

# Bioinformatics Pipeline

Yisha Tan<sup>\*1,2,3</sup>

<sup>1</sup>Archaeal Biology Centre, Synthetic Biology Research Center, Shenzhen Key Laboratory of Marine Microbiome Engineering, Key Laboratory of Marine Microbiome Engineering of Guangdong Higher Education Institutes, Institute for Advanced Study, Shenzhen University, Shenzhen 518060, China

<sup>2</sup>HKUST Shenzhen-Hong Kong Collaborative Innovation Research Institute, Futian, Shenzhen, China

<sup>3</sup>Department of Ocean Science, The Hong Kong University of Science and Technology, Hong Kong, SAR China

## Introduction

This is the bioinformatics pipeline used for the manuscript titled: *“Genetic diversity of cyanobacteria and their viruses reveals features of coevolution in cultured and natural environments.”*

We focus on analyzing metagenomic data to detect and characterize cyanobacteria and associated cyanophages, especially in estuarine environments. The major steps are:

1. Metagenomic assembly, gene prediction, binning, and taxonomy.
2. Viral prediction, clustering, taxonomic/functional annotation.
3. Whole-genome analyses such as mapping, variant calling, and phylogenetics.

Below is a detailed explanation of each step, sample commands, and usage notes.

## 1 Metagenomic Assembly and Genome Binning

### 1.1 Trimming Sequences

```
1 sickle pe \  
2   -f raw_data1.fastq \  
3   -r raw_data2.fastq \  
4   -t sanger \  
5   -q 25 \  
6   -o raw_data1.trim.fastq \  
7   -p raw_data2.trim.fastq
```

---

\*ytanap@connect.ust.hk

## 1.2 Assembly

```
1 ~/SeqTools/interleave.pl \  
2     -fwd raw_data1.trim.fastq \  
3     -rev raw_data2.trim.fastq \  
4     -o total.trim.fa  
5  
6 idba_ud \  
7     --mink 65 \  
8     --maxk 145 \  
9     --step 10 \  
10    -r total.trim.fa \  
11    -o k65-145
```

## 1.3 Gene Prediction

```
1 prodigal \  
2     -f gff \  
3     -i sample.fa \  
4     -o sample.gff \  
5     -a prot.faa \  
6     -p meta  
7  
8 python2 emapper.py \  
9     -i prot.faa \  
10    -o prot.txt \  
11    -m diamond \  
12    --cpu 20
```

## 1.4 Mapping and Binning

```
1 # Build index  
2 ~/bowtie2-2.4.2/bowtie2-build sample.fa sample.fa  
3  
4 # Map reads  
5 ~/bowtie2-2.4.2/bowtie2 \  
6     -x sample.fa \  
7     -1 raw_data1.trim.fastq \  
8     -2 raw_data2.trim.fastq \  
9     -S sample.sam  
10  
11 # Convert & sort  
12 samtools view -bS sample.sam > sample.bam  
13 samtools sort sample.bam -o sample_sort.bam  
14 samtools index sample_sort.bam
```

```

15 samtools flagstat sample_sort.bam > sample_flag.txt
16
17 # Coverage
18 ~/jgi_summarize_bam_contig_depths --outputDepth sample_depth.txt
   sample_sort.bam
19
20 # Binning
21 metabat \
22     -i sample.fa \
23     -a sample_depth.txt \
24     -o Binning

```

## 1.5 Quality Assessment

```

1 checkm lineage_wf \
2     -x fa \
3     -t 32 \
4     checkm_result

```

## 1.6 Taxonomy

```

1 gtdbtk classify_wf \
2     --genome_dir genomes \
3     --out_dir output \
4     --cpus 10
5
6 blastn \
7     -db ~/SILVA/SILVA_132_SSURef_Nr99_tax_silva.fasta \
8     -query bin.fasta \
9     -out blastn_silva_result.txt \
10    -evaluate 1e-5 \
11    -outfmt 6 \
12    -max_target_seqs 1 \
13    -num_threads 32

```

## 1.7 Abundance

```

1 coverm contig \
2     -c raw_data1.trim.fastq raw_data2.trim.fastq \
3     --reference cyano_silva138.fasta \
4     -m rpkm \
5     --output-file genes.rpkm.txt \
6     -t 60

```

## 2 Viral Prediction and Analysis

### 2.1 Recovery of Viral Contigs

```
1 virsorter run \  
2   -w vs2 \  
3   -j 60 \  
4   --min-length 1500 \  
5   -i sample.fa \  
6   -d vs2_databse all  
7  
8 python dvf.py \  
9   -i sample.vir.fa \  
10  -l 1000 \  
11  -c 2  
12  
13 checkv end_to_end \  
14   sample.vir.fna \  
15   output_checkv \  
16   -t 16
```

### 2.2 Cyanophage Genomes and Clustering

```
1 cdhit \  
2   -i sample.vir.faa \  
3   -o output_cdhit \  
4   -c 0.95 \  
5   -aS 0.85 \  
6   -n 9 \  
7   -d 0 \  
8   -T 8
```

### 2.3 Viral Taxonomy and Genome Clustering

```
1 python ~/vRhyme/aux/coverage_table_convert.py \  
2   -i sample.vir.rpkm.txt \  
3   -o vir.coverage.tsv  
4  
5 vRhyme -i sample.vir.fasta -c vir.coverage.tsv  
6  
7 vcontact2_gene2genome \  
8   --proteins all_cluster.vir.faa \  
9   --output all_cluster.vir.g2g.csv \  
10  --source-type Prodigal-FAA  
11
```

```

12 vcontact \
13     --rel-mode Diamond \
14     --pcs-mode MCL \
15     --vcs-mode ClusterONE \
16     --cl-bin ~/Tools/ \
17     --db cyanophage_db \
18     --verbose \
19     --threads 30 \
20     --raw-proteins all_cluster.vir.faa \
21     --proteins-fp all_cluster.vir.g2g.csv \
22     --output-dir result

```

## 2.4 Gene Prediction and Function

```

1 prodigal \
2     -p meta \
3     -a protein_seq.fasta \
4     -m \
5     -d nucleotide_seq.fasta \
6     -o genome.gff \
7     -f gff \
8     -s stat \
9     -i genome.fasta
10
11 python emapper.py \
12     -i genome.faa \
13     -o genome.txt \
14     -m diamond \
15     --cpu 20

```

## 2.5 Viral Abundance Mapping

```

1 makeblastdb \
2     -in genome.faa \
3     -dbtype prot
4
5 blastp \
6     -query bin.faa \
7     -db genome.faa \
8     -out genome.txt \
9     -outfmt 6 \
10     -evalue 1e-3 \
11     -max_target_seqs 1 \
12     -num_threads 30
13

```

```

14 coverm contig \
15     -c raw_data1.trim.fastq -2 raw_data2.trim.fastq \
16     --reference vir.genome.fa \
17     -m rpkm \
18     --output-file genome.rpkm.txt \
19     -t 60

```

## 2.6 ANI/AAI

```

1 comparem aai_wf \
2     --cpus 40 \
3     vir.fasta \
4     ANI

```

## 2.7 Phylogenetic Trees

```

1 trimal \
2     -in combined_genome.fasta \
3     -out combined_genome.trim.fasta \
4     -gt 0.95
5
6 muscle \
7     -in combined_genome.trim.fasta \
8     -out combined_genome.trim.aln.fasta \
9     -maxiters 1 \
10    -diags -sv \
11    -distance1 kbit20_3
12
13 iqtree \
14     -s combined_genome.trim.aln.fasta \
15     -B 1000 \
16     -nt AUTO \
17     -m MFP \
18     -pre trim

```

## 2.8 BLAST Searches

```

1 makeblastdb \
2     -in cyanophage_ref.fa \
3     -dbtype nucl \
4     -out cyanophage_ref
5
6 nohup blastn \
7     -max_target_seqs 1 \

```

```

8  -db cyanophage_ref \
9  -out ref.cyanophage.blast \
10 -query lab.cyanophage.fa \
11 -num_threads 10 \
12 -outfmt 6 qseqid qlen qstart qend salltitles sseqid slen
    sstart send qcovs bitscore evalue pident

```

## 3 Whole-Genome Analyses

### 3.1 Alignment

```

1  ~/bowtie2-2.4.2/bowtie2 \
2  -x sample.fa \
3  -1 raw_data1.trim.fastq \
4  -2 raw_data2.trim.fastq \
5  -S sample.sam \
6  -p 10

```

### 3.2 Deduplicate Sequences

```

1  java -jar ~/Picard/picard.jar SortSam \
2  -I sample.sam \
3  -O sample_sort.sam \
4  -SO coordinate
5
6  java -jar ~/Picard/picard.jar MarkDuplicates \
7  -I sample_sort.sam \
8  -O sample_mdup.bam \
9  -M sample_mdup.metrics
10
11 java -jar ~/Picard/picard.jar AddOrReplaceReadGroups \
12 -I sample_mdup.bam \
13 -O sample_mdup_rg.bam \
14 -SO coordinate \
15 -LB sample \
16 -PL illumina \
17 -PU sample \
18 -SM sample
19
20 ~/samtools-1.12/samtools index sample_mdup_rg.bam

```

### 3.3 Mutation and Functional Analyses

```
1 ~/gatk-4.2.0.0/gatk HaplotypeCaller \  
2   -I sample_mdup_rg.bam \  
3   -O sample.g.vcf \  
4   -R ref.fa \  
5   --sample-ploidy 1 \  
6   -ERC GVCF
```

## Summary

This pipeline provides a comprehensive approach for:

- **Preprocessing & Assembly:** Trimming (Sickle), assembly (IDBA-UD), gene prediction (Prodigal).
- **Binning & Taxonomy:** Coverage mapping (Bowtie2, Samtools), binning (MetaBAT), and taxonomy (CheckM, GTDB-Tk).
- **Viral Discovery:** Identifying viral contigs (VirSorter, DeepVirFinder), evaluating completeness (CheckV), clustering and binning (CD-HIT, vRhyme, vContact2).
- **Abundance & Annotation:** RPKM calculations (CoverM), functional annotation (EggNOG-mapper).
- **Phylogenetic & Variant Analyses:** Tree building (trimal, MUSCLE, IQ-TREE), variant calling (GATK).

For questions or troubleshooting, please contact **Yisha Tan** at [ytanap@connect.ust.hk](mailto:ytanap@connect.ust.hk).