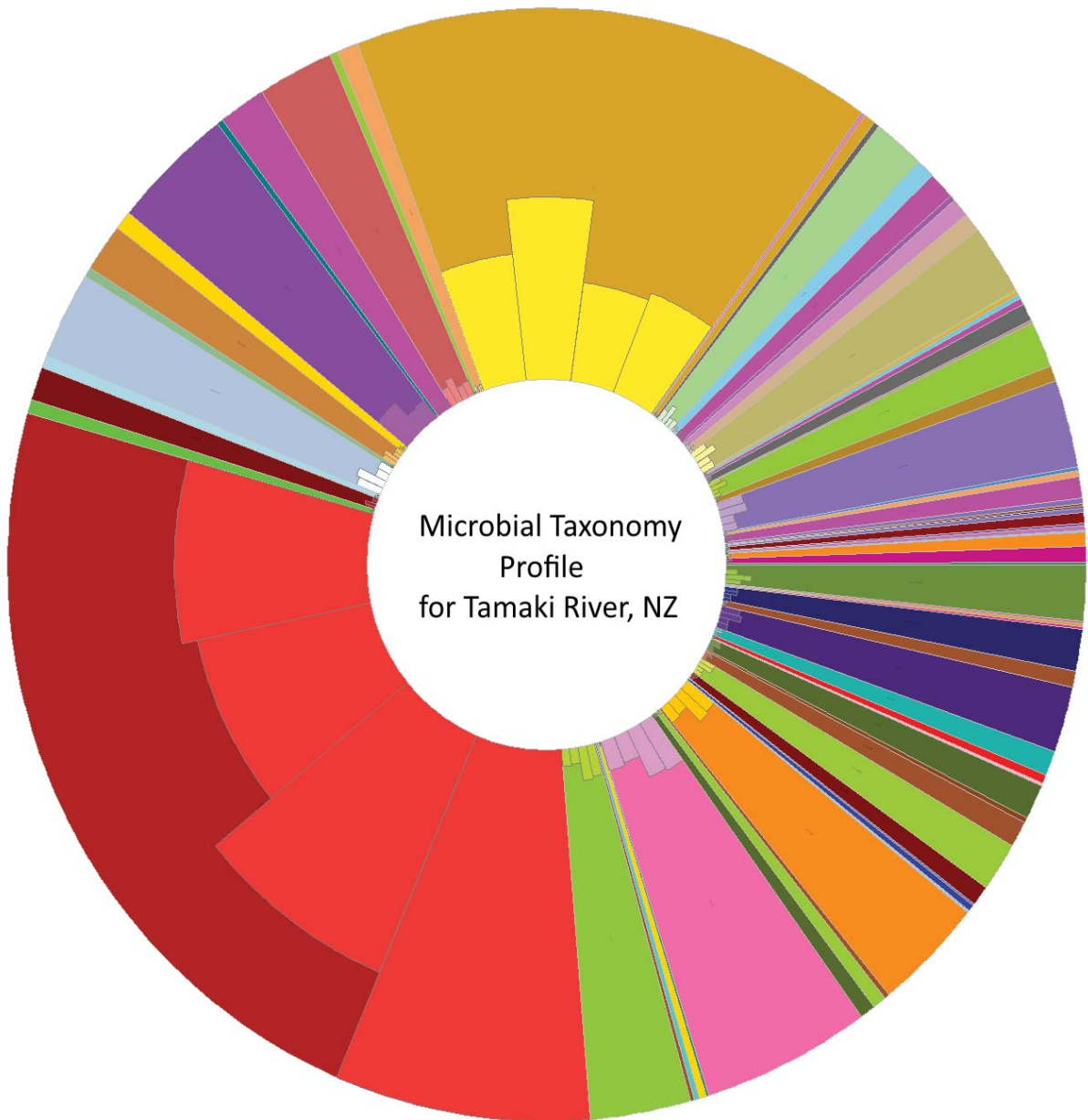# Institute of Fundamental Sciences

## A comparison of next-generation sequencing protocols for microbial profiling

**A thesis**

**submitted in partial fulfillment**

**of the requirements for the degree of**

**Master of Science in Genetics**

**By**

**Yang Fong (Richard)**

**2016**

**Massey University, Palmerston North**

**New Zealand**

Microbial Taxonomy
Profile
for Tamaki River, NZ

One of the responsibilities faced by the Environmental Genome Project is to provide the science base upon which society can make better informed risk management decisions.

-Samuel Wilson-

## Abstract

The introduction of massive parallel sequencing has revolutionized analyses of microbial communities. Illumina and other Whole Genome Shotgun Sequencing (WGS) sequencing protocols have promised improved opportunities for investigation of microbial communities. In the present work, we compared and contrasted the findings from different NGS library preparation protocols (Illumina Nextera, Nextera-XT, NEXTFlex PCR-free and Ion-Xpress-400bp) and two sequencing platforms (MiSeq and Ion-Torrent). Short reads were analysed using the rapid database matching software PAUDA and visualization software MEGAN5, which provides a conservative approach for taxonomic identification and functional analyses. In analyses of a Tamaki River water sample, biological inferences were made and compared across platforms and protocols. For even a relatively small number of reads generated on the MiSeq sequencing platform important pathogens were identified in the water sample. Far greater phylogenetic resolution was obtained with WGS sequencing protocols than has been reported in similar studies that have used 16S rDNA Illumina sequencing protocols. TruSeq and Nextera-XT sequencing protocols produced similar results. The latter protocol offered cheaper, and faster results from less DNA starting material. Proteobacteria (alpha, beta and gamma), Actinobacteria and Bacteroidetes were identified as major microbial elements in the Tamaki River sample. Our findings support the emerging view that short read sequence data and enzymatic library prep protocols provide a cost effective tool for evaluating, cataloguing and monitoring microbial species and communities. This is an approach that complements, and provides additional insight to microbial culture "water testing" protocols routinely used for analysing aquatic environments.

## Acknowledgement

There are many people I would like to express my gratitude and cordial thanks in helping me out in preparing my Master's Thesis. This dissertation would not have been possible without your support and strong collaboration between different academia backgrounds.

## Table of Contents

## List of Acronyms

| | |
|---|---|
| **%** | Percent |
| **°C** | Degrees Celsius |
| **μl** | Microlitre(s) |
| **μM** | Micromolar |
| **100 PE** | 2 x 100 base pair paired-end read |
| **150 PE** | 2 x 150 base pair paired-end read |
| **250 PE** | 2 x 250 base pair paired-end read |
| **300 PE** | 2 x 300 base pair paired-end read |
| **A** | Adenine |
| **A260** | Nanodrop absorbance at 260 nanometres |
| **A280** | Nanodrop absorbance at 280 nanometres |
| **AFLP** | Amplified Fragment Length Polymorphism |
| **ATL** | A-Tailing Mix |
| **ATM** | Amplicon Tagment Mix |
| **ATP** | Adenosine Triphosphate |
| **BAM** | Binary Alignment Matrix |
| **BGI** | Beijing Genomics Limited |
| **BIPES** | Illumina Multiplexed Paired-end Sequencing Adapter |
| **BLAST** | Basic Local Alignment Search Tool |
| **bp** | Base pair(s) |
| **C** | Cytosine |
| **CCD** | Charge-coupled Device |
| **cDNA** | Complementary Deoxyribonucleic Acid |
| **contig** | Continuous Sequence |
| **CTA** | A-Tailing Control |

| | |
|---|---|
| **CTE** | End-Repair Control |
| **CTL** | Ligation Control |
| **ddNTP** | Dideoxy Nucleotide Triphosphate |
| **dH$_2$O** | Distilled Water |
| **DNA** | Deoxyribonucleic Acid |
| **dNTP** | Deoxy Nucleotide Triphosphate |
| **ds** | Double Stranded |
| **EB** | Elution Buffer |
| **eDNA** | Environmental Deoxyribonucleic Acid |
| **EDTA** | Ethylenediamine Tetra-Acetic Acid |
| **emPCR** | Emulsion Polymerase Chain Reaction |
| **ERP** | End-Repair mix |
| **EtBr** | Ethidium Bromide |
| **E-value** | A parameter that describes the number of expected matches when searching a sequence database of a particular size and composition |
| **FC** | Flowcell |
| **fq** | Fastq File Format |
| **g** | Gram(s) |
| **G** | Guanine |
| **Gb** | Gigabytes |
| **gDNA** | Genomic DNA |
| **HiFi** | High fidelity enyzme |
| **HMW** | High Molecular Weight |
| **HT1** | Hybridization Buffer |
| **Inc.** | Incorporated |
| **Indel** | Small Insertion or deletion |
| **ISFET** | Ion Sensitive Field Effect Transistor |

| | |
|---|---|
| **KEGG** | Kyoto Encyclopedia of Genes and Genomes |
| **LCA** | Lowest Common Ancestor |
| **LIG** | Ligation Mix |
| **Log$^{10}$** | Logarithm to the base 10 |
| **M** | Molar |
| **Mb** | Megabytes |
| **MDA** | Multiple Displacement Amplification |
| **MEGAN** | Metagenome Analyzer |
| **MGS** | Massey Genome Service |
| **min** | Minute(s) |
| **ml** | Millilitre(s) |
| **mm** | Millimetre(s) |
| **mM** | Millimolar |
| **MPSS** | Massive Parallel Signature Sequencing |
| **mRNA** | Messenger Ribonucleic Acid |
| **mtDNA** | Mitochondrial Deoxyribonucleic Acid |
| **ng** | Nanogram(s) |
| **NGS** | Next-generation Sequencing |
| **No** | Number |
| **NPM** | Nextera PCR Master Mix |
| **NPS** | Non-point Source |
| **nt** | Nucleotide |
| **NT** | Neutralize Tagment Buffer |
| **NZGL** | New Zealand Genomics Limited |
| **OTU** | Operational Taxonomic Unit |
| **PAUDA** | Protein Alignment Using a DNA Aligner |

| | |
|---|---|
| **PCoA** | Principal Coordinate Analysis |
| **PCR** | Polymerase chain reaction |
| **pDNA** | Pseudo DNA |
| **PE** | Paired-end |
| **PGM** | Personal Genome Machine |
| **PhiX** | Bacteriophage PhiX174 |
| **PMM** | PCR Master Mix |
| **pmol** | Picomole(s) |
| **PPC** | PCR Primer Cocktail |
| **PP$_i$** | Pyrophosphate |
| **Q$_{10}$** | Phred Quality Score 1 error in 10 |
| **Q$_{20}$** | Phred Quality Score 1 error in 100 |
| **Q$_{30}$** | Phred Quality Score 1 error in 1000 |
| **QC** | Quality Control |
| **qPCR** | Quantitative Polymerase Chain Reaction |
| **Q-score** | Phred Quality Score |
| **RNA** | Ribonucleic Acid |
| **rpm** | Revolutions per Minute |
| **rRNA** | Ribosomal Ribonucleic Acid |
| **RSB** | Resuspension Buffer |
| **RTA** | Real-Time Analysis |
| **s** | Second(s) |
| **SAM** | Sequence Alignment Map |
| **SBS** | Sequencing by Synthesis |
| **SCIMM** | Sequence Clustering with Interpolated Markov Models |
| **SEED** | Database infrastructure for comparative genomics in MEGAN5 software |

| | |
|---|---|
| **SMRT** | Single Molecule Real Time |
| **spp.** | Species |
| **SPRI** | Solid Phase Reversible Immobilization |
| **ss** | Single Stranded |
| **STL** | Stop Ligation Buffer |
| **T** | Thymine |
| **TAE** | Tris-Acetate EDTA buffer |
| **TAP** | Taxonomic Assignment Pipeline |
| **Taq** | *Thermus aquaticus* |
| **TB** | Tuberculosis |
| **TD** | Tagmentation Buffer |
| **TE** | Tris EDTA Buffer |
| **V** | Volts |
| **WGS** | Whole Genome Shotgun Sequencing |

## List of Figures

## List of Tables