

NGS – part 2: applications

Tobias Österlund
tobiaso@chalmers.se

NGS part of the course

Week 4	Friday 13/2	15.15-17.00	NGS lecture 1: Introduction to NGS, alignment, assembly
Week 6	Thursday 26/2	08.00-09.45	NGS lecture 2: RNA-seq, metagenomics
Week 6	Thursday 26/2	10.00-11.45	NGS computer lab: Resequencing analysis
Week 7	Thursday 5/3	10.00-11.45	Marcela: Exome sequencing
Week 8	Monday 9/3	17.00	Deadline: Essay on NGS and metagenomics
Week 8	Thursday	08.00-09.45	Fredrik: HMMer and Metagenomics

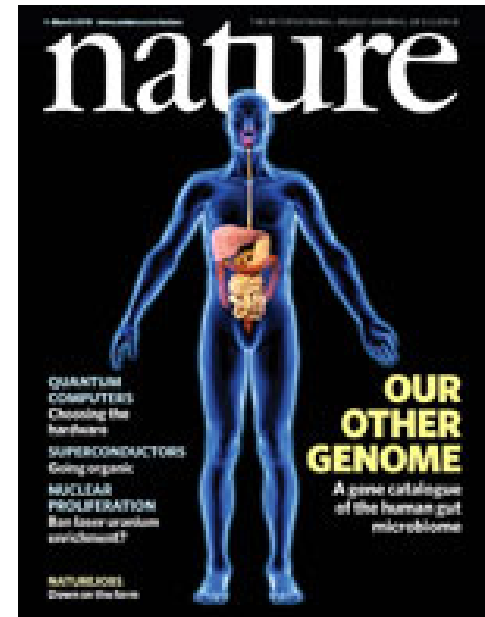
Today's lecture

- Metagenomics analysis
 - On the species level: Who's there?
 - On the gene/functional level: What are they doing?
- RNA-seq analysis
 - Data normalization
 - Finding differentially expressed genes
- Computer exercise
 - Whole genome sequencing for variant detection

Metagenomics

- Some facts about microbes

Number of microbes on Earth	5×10^{30}
Number of microbes in all humans	6×10^{23}
Number of stars in the universe	7×10^{21}
Number of bacterial cells in one human gut	10^{14}
Number of bacterial cells in one human gut	10^{13}
Number of bacterial genes in one human gut	3,000,000
Number of genes in the human genome	21,000

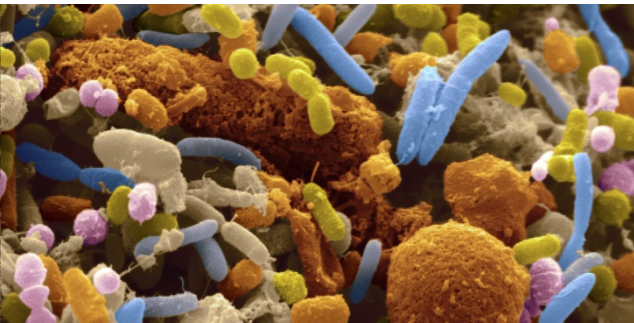


Microbial diversity

- Bacteria are present in every habitat on Earth
- There are up to 100 million bacterial species
 - only a small fraction of these are known
- More than 99% of all bacteria are uncultivable and can not be grown in laboratories

Shotgun metagenomics

- Analysis of all DNA in a sample – the metagenome



Sample with high diversity



High throughput sequencing

```
GCAACAGTTTGGCGTAATTCAATTGT  
CAGTTTACGGATTCCCTTGATTGGATAA  
TCCAGTCTGCCCCAGGCTGCAGTTGC  
AAAAGAAAGAAACGACTATGAATAAAC  
GACTTCGGATCATTGGACTGTTTGCTG  
TGTTCTTTGGCCAGATGATCCACGCGC  
AGACCACAGCGTTCCTTATCAGGGGC  
GTCTCAATGACAACGGCGCGCTGGCCA  
ACGGCATTATGATTGAAATTTTCAC  
TATACACCGTGGCGACCAATGGCAGTG  
CCTCATCGTCGCGGTCAAATGCCGCCA  
CCGTCGTCAG
```

Metagenome

Metagenomics

- Metagenomics is used to study the unculturable organisms and viruses
 - ~50% of human gut bacteria are unculturable
 - <1% of environmental bacteria are unculturable
- Metagenomes are highly fragmented and undersampled
- The majority of DNA found in metagenomes is usually very hard to annotate

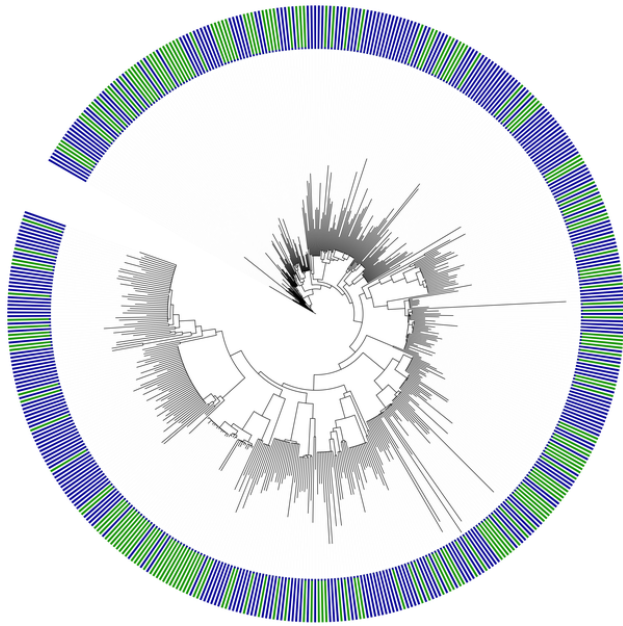
Two types of questions

Who's there?

- Identification of species, phylum etc.
- Estimation of species abundance

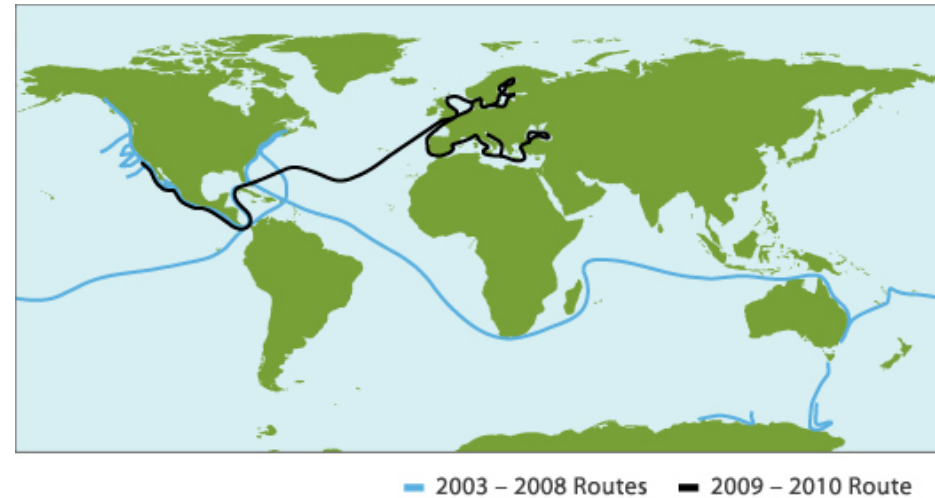
What are they doing?

- Functional annotation (gene families / pathways)
- Estimation of gene/ pathway abundance



The global ocean sampling

- Investigating microbial diversity in the ocean
- A sailing boat equipped with a sequencer

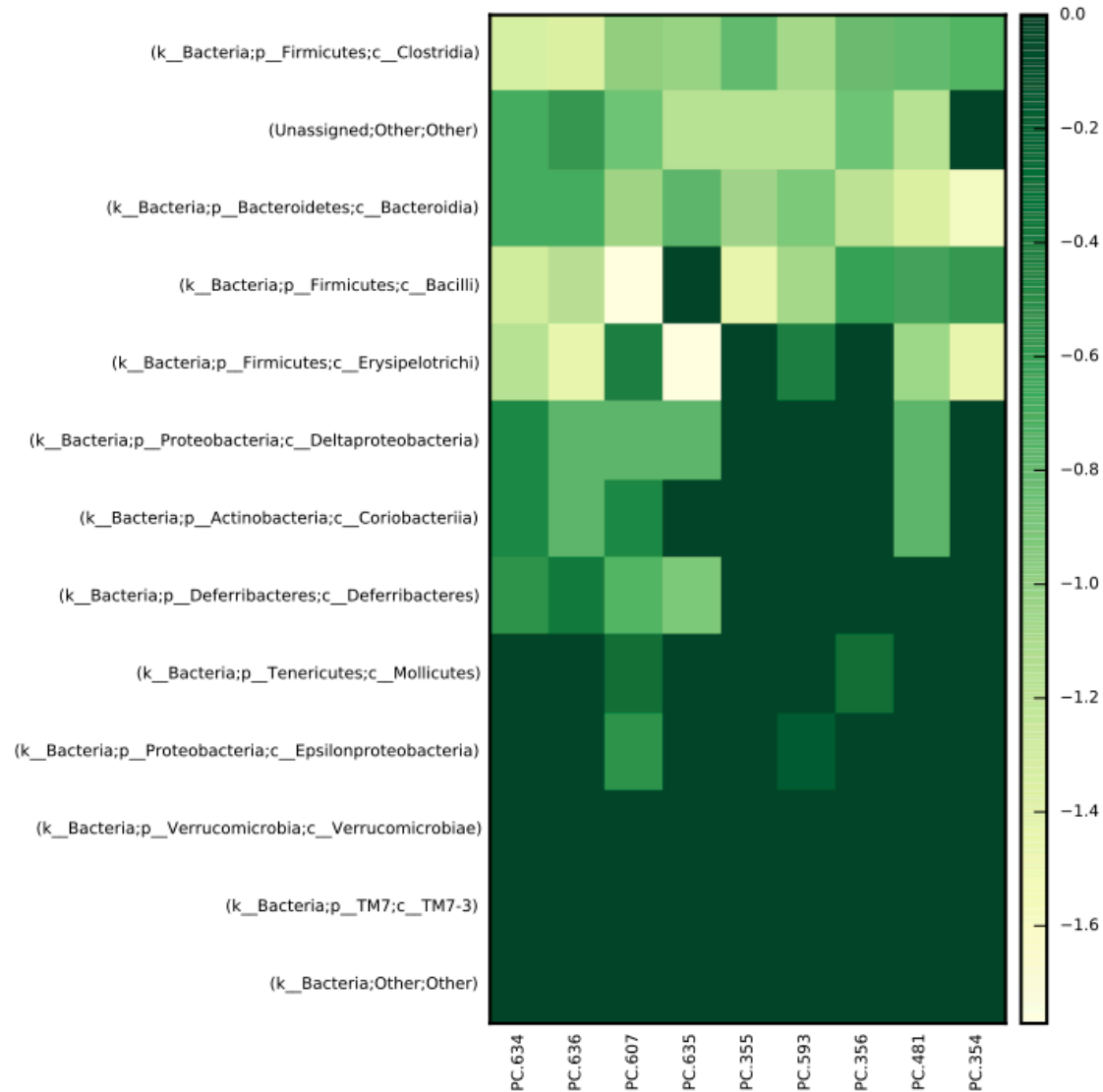


Species identification

- Prokaryots:
 - 16s rRNA gene
- Eukaryots:
 - 18s rRNA gene
- Can be amplified using PCR (amplicon sequencing)
- Sequences mapped to known species using BLAST
- Operational taxonomic unit (OTU):
 - 97% sequence similarity for the 16s rRNA gene
 - Cluster based on sequence similarity using UCLUST

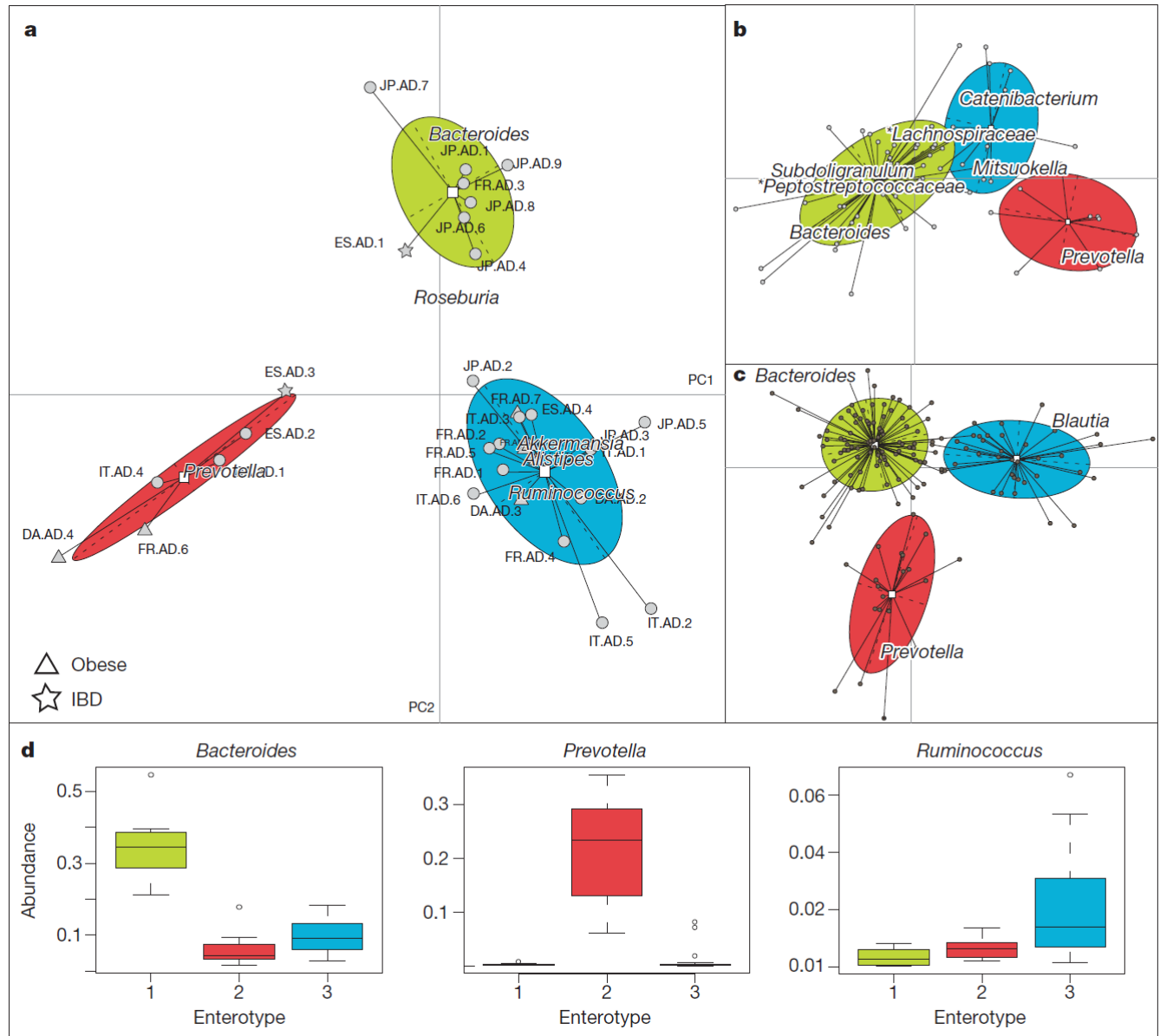
Species abundance

- Qiime
 - Bioinformatics program available at qiime.org
 - Pick OTUs
 - Analysis of species abundance
 - Bioinformatics analysis



Enterotypes of the human gut

- Map reads to a gene catalog with 1500 known species
- Cluster based on species abundance



Metagenome assembly

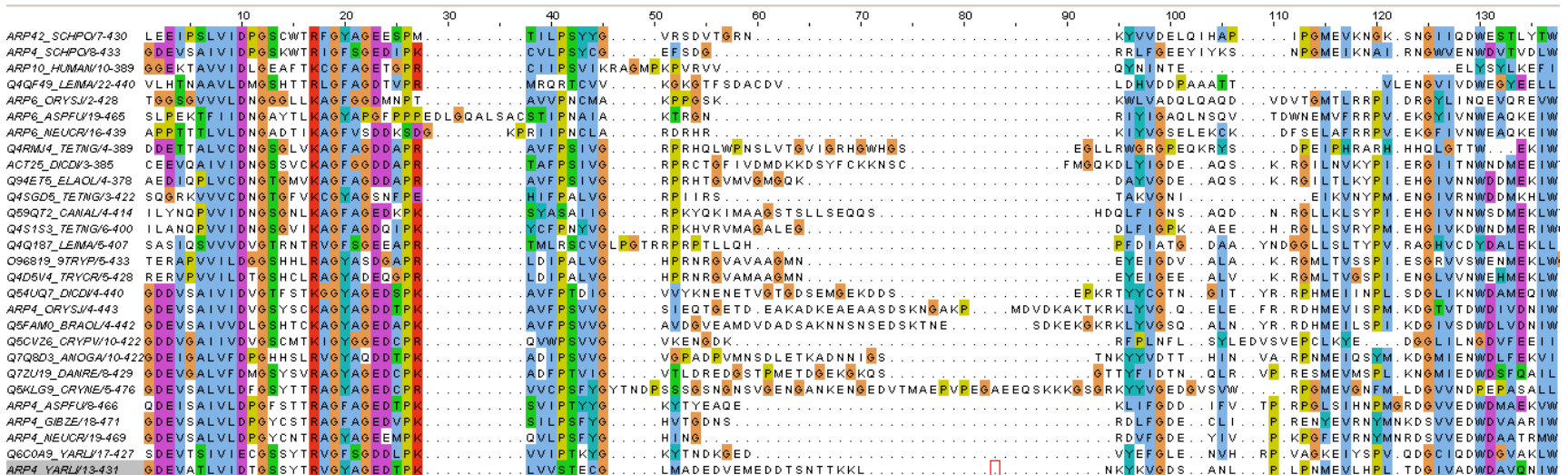
- Reference genomes are missing for majority of the bacteria we are studying
- Need to annotate the sequences
- Annotation (using e.g. blast) is easier for longer fragments
- De novo assembly of reads into longer contigs
- Can be hard due to large amount of data
 - Need large amount of memory
- Can be hard due to undersampling of the metagenome
 - Can't get a complete assembly

Metagenome assembly software

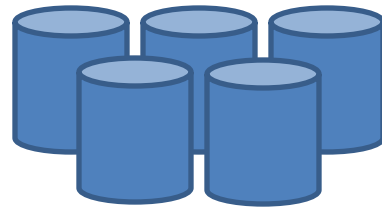
- Velvet
 - Metavelvet
 - MAQ
 - SOAP *de novo*
 - Etc.
-
- Most assemblers uses deBruijn graphs
 - Kmers
 - Need to specify k

Functional analysis

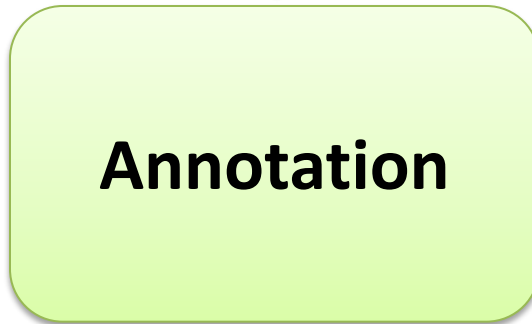
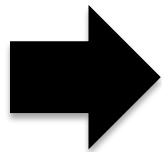
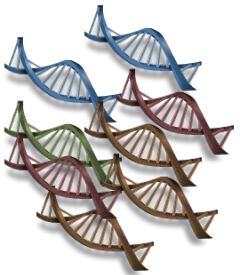
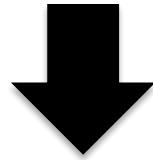
- "Gene centric analysis" (What are they doing?)
- Only a small fraction of the bacterial genomes have been sequenced.
- Annotation done using protein profiles catching the variability (PFAM, TIGRFAM, COG, etc)



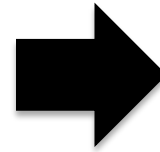
PFAM domain for actin.



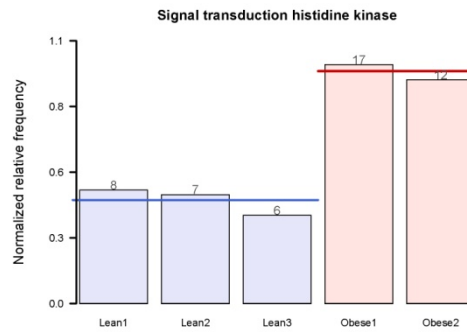
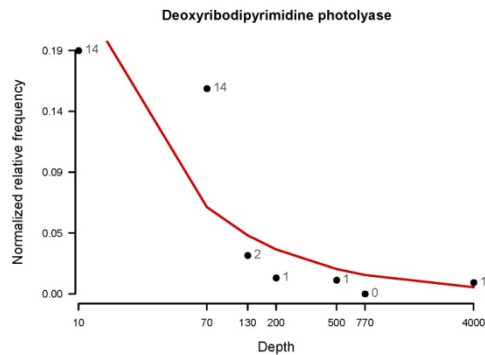
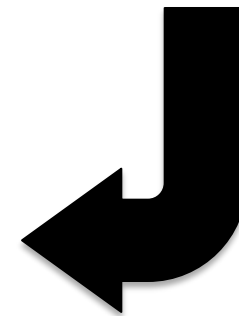
Reference database

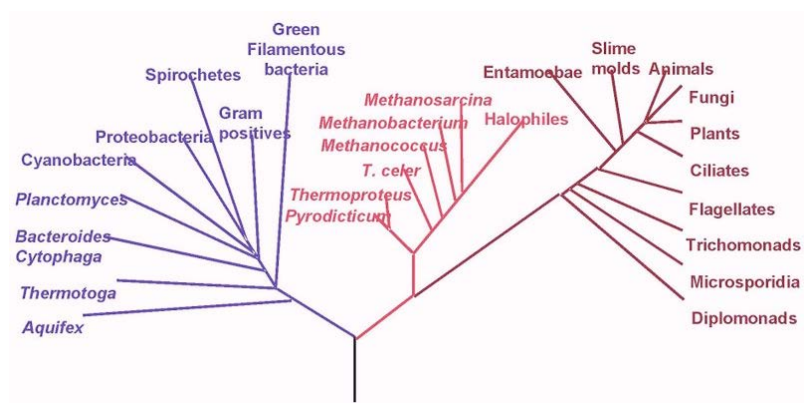
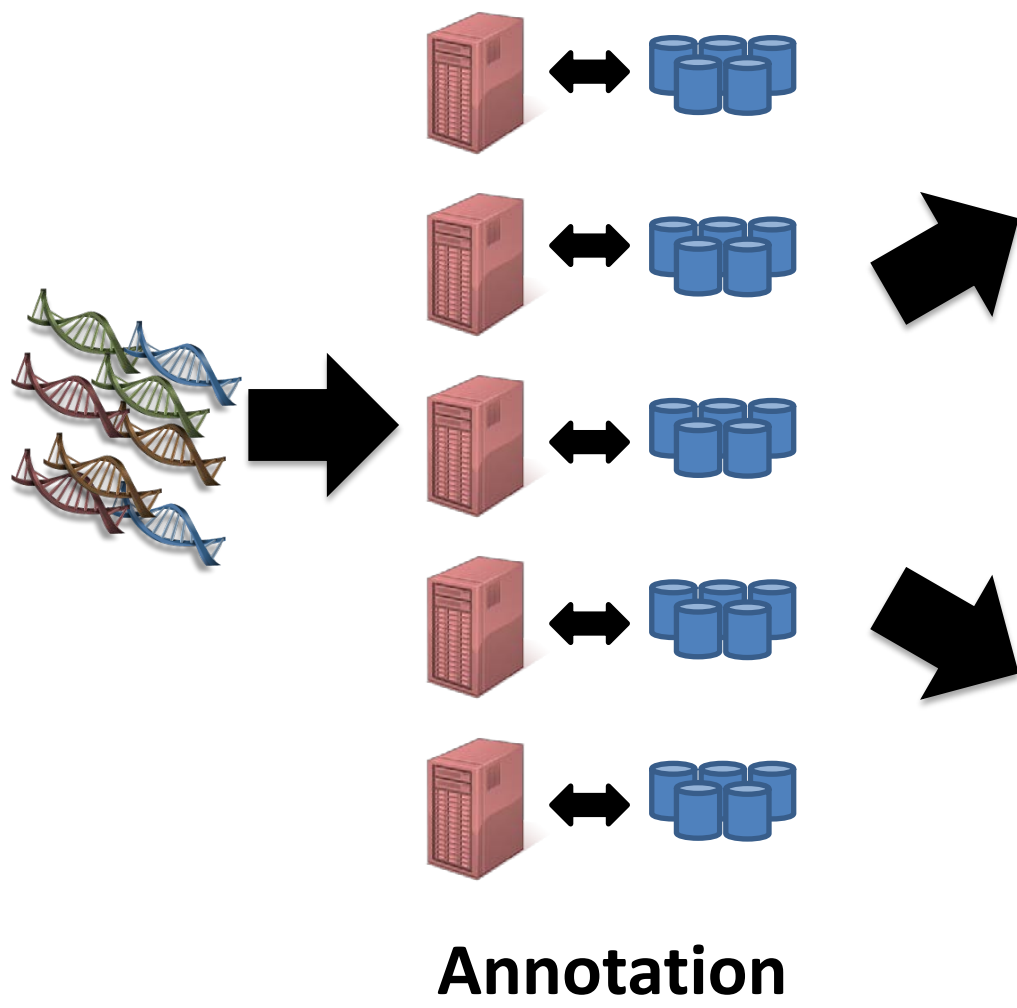


ShotgunAnnotatorR



ShotgunFunctionalizeR





Taxonomic affiliation

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5
Gene1	591	536	1260	284	19
Gene2	28	21	19	36	10
Gene3	53	51	97	118	36
Gene4	106	149	266	47	11
....					
....					

Gene occurrences

Identification of significant genes

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5
Gene1	591	536	1260	284	19
Gene2	28	21	19	36	10
Gene3	53	51	97	118	36
Gene4	106	149	266	47	11
....					
Gene1312	243	362	163	258	423
Gene1313	13	43	23	67	34
....					
Total	132 567	80 456	197 723	73 491	134 513

Normalization

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5
Gene1	591	536	1260	284	19
Gene2	28	21	19	36	10
Gene3	53	51	97	118	36
Gene4	106	149	266	47	11
....					
Gene1312	243	362	163	258	423
Gene1313	13	43	23	67	34
....					
Total	132 567	80 456	197 723	73 491	134 513

$X_{i,j}$ (arrow pointing to Gene4, Sample 1)

n_j (arrow pointing to Total, Sample 1)

$X_{i,j}$ - number of reads matching gene i in sample j

n_j - normalization factor per sample

$$R_{i,j} = \frac{X_{i,j}}{n_j}$$

Normalization

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5
Gene1	0.004458	0.006662	0.006373	0.003864	0.000141
Gene2	0.000211	0.000261	9.61E-05	0.00049	7.43E-05
Gene3	0.0004	0.000634	0.000491	0.001606	0.000268
Gene4	0.0008	0.001852	0.001345	0.00064	8.18E-05
....					
Gene1312	0.001833	0.004499	0.000824	0.003511	0.003145
Gene1313	9.81E-05	0.000534	0.000116	0.000912	0.000253
....					
Total	1	1	1	1	1

How to normalize metagenomic data?

$$R_{i,j} = \frac{X_{i,j}}{n_j}$$

- n_j – normalization factor per sample
- Divide with total number of reads mapped in each sample?
- Divide with the total number of reads in each sample
- Divide with the total number of reads mapping to the 16s rRNA gene in each sample?
- More advanced method?

Identification of significant genes

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5
Gene1	591	536	1260	284	19
Gene2	28	21	19	36	10
Gene3	53	51	97	118	36
Gene4	106	149	266	47	11
....					
Gene1312	243	362	163	258	423
Gene1313	13	43	23	67	34
....					
Total	1 32 567	80 456	1 97 723	73 491	1 34 513

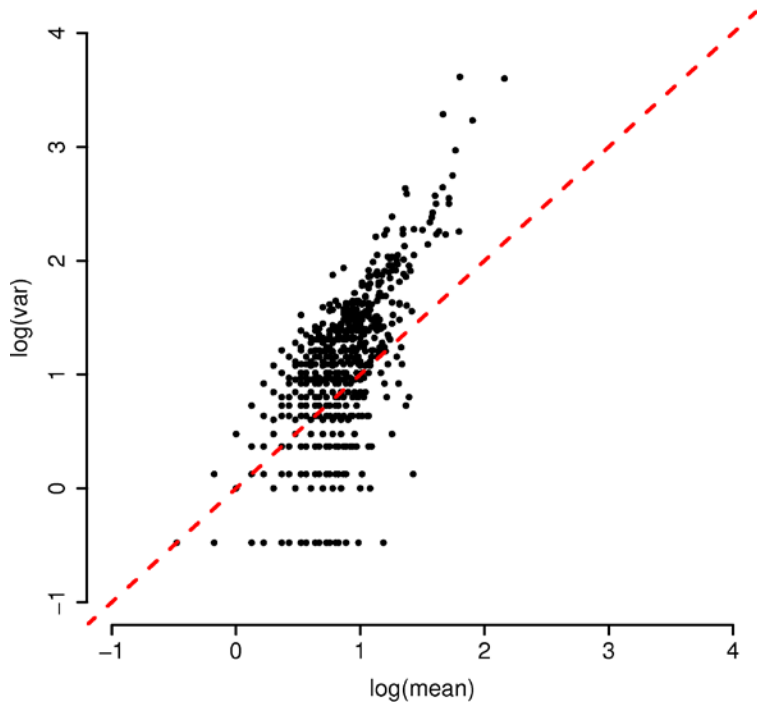
$$\log \left(\frac{E[X_{i,j}]}{n_j} \right) = \alpha_0 + \sum \alpha_k y_k$$

Baseline Covariates (groups)

Statistical analysis

- Data from metagenomics is discrete (counts per gene/species)
- Not normally distributed
- $X_{i,j} \sim \text{Poisson}(\lambda_i)$
 $E[X_{i,j}] = \lambda_i$
 $\text{Var}[X_{i,j}] = \lambda_i$

Statistical analysis



- $\text{Var}[X_{i,j}] > \text{E}[X_{i,j}]$
- Overdispersed data!

$$\text{Var}[X_{i,j}] = \phi \lambda_i$$

Estimated from the
total residual sum

- The proportion of false positives are estimated using Benjamini-Hochberg's false discovery rate.

Summary metagenomics

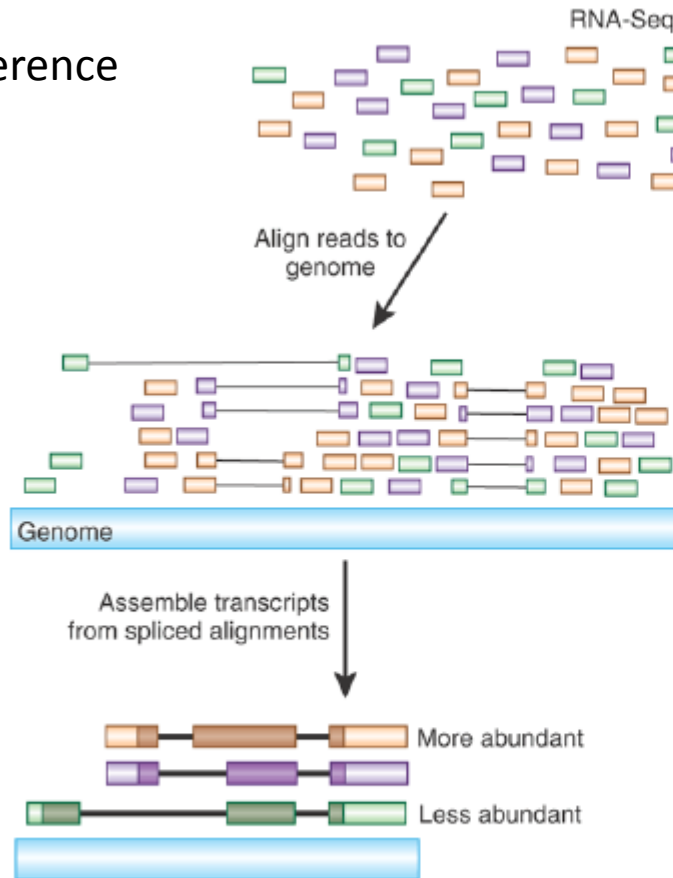
- Metagenomics provides a powerful way to do culture-independent analysis of bacterial communities
- The low cost of next generation sequencing have increased the power of metagenomics substantially
- Examples of metagenomics studies of microbial communities in the human gut and from environmental samples

RNA-seq

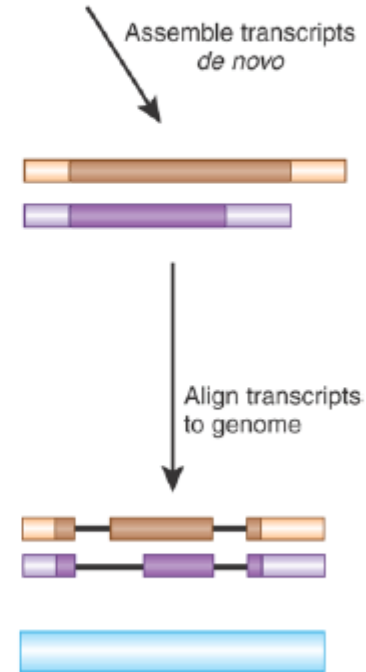
- Large-scale mRNA quantification
 - Identification of differentially expressed genes
 - Sequence all mRNA and map to reference sequence
- De novo transcriptome assembly
 - Find new transcripts
 - Alternative splicing
 - When no reference sequence is available
 - Map the reads back to the newly assembled contigs
 - Can help in genome annotation

RNA-seq analysis strategy

Good reference



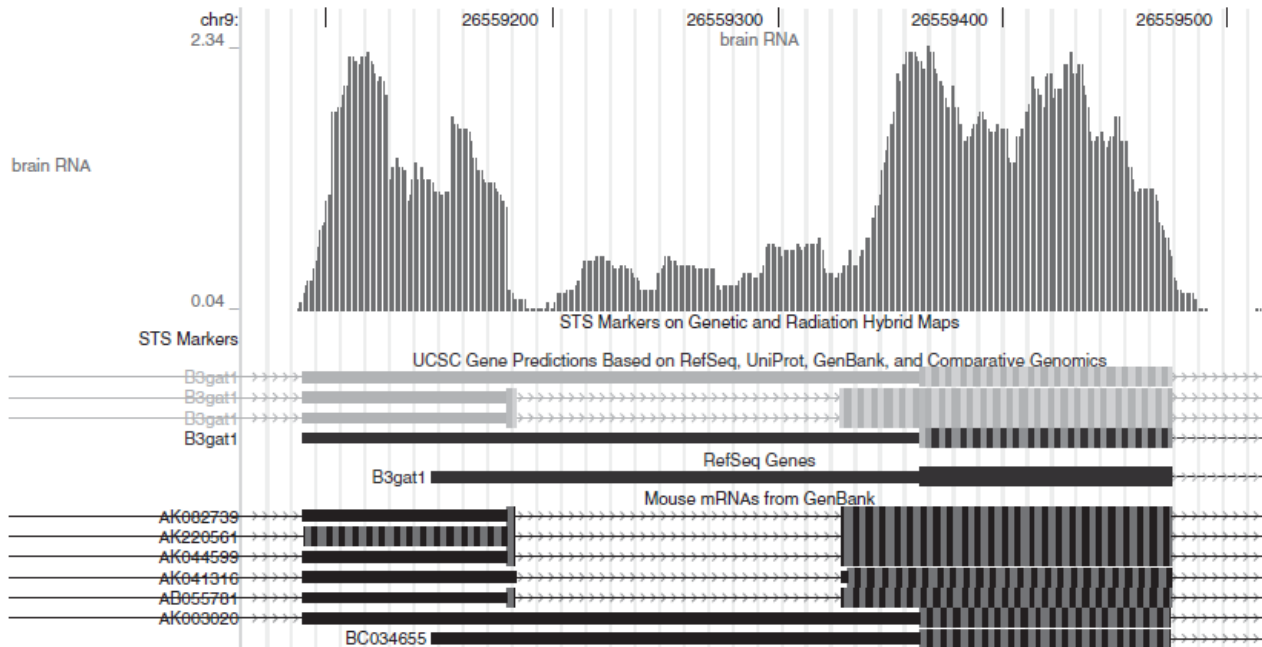
No genome



Haas and Zody, Nature Biotechnology 28, 421–423 (2010)

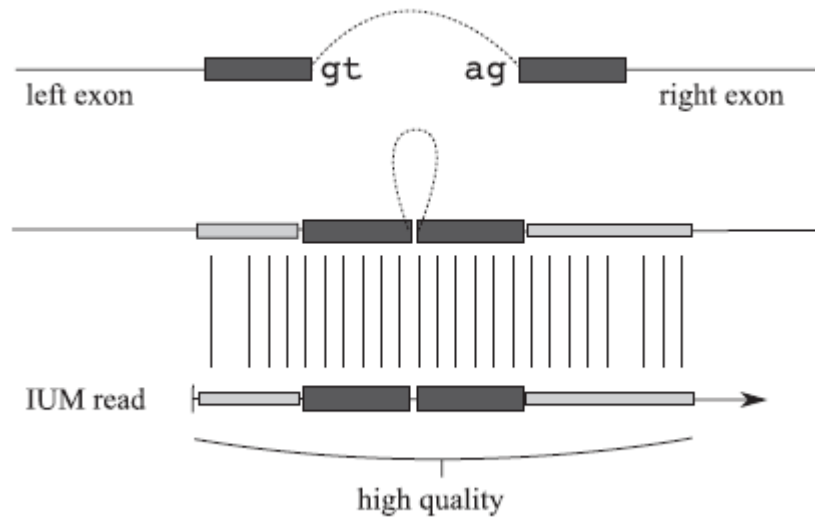
Alignment

- Using a splice-aware aligner



TopHat aligner (Trapnell et al. Bioinformatics 2009)

Alignment



TopHat aligner (Trapnell et al. Bioinformatics 2009)

De novo transcriptome assembly



Trinity command line example:

```
Trinity --seqType fq --left reads_1.fq --right reads_2.fq --CPU 6 --max_memory 20G
```

- Inchworm assembles the transcripts
- Chrysalis and Butterfly estimates possible splice variants from the data

Statistical analysis

- Data from RNA-seq comes as reads/fragments per gene
 - $X_{i,j}$ = number of reads matching gene i in sample j

	Treatment A			Treatment B		
	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
Gene1	66489	29192	18643	21721	84669	80540
Gene2	11288	2899	1062	6130	9581	17251
Gene3	44979	12906	14604	10378	85043	39478
Gene4	7133	4772	1124	319	6863	7286
Gene5	34282	14379	13748	6133	12648	7620
Gene6	6531	7184	1962	651	1334	13125
Total	170702	71332	51143	45332	200138	165300

Data normalization

$$R_{i,j} = \frac{X_{i,j}}{n_j}$$

- n_j – normalization factor per sample
- Divide with total number of reads mapped in each sample?
- House keeping genes have a large influence on the normalization
- Robust scaling (Anders and Huber 2010)

$$n_j = \text{median}_i \frac{X_{i,j}}{\left(\prod_{j=1}^m X_{i,j} \right)^{1/m}}$$

RNA-seq is semi-quantitative

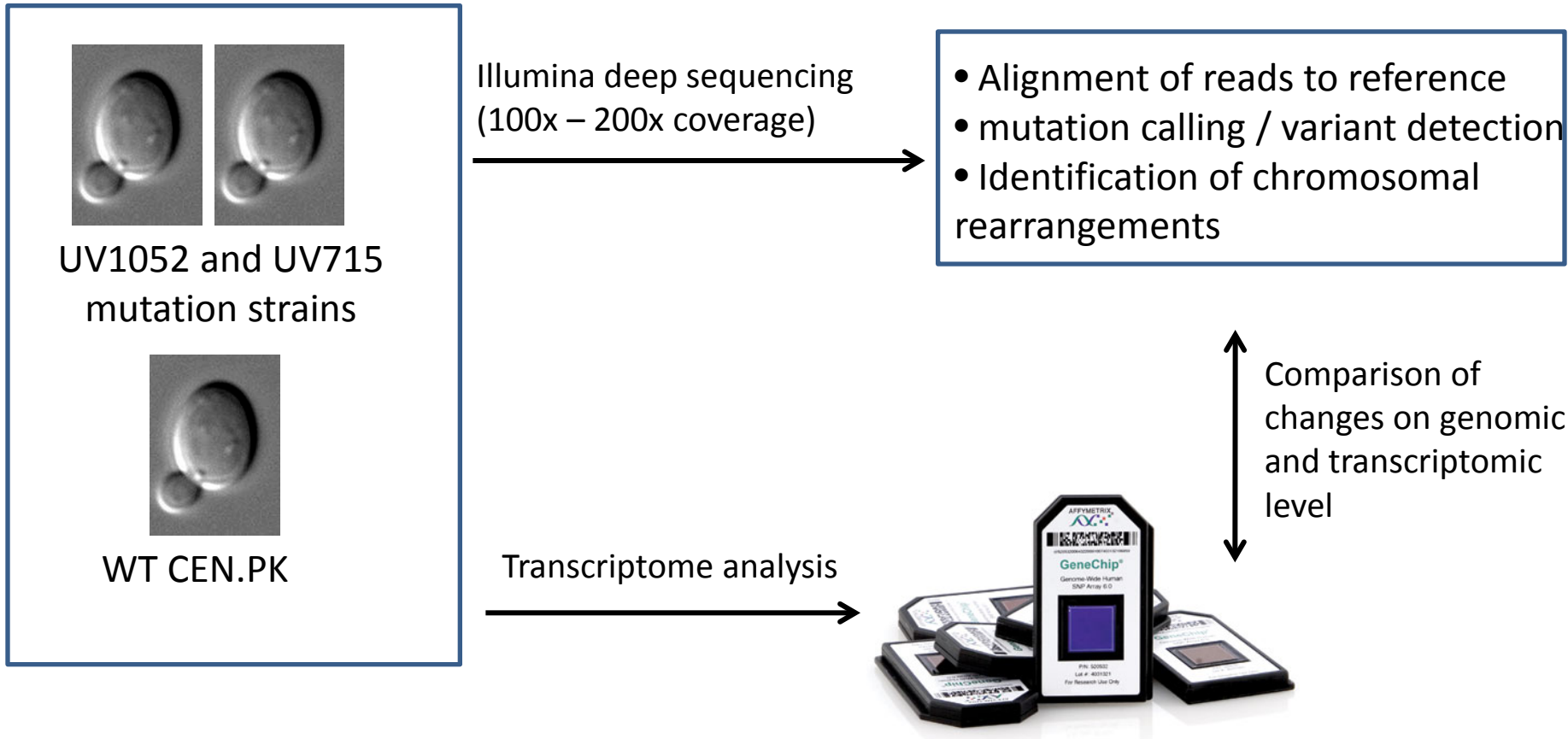
- Compare the same gene over different conditions
 - calculate fold-change and p-value
- Difficult to compare two genes from the same samples
 - Genes have different lengths
 - Genes have different GC-content (PCR-bias)

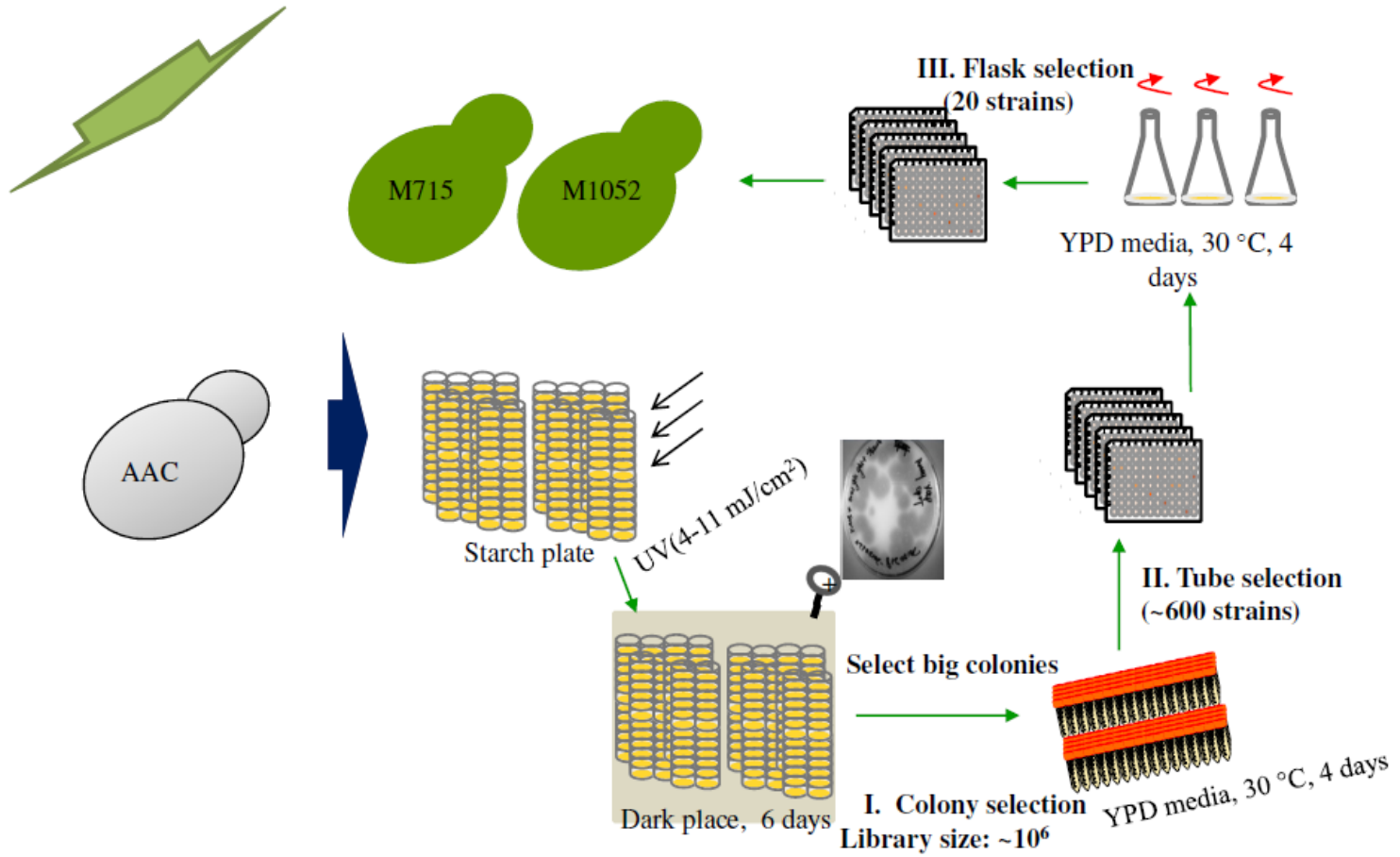
Study design

- How much should I sequence?
 - Depends on your question
 - Metagenomics: Sequence as much as possible
 - Your metagenome will still be undersampled
 - Need a lot of sequence to do assembly
 - RNA-seq: Sequence deep enough (enough coverage) to be able to detect both highly expressed transcript and rare transcripts
- Biological Replicates!!!

Sequencing lab

Genome sequencing of amylase producing yeast strains





Software used in lab

- Fastx toolkit – programs for preprocessing and quality control of Fastq and fasta files
- BWA – short read aligner
- Samtools – handling SAM and BAM files
- Integrative Genomics Viewer (IGV) – A genome browser viewing alignments (BAM-files)