

From .fastq to OTU: an overview

GC

TTCGAGGACTACTGGGGTATCTAATCCCGTTCGCTCCCTAGCTTCGC
CTTGATATCTACGCATTCACCGCTACACCAAGAATTCCAGTTGCCCT
GATTTCACAAAGTGACTTGTGCGCCGCCTACGCGCCCTTACGCCAG

FF:F,FF,FFF::FFFFF:F:FFFFF:FFFF:F,FF,FFFFFFFFFFFF:
FFFFFFFFFFFF:F:FFFFF,F:FFF:FF:::F,,FF:FFFFFF:F
F:FFF,:F,FFF:FFF:FFF,:FFF:FFF:FFFFFFF:,F,F:F
1010:95:HLWN2DRXX:1:2101:1090:1000 2:N:0:1
ATACCAGGACTACAGGGGTTCTAATCCCGTTCGCTCCCTGGCTTCG
TCTCGATATCTACGCATTCACCGCTACACCGAGAATTCCAGTAGCCCC
AGATTTCACAAAGTGACTTGTGCGCCGCCTACGCGCCCTTACGCCA

FF,F:FFF,FF:FFFFF:FFFFFFFFFFFFF:F:FFFFFFF
F,FFFFFFF:FFFFF:FFFFFFFFFFFF:FFFFF,FFFFF:FFFFFFF
FFFFF:FFFFFFFFFFFF,FFF,FFFFF:FFFFFFF,FFFF::FFF,
1010:95:HLWN2DRXX:1:2101:1127:1000 2:N:0:1

TTGGAGATGCGTTCTTCATCTATTCCAGAACGAAAGAGATCC
TTTGAAGTGTGACCCGTTGCACACGACCGAACGGAGGGTC
GCGTTGTTATTGTTCACTCCGCAGGTTCTCCTACGGTTA

	B1	B2	B3	B4	B5
AGTTGA	6	0	13	10	15
GTGGAC	4	3	7	2	15
CTGTAA	3	3	0	10	12

Many pipelines exist

- QIIME
- QIIME2 (uses DADA2)
- Usearch
- Vsearch
- DADA2
- PEAR
- And many more

Overview of steps within DADA2 pipeline

Library prep and sequencing

Demultiplexing

Check for remaining primers

Filter and Trim: `fastqFilter()` or `fastqPairedFilter()`

Infer sample composition: `dada()`

Merge paired reads: `mergePairs()`

Dereplicate: `derepFastq()`

Make sequence table: `makeSequenceTable()`

Remove chimeras: `isBimeraDenovo()` or `removeBimeraDenovo()`

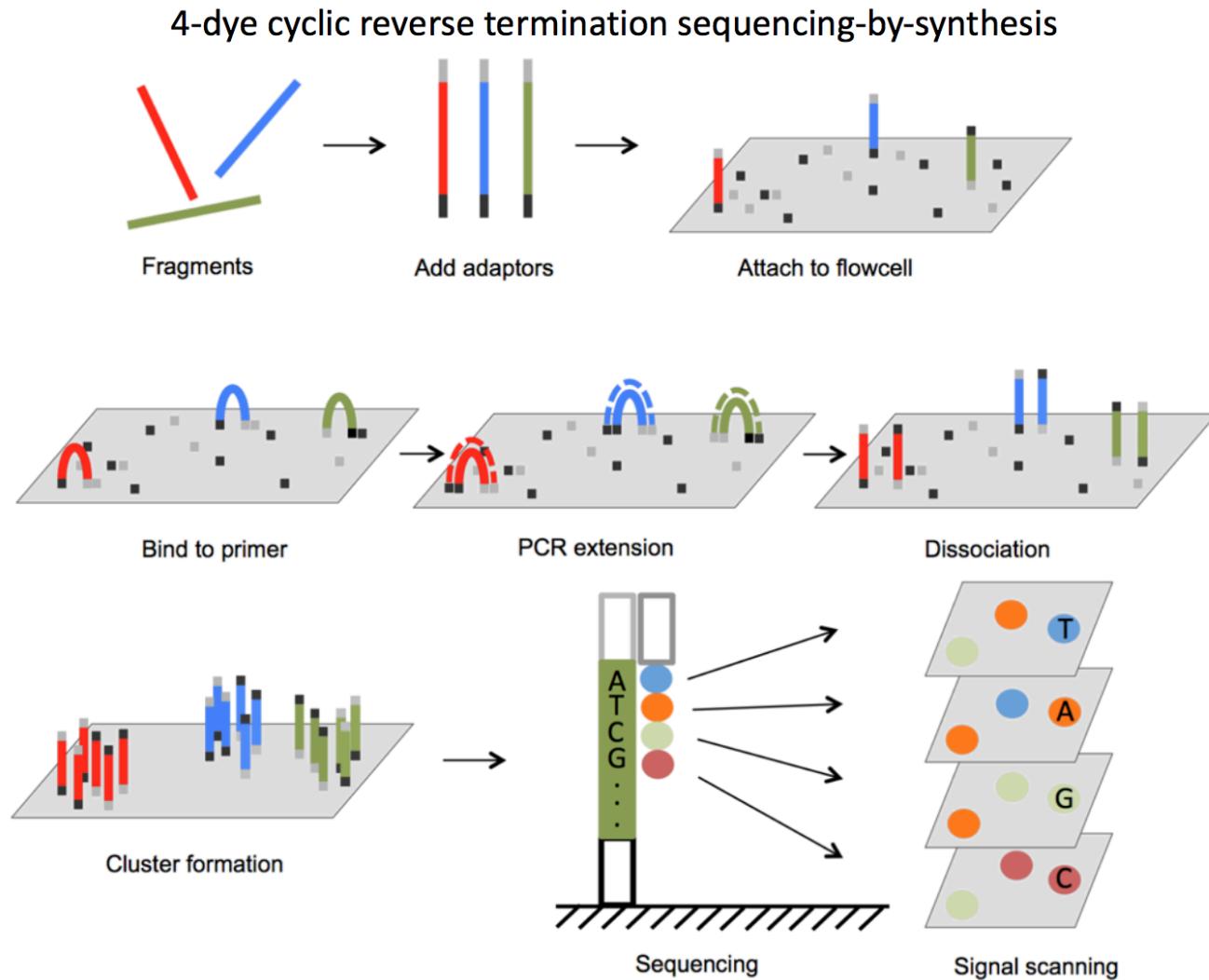
Assign taxonomy

Library Prep

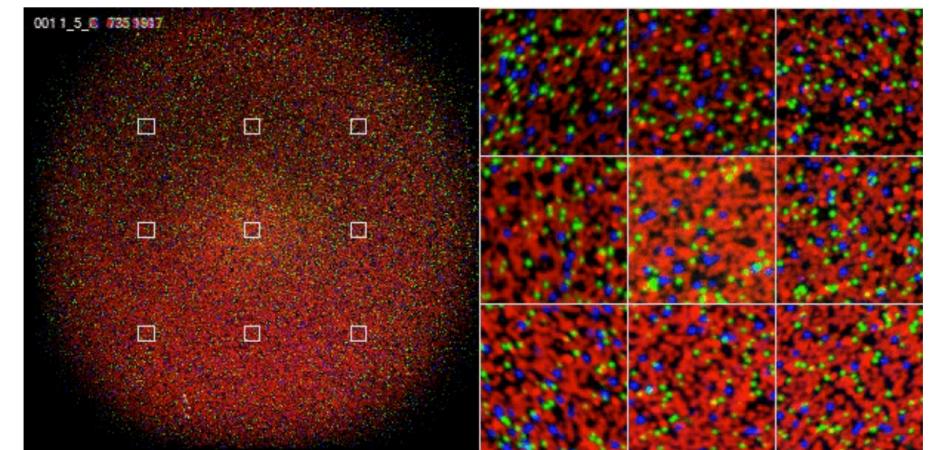


Sequencing

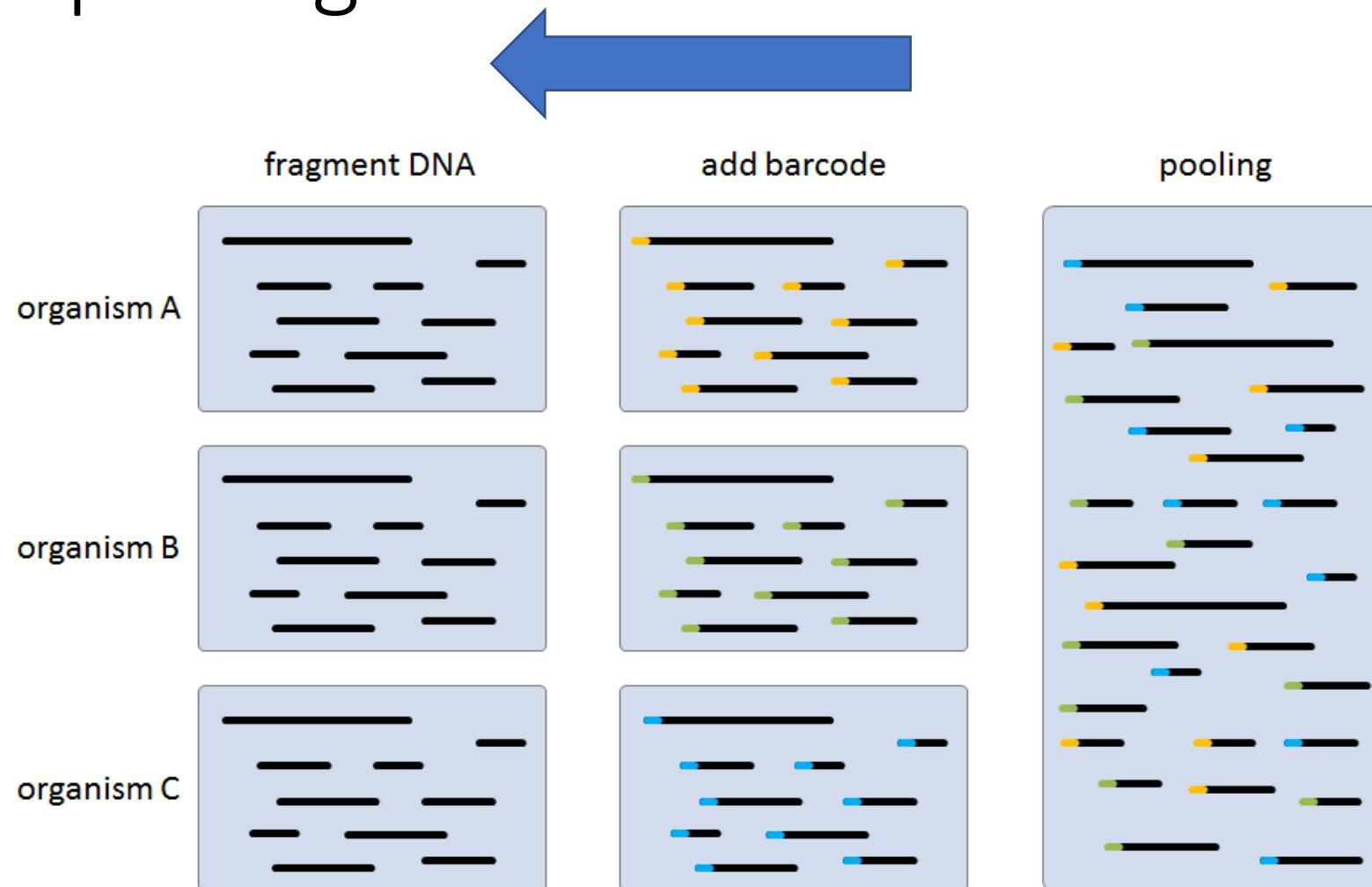
Illumina sequencing: how it works



Amplicon sequencing run on Illumina MiSeq



Demultiplexing



Removing non-biological nucleotides

- If not removed, non-biological nucleotides can cause problems with OTU picking
- Non-specific binding of primers can mask actual sequence of data

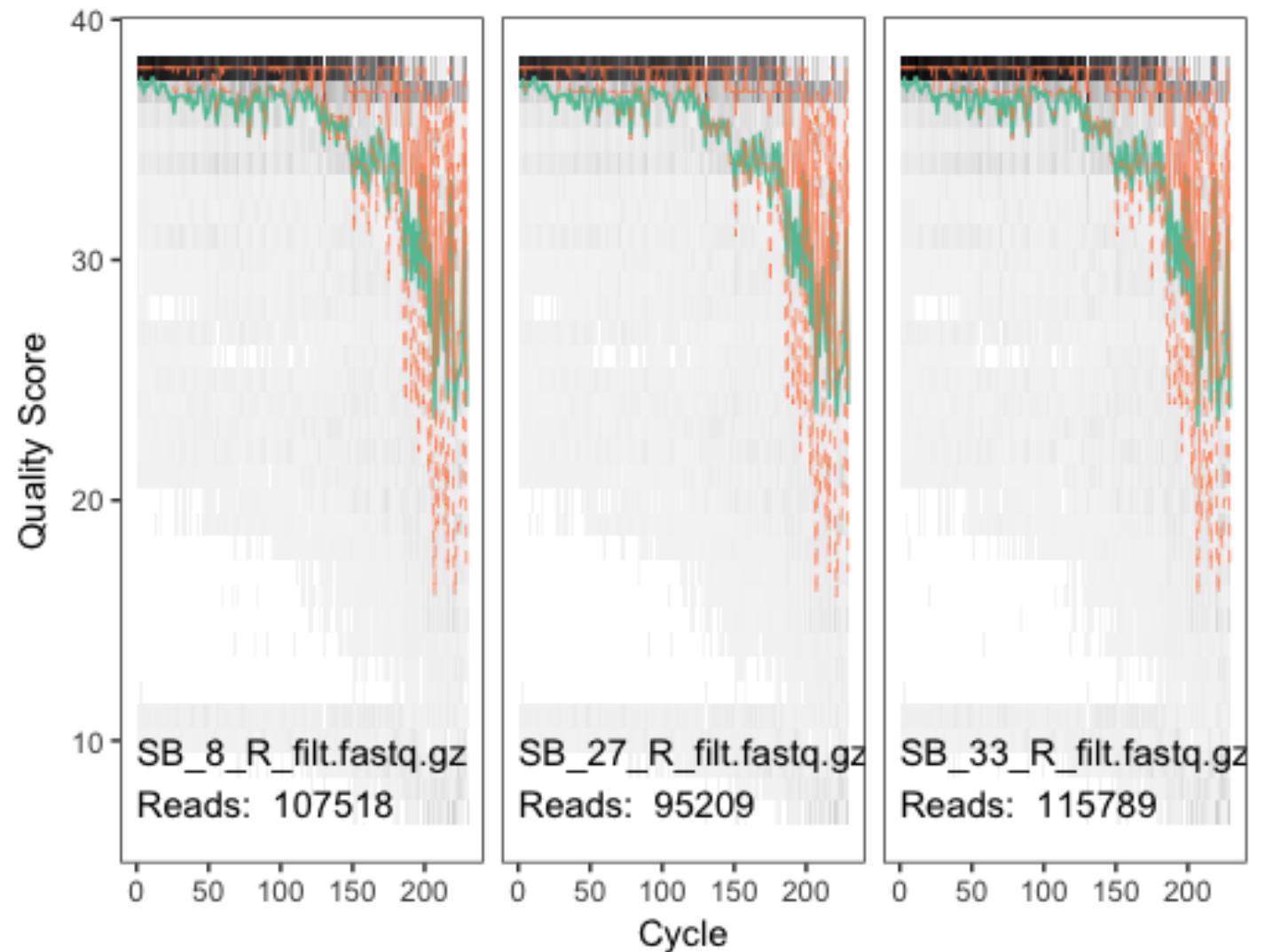
```
65 - `}`{r}  
66 FWD <- "GTGYCAGCMGCCGCGTAA"  
67 REV <- "GGACTACNVGGGTWTCTAAT"  
68 ``"
```

ATTCTGA ✓
TAAGCT

ATTCTGA x
TATGCT

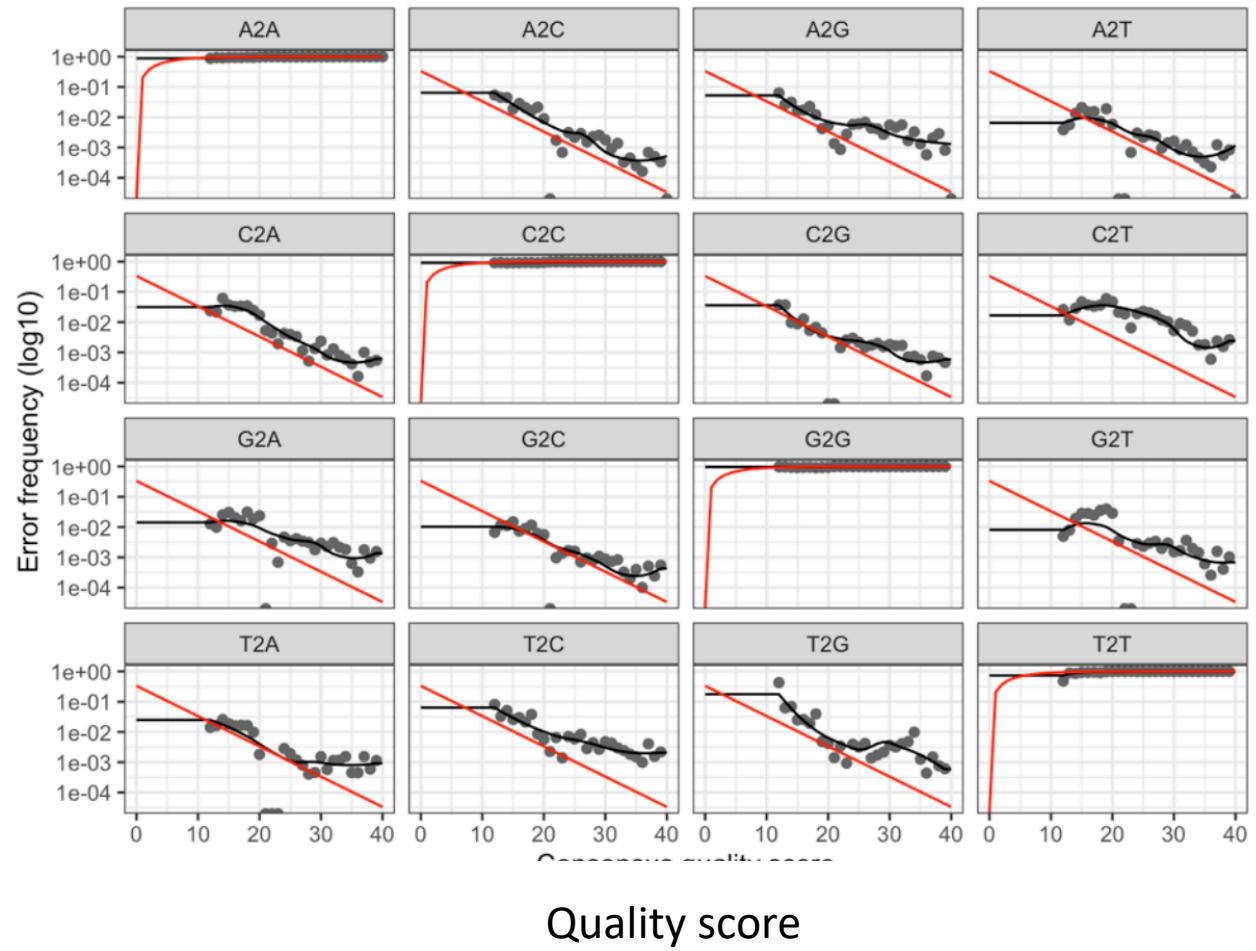
Quality filtering

- Different ways to accomplish this...
 - Minimum quality score required for inclusion
 - Minimum or maximum length



Infer sample composition

- Plotting probability of errors
- "Fixes" sequencing errors
 - # of occurrences of each sequence



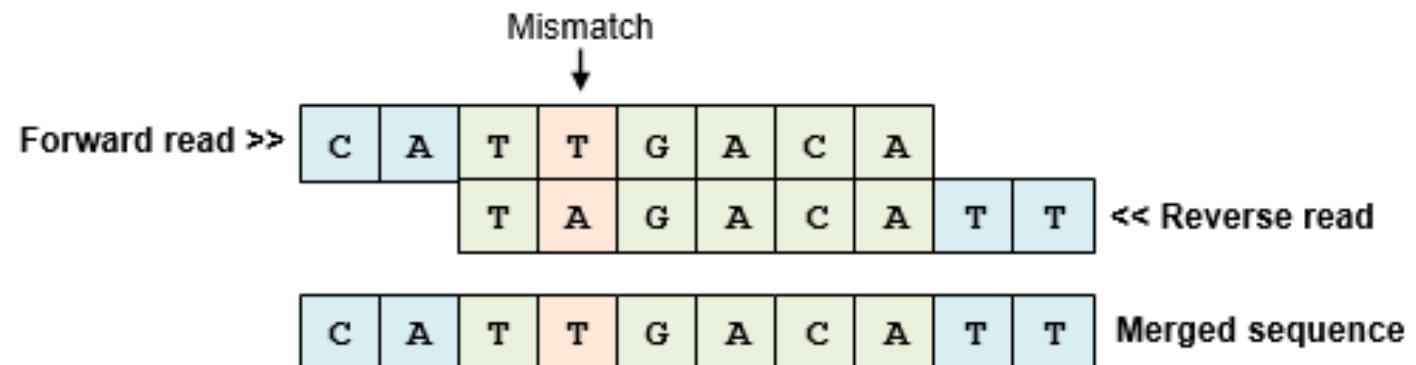
Merge reads

Why merge?

Serial No. of dollar bill

Considerations:

- Minimum number of overlapping base pairs
- Maximum number of mismatches
- What if you can't get your reads to merge?



Dereplication

- During dereplication, we condense the data by collapsing together all reads that encode the same sequence.
- This step occurs in all pipelines, but DADA2 retains a summary of quality info here.
- Before there may be 10000s of repeats. After dereplication, we condense our dataset to a single representative sequence with counts.
 - Significantly reduces computation time

	AGT TGA
B1	0
B2	3
B3	3
B4	0

	GTG GAC
B1	0
B2	3
B3	0
B4	3

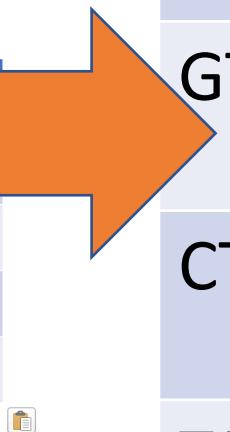
	CTGT AA	AATA GT
B1	0	0
B2	3	3
B3	3	3
B4	0	0

Making sequence table

	AGT TGA
B1	0
B2	3
B3	3
B4	0

	GTG GAC
B1	0
B2	3
B3	0
B4	3

	CTGT AA
B1	0
B2	3
B3	3
B4	0



	B1	B2	B3	B4	B5
AGTTGA	6	0	13	10	15
GTGGAC	4	3	7	2	15
CTGTAA	3	3	0	10	12
TAAAGT	1	0	1	0	1

But we're not quite done...

Chimeras!

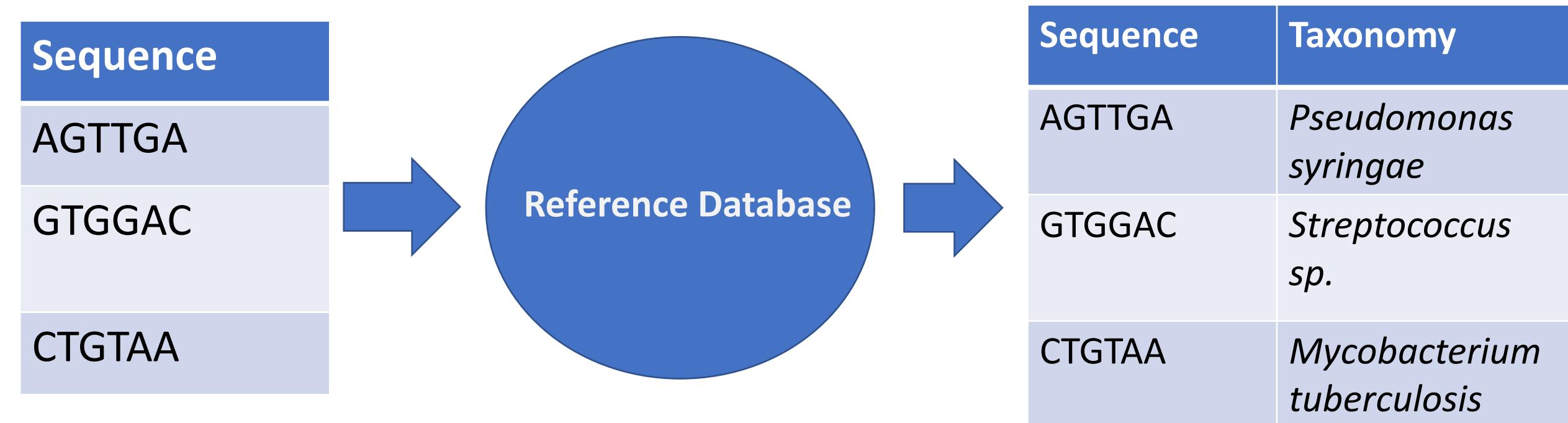
- What is a chimera?
- Why do we care about this in sequencing?



Chimera removal

	B1	B2	B3	B4	B5
AGT TGA	6	0	13	10	15
GTGGAC	4	3	7	2	15
CTG TAA	3	3	0	10	12
X TAA AGT	1	0	1	0	1

Assigning taxonomy



	B1	B2	B3	B4	B5	Taxonomy
AGTTGA	6	0	13	10	15	<i>Pseudomonas syringae</i>
GTGGAC	4	3	7	2	15	<i>Streptococcus sp.</i>
CTGTAA	3	3	0	10	12	<i>Mycobacterium tuberculosis</i>

DADA2 vignette

[https://www.bioconductor.org/packages/devel/bioc/vignettes/dada2/inst/doc/dada2-
intro.html#derePLICATE](https://www.bioconductor.org/packages/devel/bioc/vignettes/dada2/inst/doc/dada2-intro.html#derePLICATE)