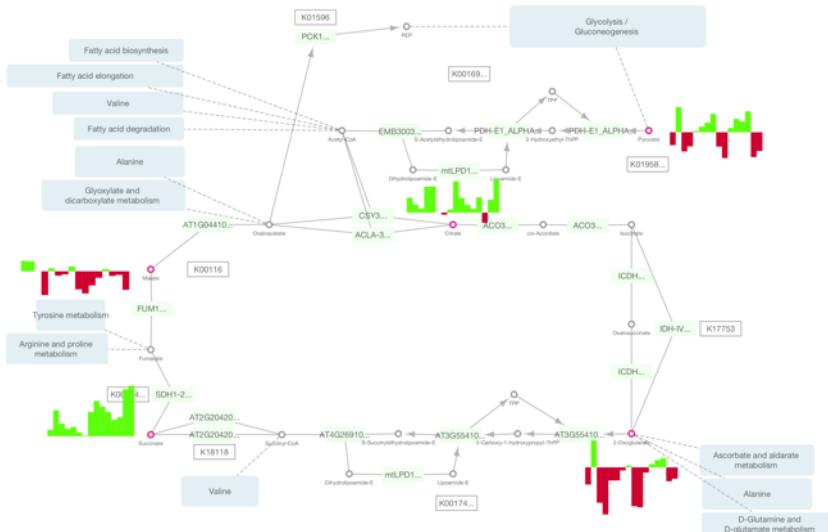
Cytoscape

www.cytoscape.org

MeV

www.tm4.org/mev.html

R
D3.js



Genes involved in pathways associated with cell wall biogenesis

