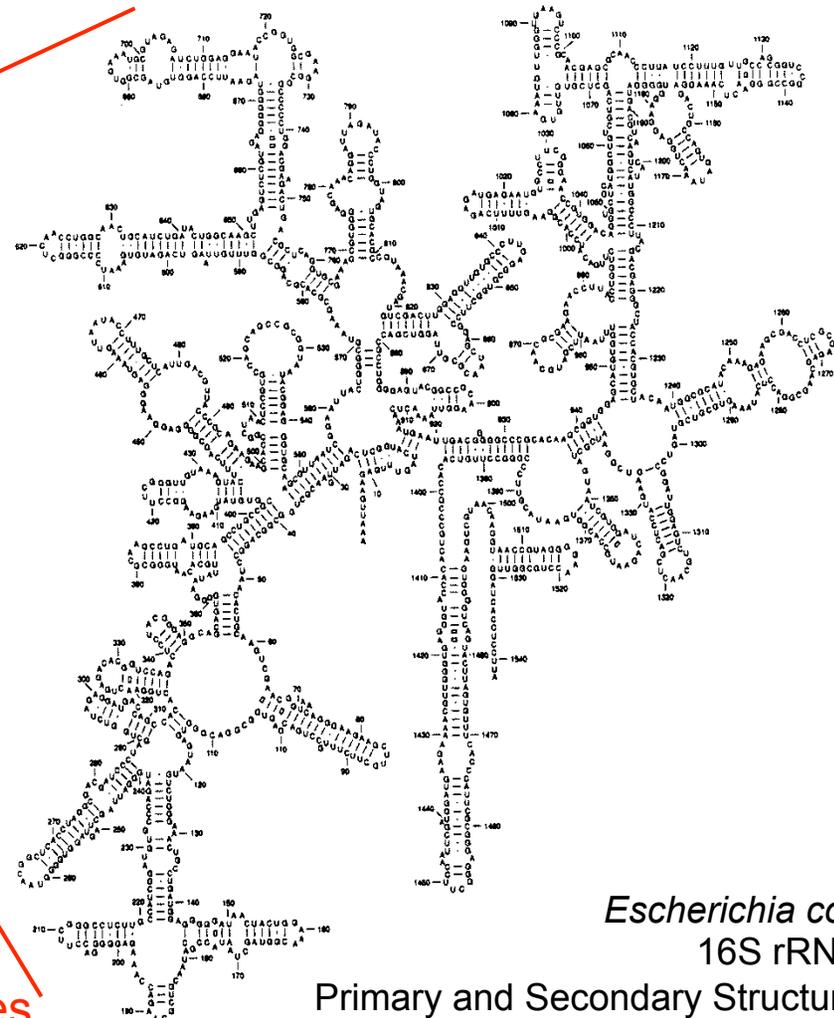
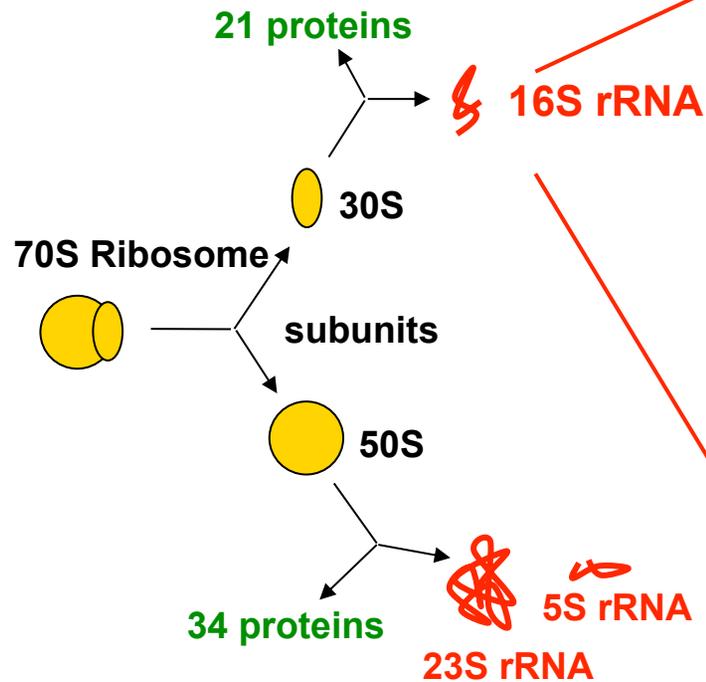




Evaluation of multiplexed 16S rRNA microbial population surveys using Illumina MiSeq platform

Julien Tremblay, PhD
SFAF, Santa Fe, NM – June 7th 2012

16S rRNA as phylogenetic marker gene

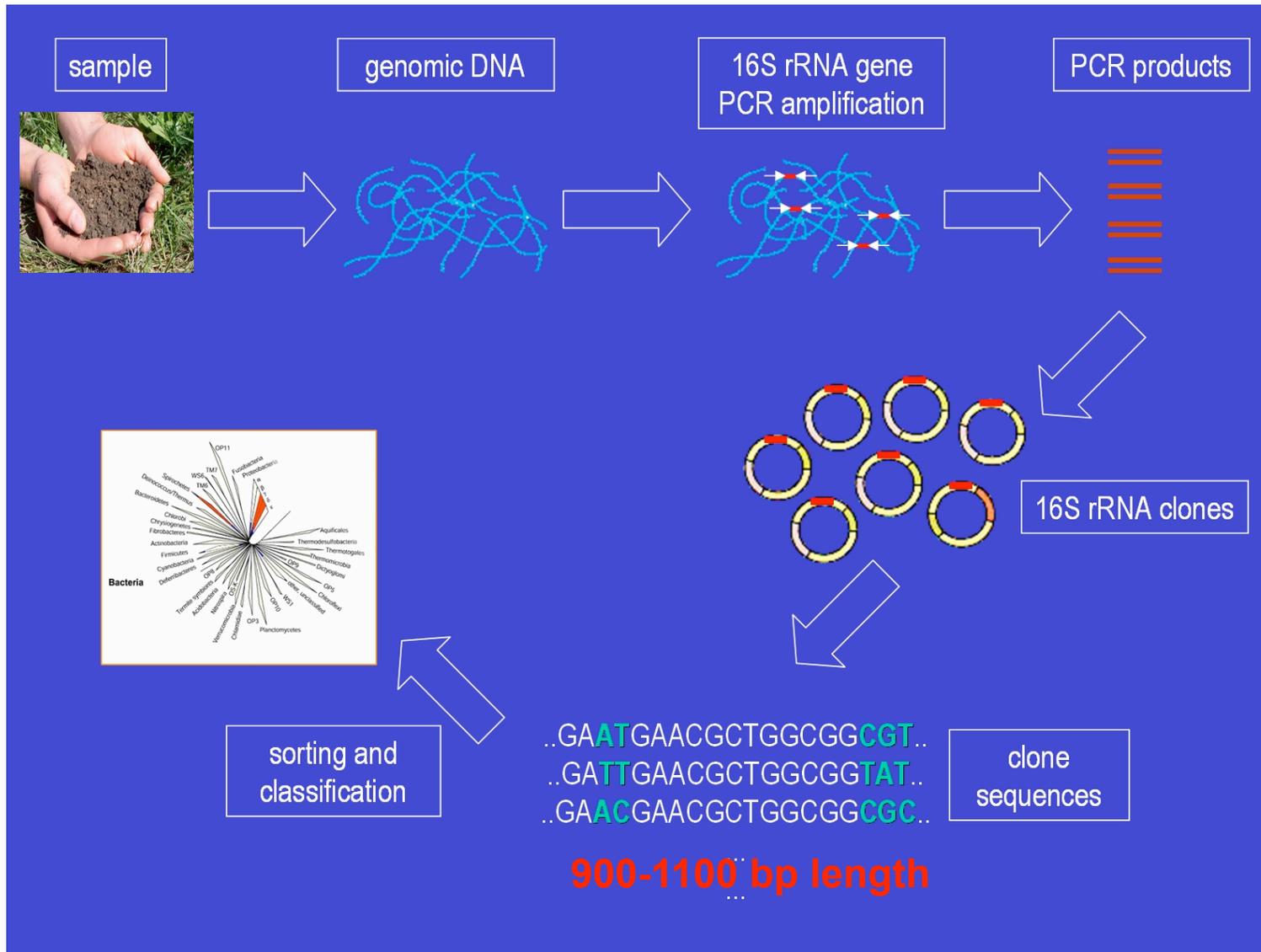


Escherichia coli
16S rRNA

Primary and Secondary Structure

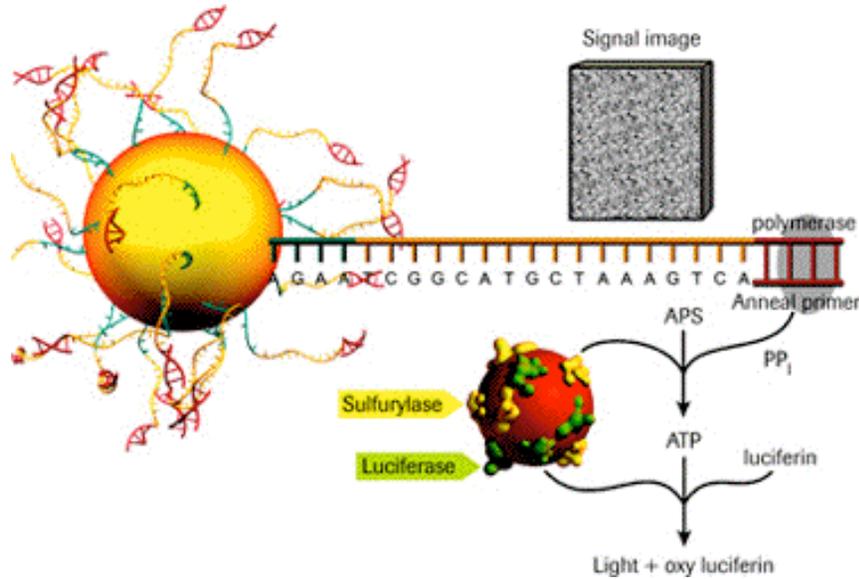
highly conserved between different species
of bacteria and archaea

16S rRNA in environmental microbiology (Sanger clone libraries)



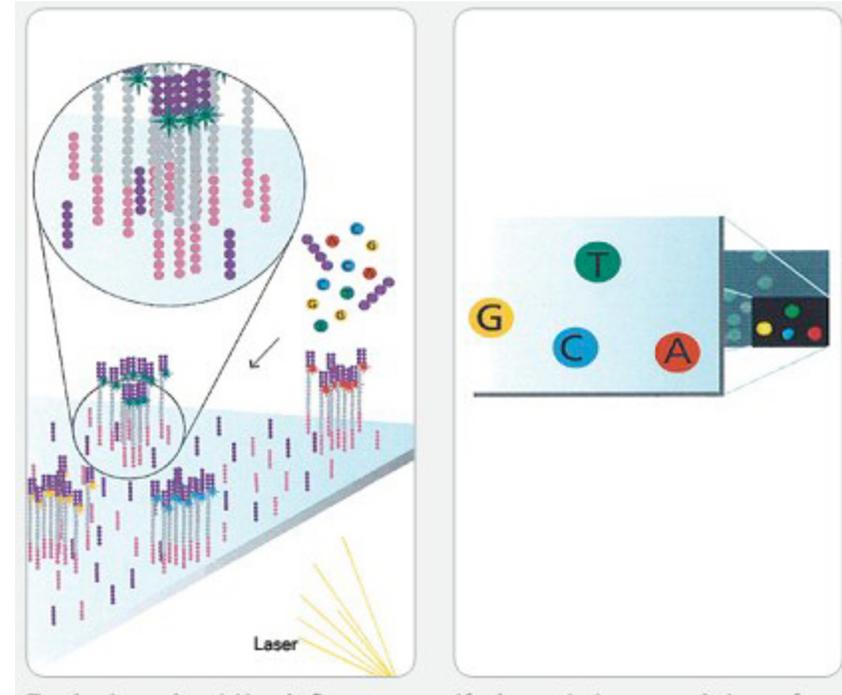


454

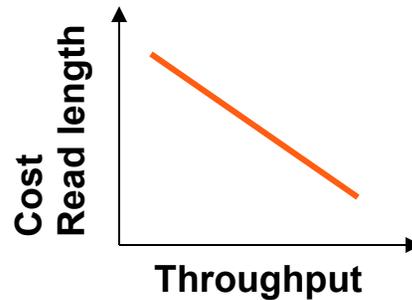


0.5M 450bp reads
\$\$

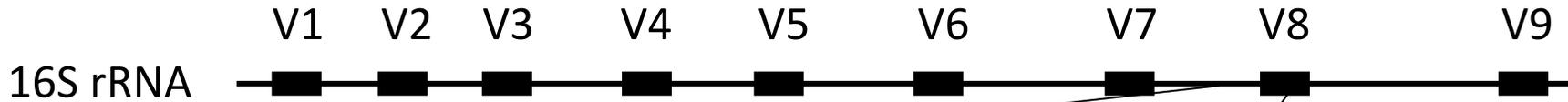
Illumina



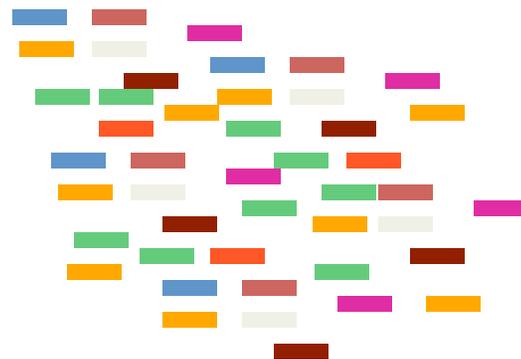
10-350M 150bp reads/lane
\$



Surveying microbial diversity with short 16S amplicons



Generate amplicons of a given variable region from bacterial community (many millions of sequences)



- █ X 10
- █ X 1
- █ X 1,000
- █ X 2,000
- █ X 200
- █ X 1,200
- █ X 800
- █ X 10,000

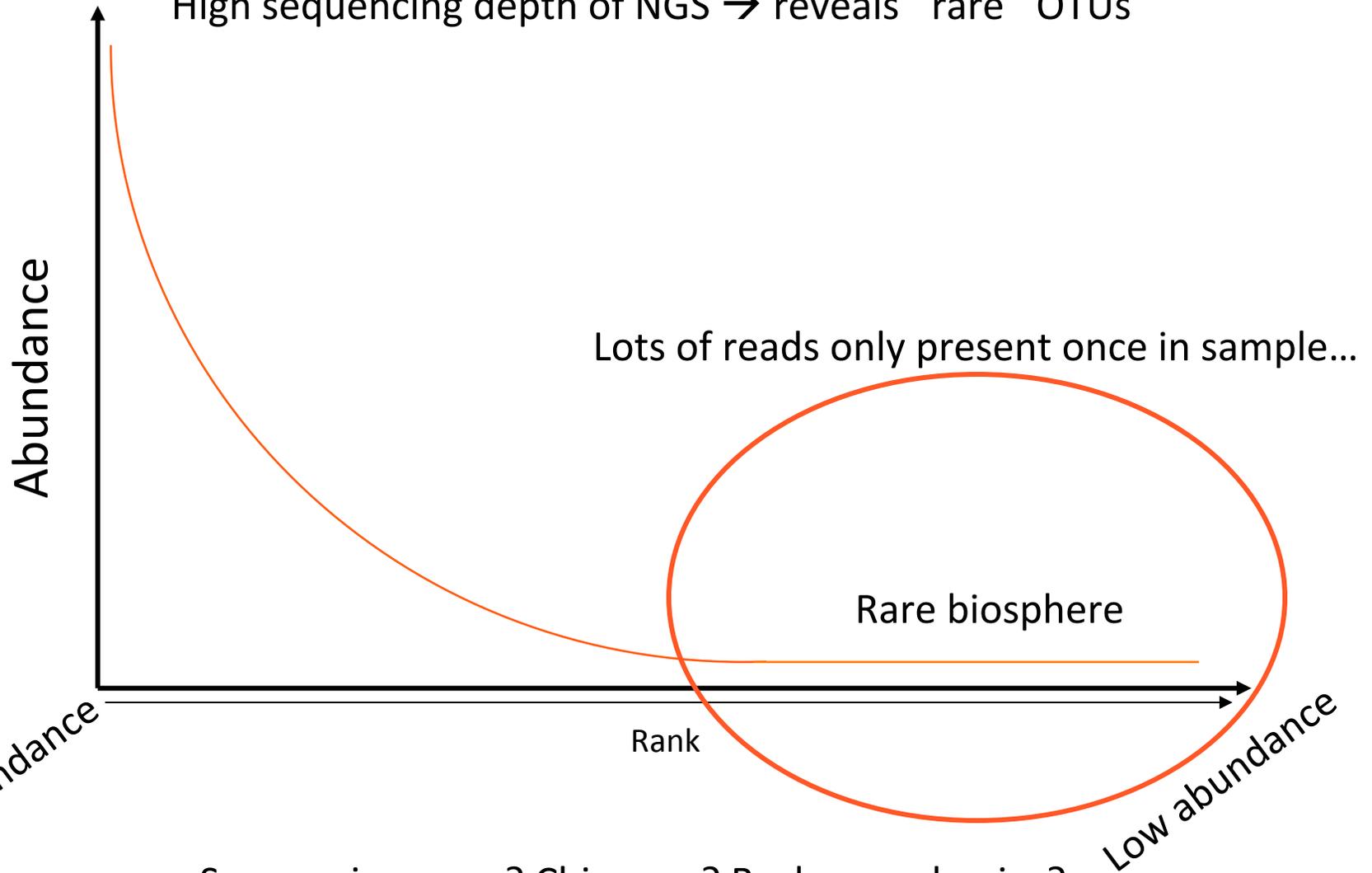
Reduce dataset by dereplication/clustering

**Why amplicon tags ?
Deeper sequencing,
cheaper, faster**

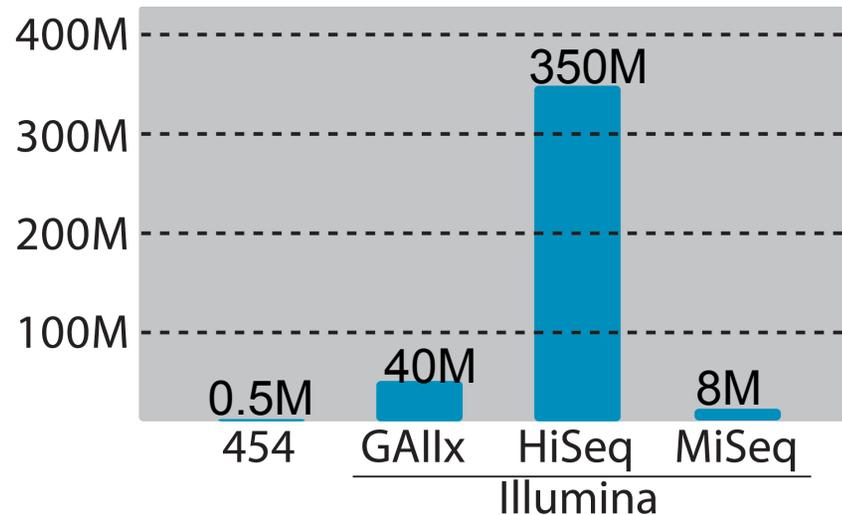
Identification
(BLAST, RDP classifier)



High sequencing depth of NGS → reveals “rare” OTUs

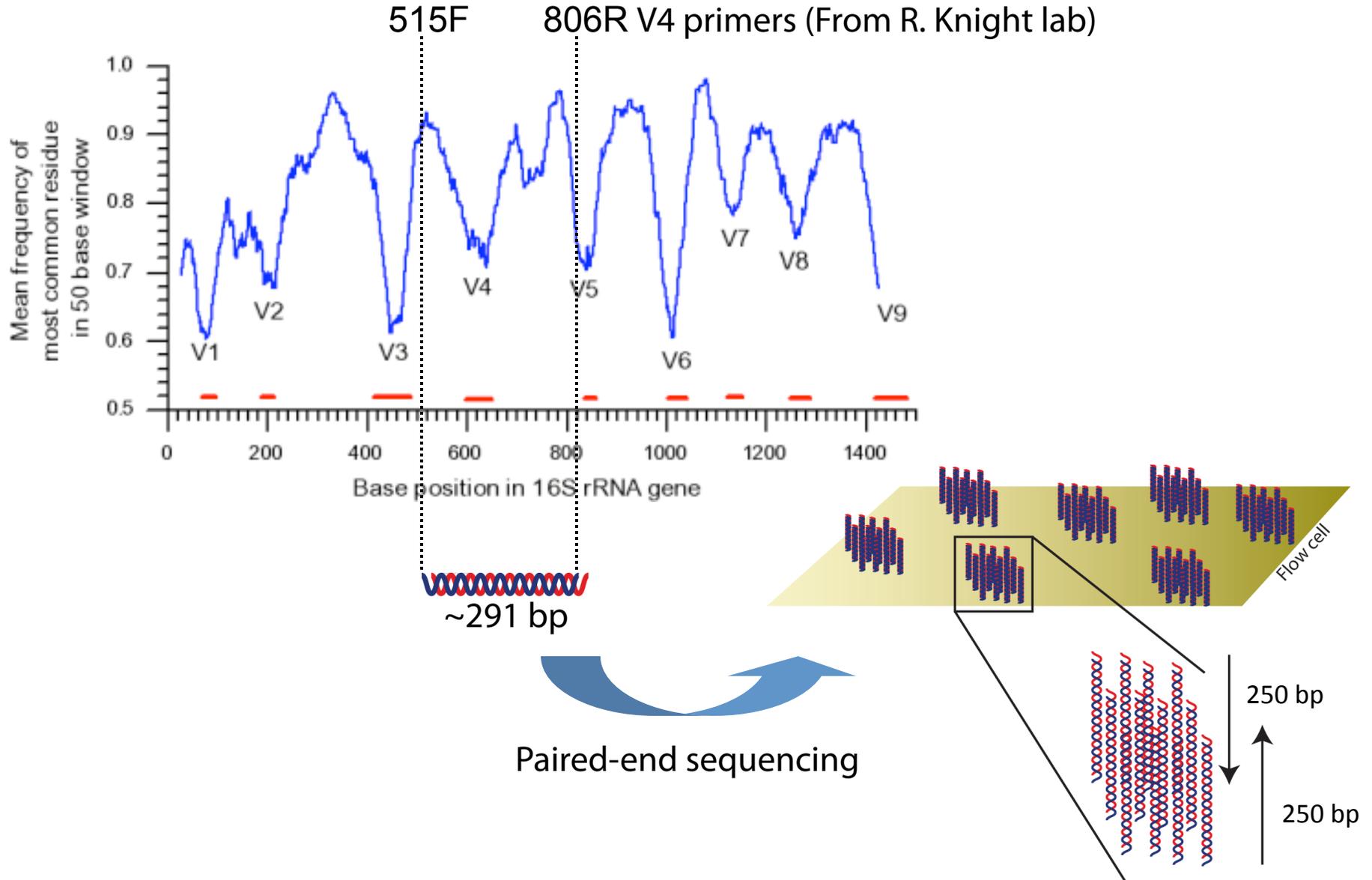


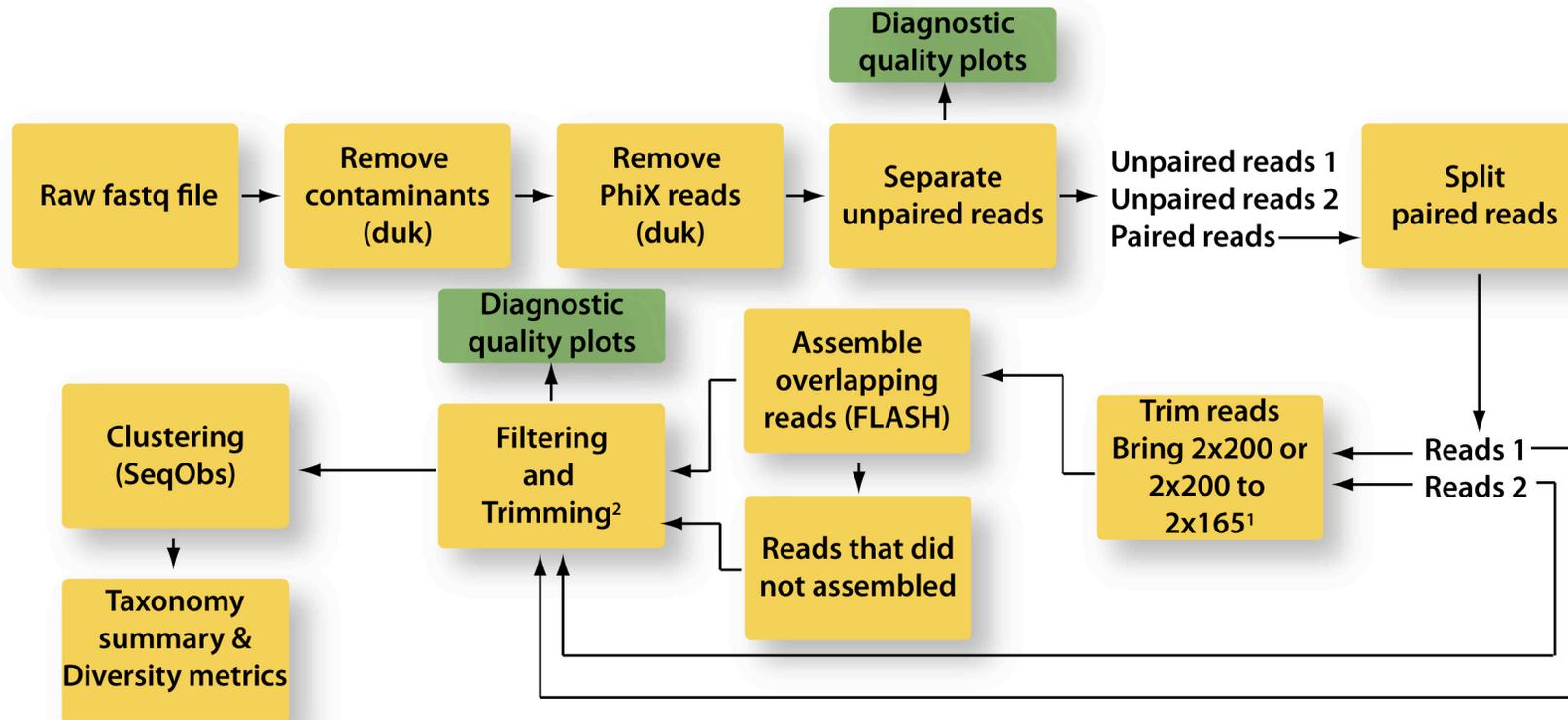
Platforms throughput



- **Move 16S tags sequencing from 454 to Illumina MiSeq platform**
 - \$ and higher quality reads with Illumina
 - HiSeq = Slow (18 days) and huge output compared to 454 (suitable for big projects 1000+ indexes(barcodes)/libraries)
 - MiSeq = moderately high throughput (More suitable for standard projects) and fast (~30 hrs run)
 - Longer reads with MiSeq (250 bp) compared to HiSeq (150 bp)
- ↑throughput → Clustering algorithm development.

Sequencing design for 16S tags on MiSeq



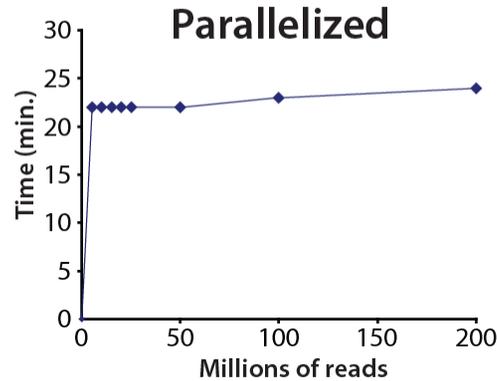


- **Remove low quality reads**
- **Reduce dataset**
- **Provide diversity metrics and taxonomic classification**

16S tags clustering



Sequences sorted by alphabetical order
Counts of each unique sequences (100% dereplication)

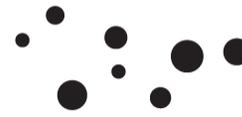


Dereplicated sequences binned by abundance:

High abundance



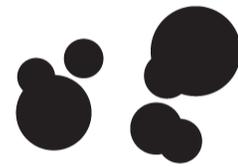
Low abundance



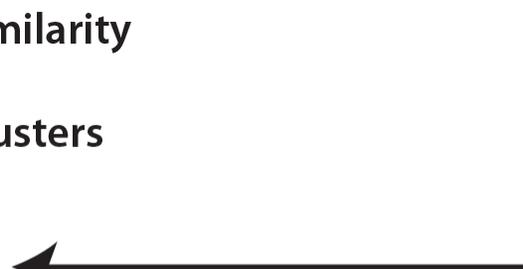
Clustered at 97% similarity



High abundance clusters



High abundant clusters to absorb low abundant sequences (97% similarity)



Final clusters





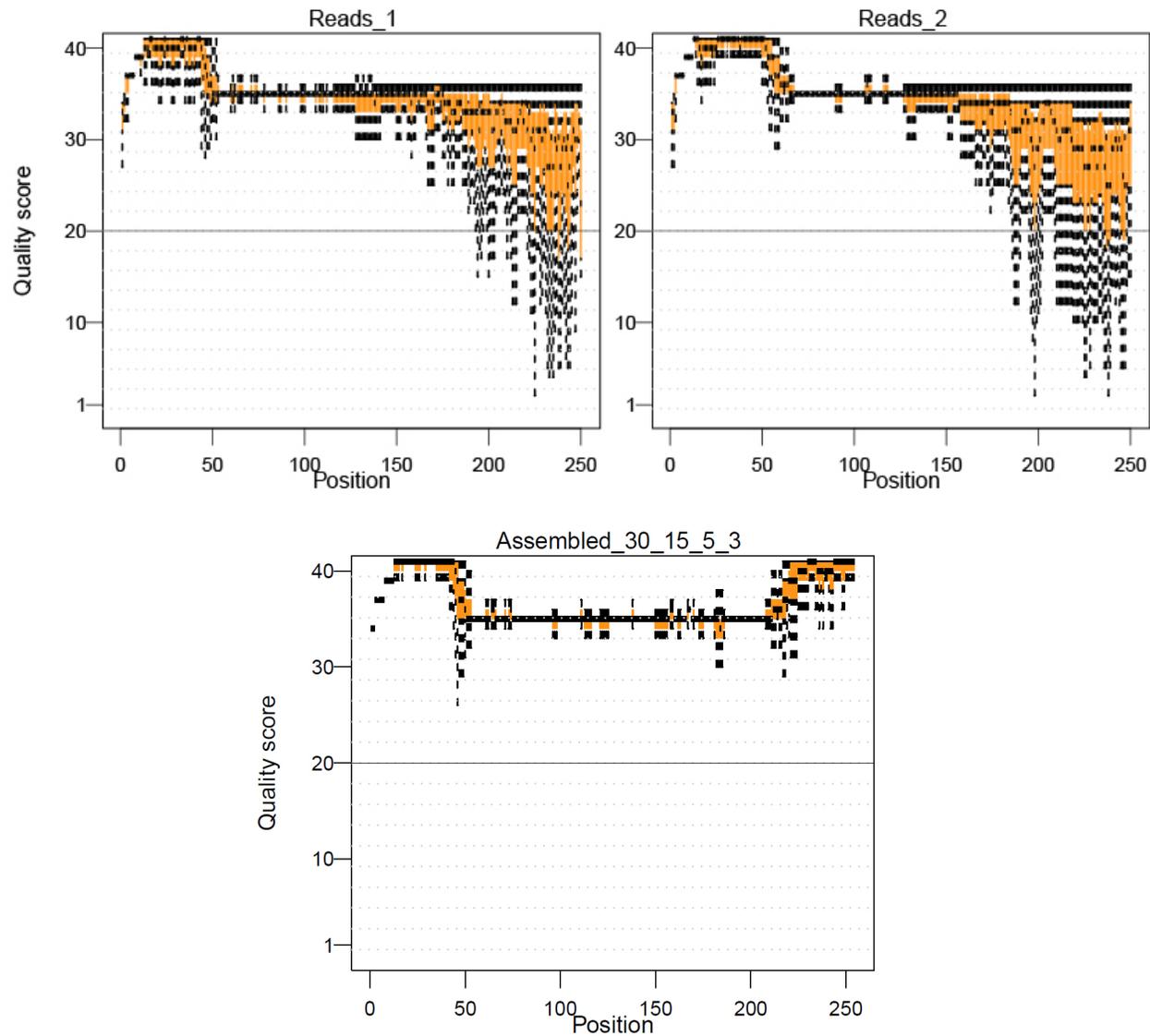
- MiSeq 96 barcodes typical run**

Total reads	7.6M
Contaminants reads	0.12M
PhiX reads	2M
Paired reads 1 (non-PhiX)	2.7M
Paired reads 2 (non-PhiX)	2.7M
Assembled reads	2.5M
Assembled reads QC passed	2.1M
Assembled reads / barcodes	22,574 ± 8,244

How is quality?

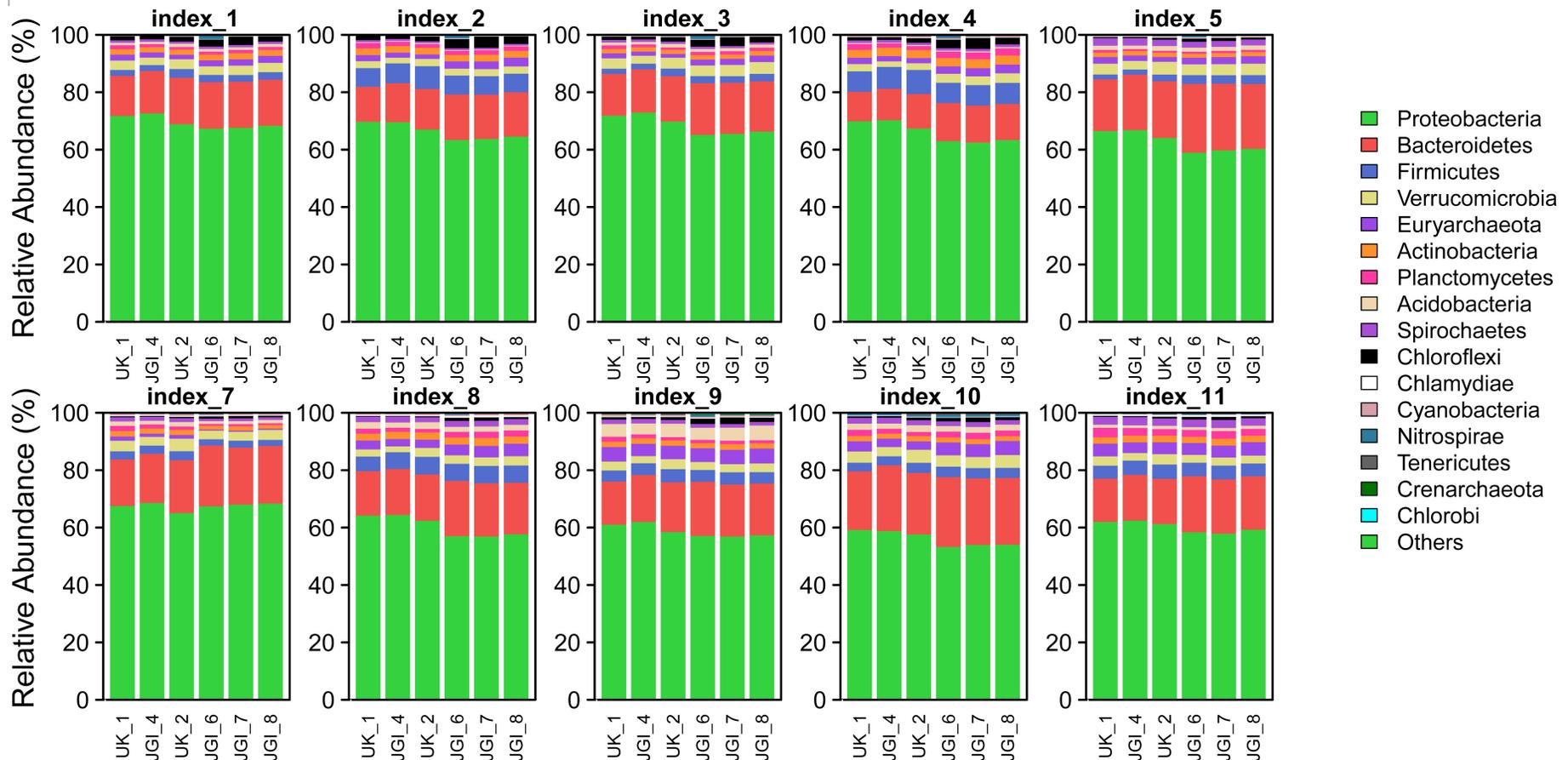


- **2x250 runs**



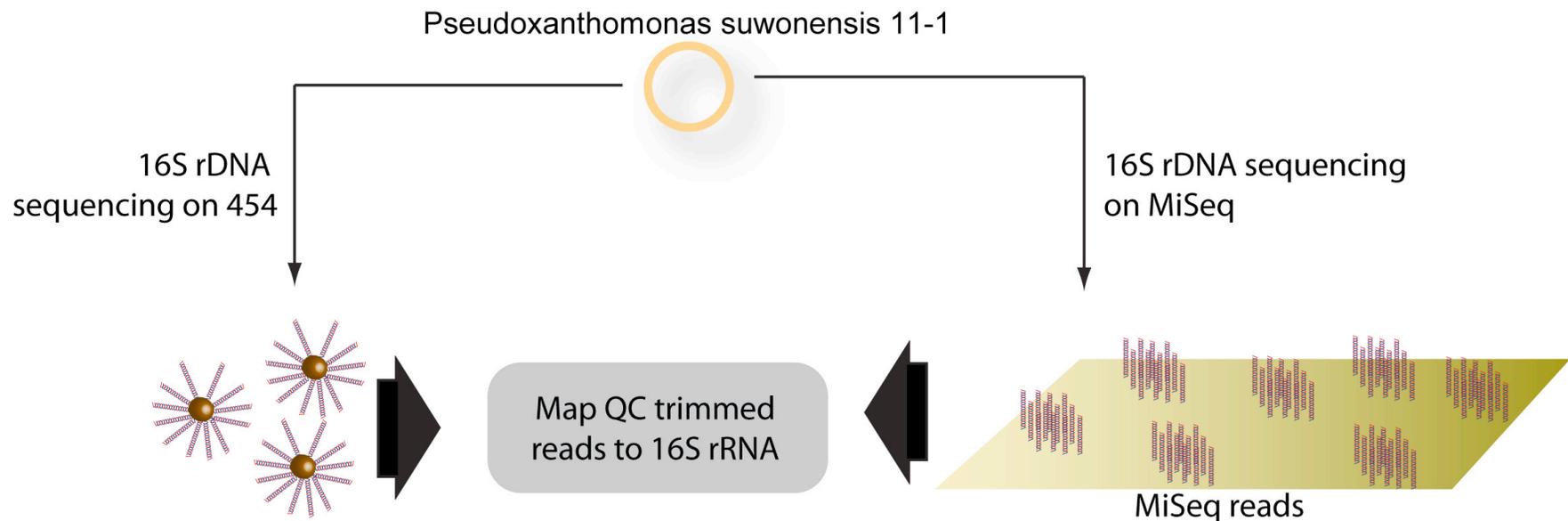


- Multiple runs with same sample
- Highly similar taxonomy

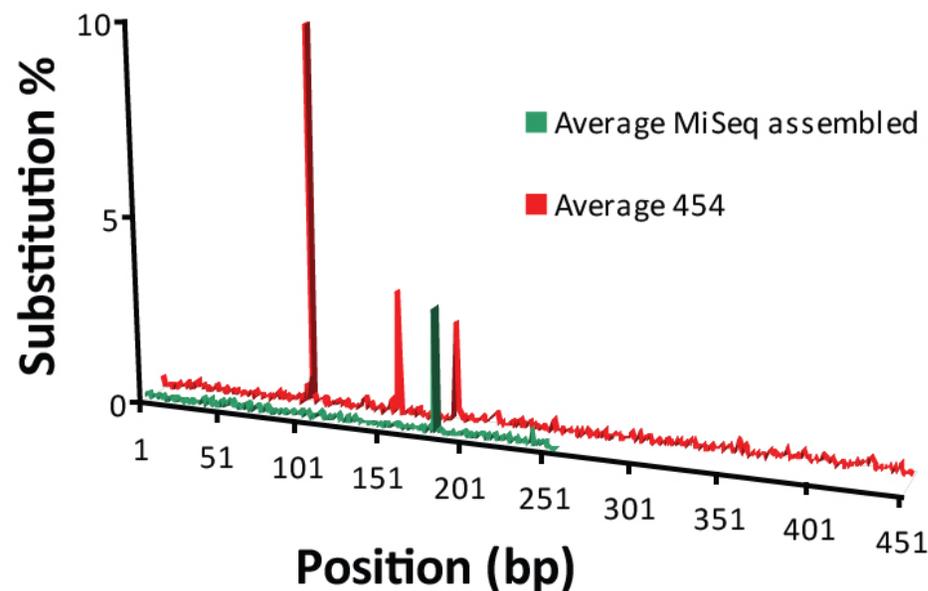
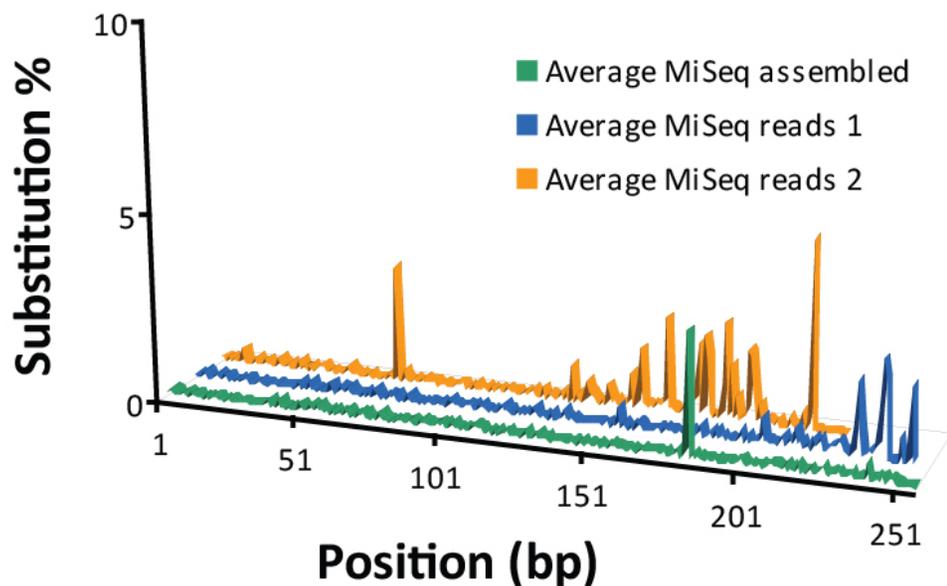




- **Is quality better on MiSeq?**
 - **Sequence 16S rDNA amplicons of a single microbial genome on MiSeq and 454**
 - **Map reads to 16S rDNA and estimate error**



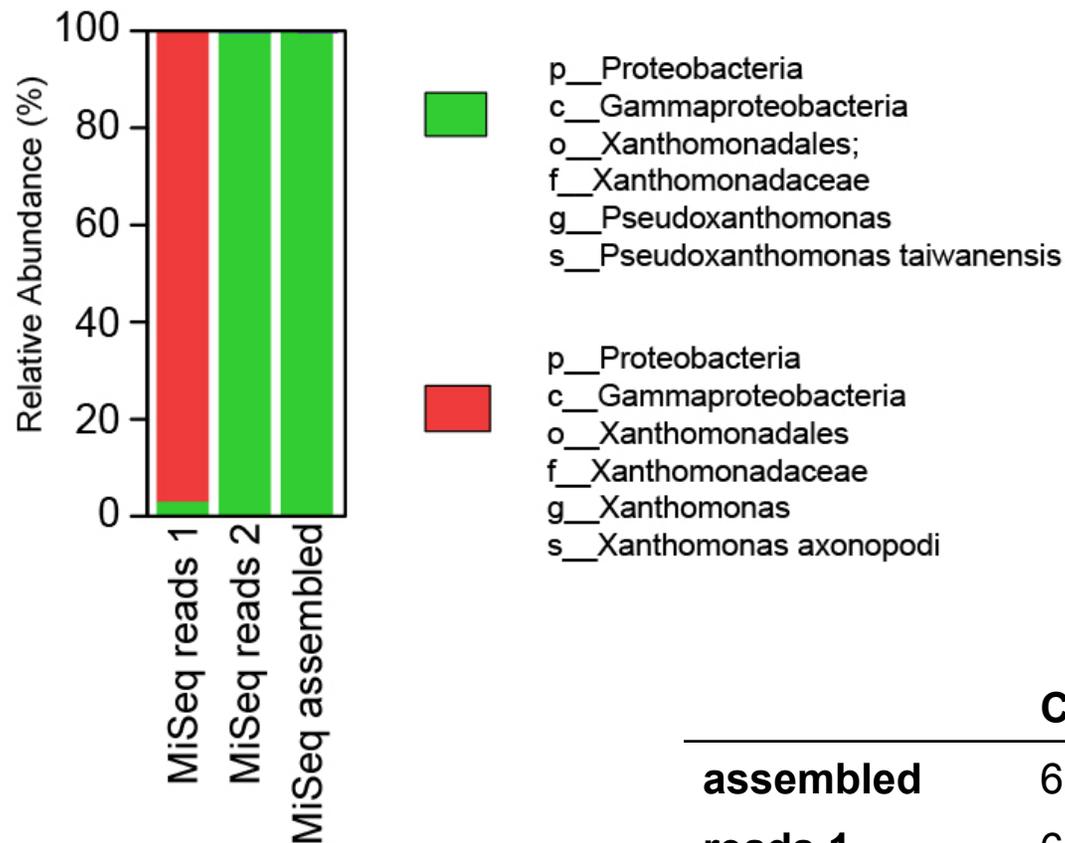
Each experiment done in triplicate



Error type (error/ 1,000,000 sequenced bases)	MiSeq reads 1	MiSeq reads 2	MiSeq assembled	454
Insertion rate	0.49	0.14	0.39	2,365
Deletion rate	82	62	13	2,816
Substitution rate	1,338	1,827	1,163	4,392



• Advantage of using assembled reads?



So why bother assembling reads?

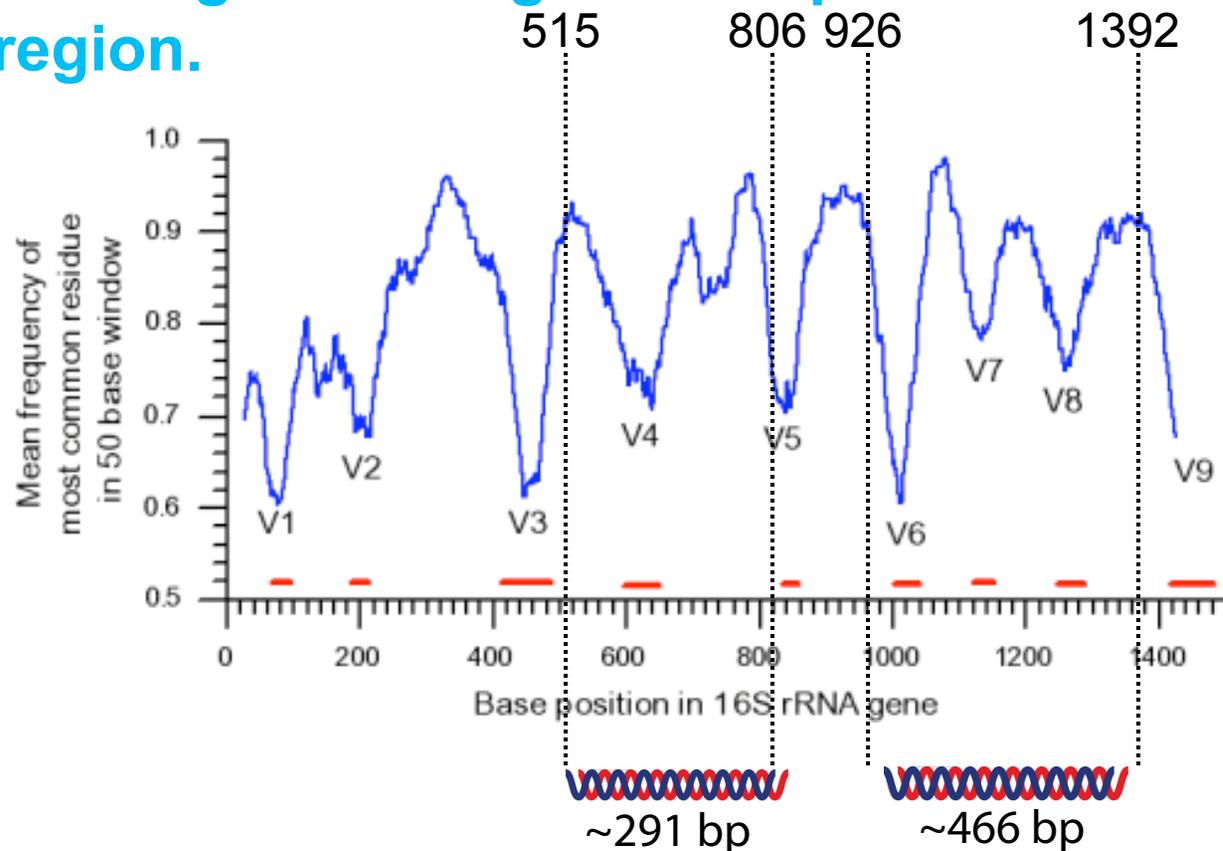
	Count	Average length
assembled	63,352	250 ± 1.3
reads 1	64,070	188 ± 36.2
reads 2	57,656	153 ± 24.3



- **MiSeq offers deeper sequencing than 454**
 - Better (assembled) reads quality on MiSeq
 - Lower error rate
 - Higher throughput
 - lower cost
 - Higher multiplexing
- **Cons**
 - Lower read length (Compensated by high quality of reads?)
 - Longer primers (more expensive)
 - Sequencing run is fast, but library preparation time is long.



- Compare runs of 454 and MiSeq of same sample
 - Although challenge to compare V4 with V6-8 region.

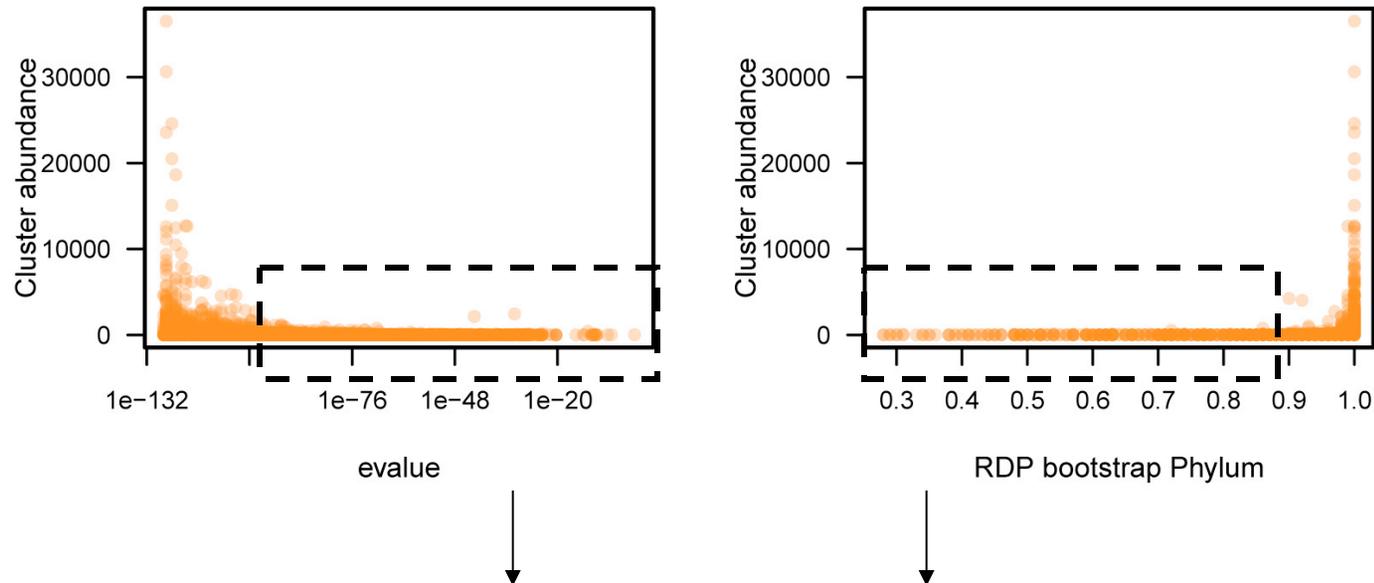




- **Susannah Tringe**
- **Feng Chen**
- **Kanwar Singh**
- **Ed Kirton**
- **Alison Fern**
- **Chris Daum**
- **Christine Naca**
- **And more!**



MiSeq wetlands samples



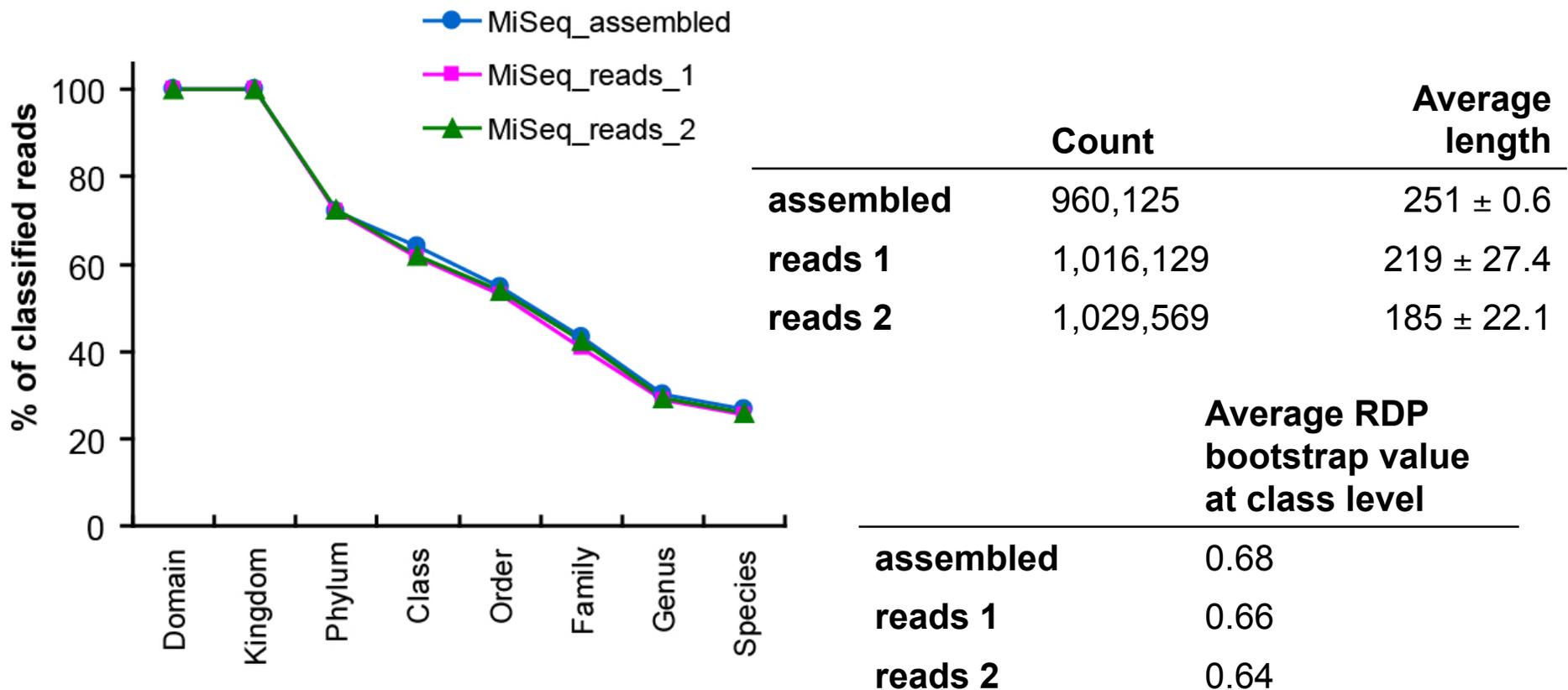
Low abundant reads consistently shows low confidence in Classification.

Low abundant reads = errors, artifacts?

Low abundant reads are underrepresented in databases?



- **Advantages of using assembled reads?**
 - **Test samples from a wetlands site (highly complex environment)**

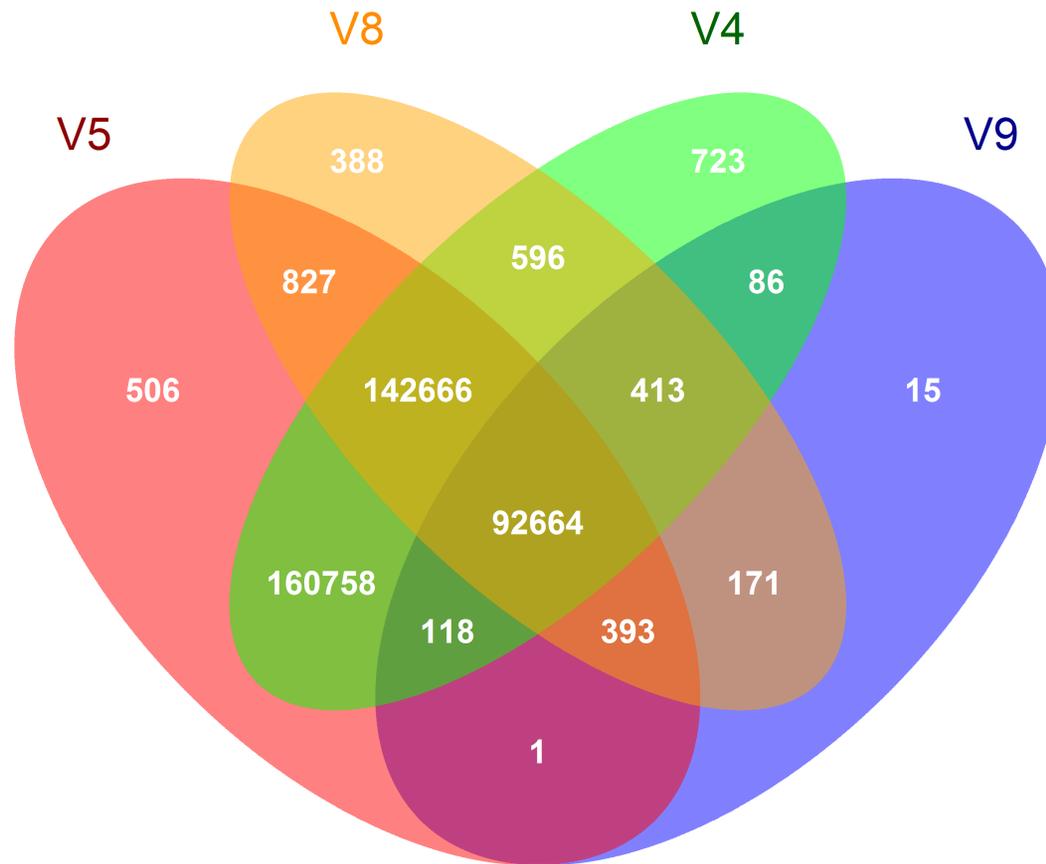




- **Short size of amplicon**
- **What filtering parameters to use (stringency level)?**
 - → **balance between stringency filter and keeping as much data as we can**
- **Whole new dimension for rare biosphere?**
- **Handling large numbers of sample (tens of thousand magnitude)**
- **Sequencing run is fast, but library preparation time is long.**
- **Cost of barcoded primers (will need lots of barcodes), handling.**
- **Huge ammount of samples → statistics models...**



- **Primer pair of variable region is likely to affect outcome of results.**



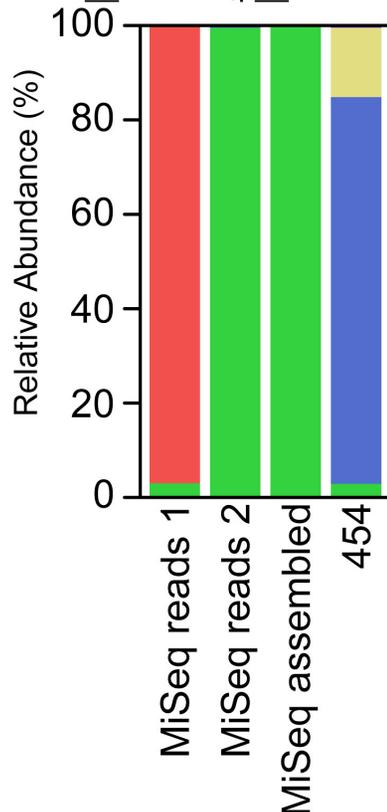
In silico PCR on 16S Greengenes database.



• Advantage of using assembled reads?

Raw reads → assembled → QC → clustering → classification with RDP classifier

- k__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Xanthomonadales;f__Xanthomonadaceae;g__Pseudoxanthomonas
- k__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Xanthomonadales;f__Xanthomonadaceae;g__Xanthomonas
- k__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Xanthomonadales;f__Xanthomonadaceae;
- k__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Xanthomonadales;f__Xanthomonadaceae;g__Lysobacter



	MiSeq
Cluster #1	21,106
Cluster #2	5
Cluster #3	4
Cluster #4	-
Cluster #5	-
Cluster #6	-
Cluster #7	-
Cluster #8	-

	454
Cluster #1	11,071
Cluster #2	1,756
Cluster #3	39
Cluster #4	37
Cluster #5	10
Cluster #6	8
Cluster #7	1
Cluster #8	1