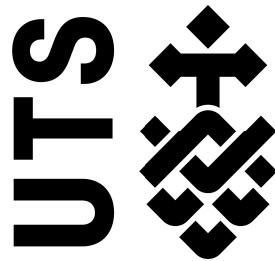


# **The integration of fine-scale DNA-DNA associations by inclusion of Hi-C DNA cross-linking information into metagenomic community analysis**

by

**Matthew Zachariah DeMaere**

A Thesis Submitted In  
Fulfillment of the Requirements  
for the Degree of  
Doctor of Philosophy



University of Technology Sydney  
Sydney, NSW, Australia

February, 2019



Nature uses only the longest threads to weave her patterns, so that each small piece of her fabric reveals the organisation of the entire tapestry.

*Richard Feynman - The Character of Physical Law: Messenger Lectures*

# Certificate of Original Authorship

I, Matthew Zachariah DeMaere, declare that this thesis is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the Faculty of Science at the University of Technology Sydney. This thesis is wholly my own work unless otherwise reference or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis. This document has not been submitted for qualifications at any other academic institution. This research is supported by the Australian Government Research Training Program.

Production Note:  
Signature removed  
prior to publication.

Date: 2019-03-11

Matthew Z. DeMaere

# Acknowledgments

I would like to thank my principal supervisor Associate Professor Aaron Darling for his time, advice and support during my candidature. In particular, for supporting the original initiative, a patient ear for long digressions and tolerance for random-walk problem solving. In addition, I would like to thank my co-supervisor Professor Steven Djordjevic for his time, advice and crucial support. I would like to thank the institute for fostering an environment conducive to quality research and their continual pursuit of a positive postgraduate student experience. I would also like to thank Professor Ian Charles for his valuable assistance in initiating and arranging the candidature.

I would like to thank my mother Jennifer, father James and sister Kathleen for their enthusiasm and continued support in what has been a *long* journey.

Lastly, I would especially like to thank my wonderful wife, Klára, for her ceaseless support despite the increasing complication of our lives and her ever expanding contribution to this endeavor.

Matthew Z. DeMaere  
Sydney, Australia, 2019.

# Thesis Format

A thesis by publication.

# List of Publications

## Peer reviewed

1. **DeMaere, M. Z., & Darling, A. E.** (2016). Deconvoluting simulated metagenomes: the performance of hard- and soft- clustering algorithms applied to metagenomic chromosome conformation capture (3C). *PeerJ*, 4, e2676.  
<https://doi.org/10.7717/peerj.2676>
2. **DeMaere, M. Z., & Darling, A. E.** (2018). sim3C: simulation of Hi-C and Meta3C proximity ligation sequencing technologies. *GigaScience*, 7(2).  
<https://doi.org/10.1093/gigascience/gix103>
3. **DeMaere, M. Z., Darling, A. E.** (2019). bin3C: exploiting Hi-C sequencing data to accurately resolve metagenome-assembled genomes. *Genome Biology*, 20(1), 46.  
<https://doi.org/10.1186/s13059-019-1643-1>

# Contents

<b>Certificate of Authorship</b>	iv
<b>Acknowledgments</b>	v
<b>Thesis Format</b>	vi
<b>List of Publications</b>	vii
<b>List of Figures</b>	xii
<b>List of Tables</b>	xiii
<b>Abstract</b>	1
<b>1 Introduction</b>	3
1.1 Background . . . . .	3
1.1.1 The uncultured majority . . . . .	3
1.1.2 Metagenomics . . . . .	4
1.1.3 Metagenome-assembled genomes . . . . .	5
1.1.4 Validation of metagenome-assembled genomes . . . . .	6
1.1.5 Chromosome conformation capture . . . . .	7
1.1.6 Addressing Metagenomics with Hi-C . . . . .	8
1.1.7 Graphical Model . . . . .	10
1.1.8 Community Detection . . . . .	11
1.1.8.1 Validation Testing . . . . .	12
1.1.8.2 Resolution Limits . . . . .	13
1.2 Research Aim . . . . .	13
1.3 Outline of Thesis . . . . .	14
1.3.1 Objectives . . . . .	14
1.3.2 Chapter Summaries . . . . .	14



1.4	Further Work . . . . .	16
1.5	List of Abbreviations . . . . .	17
1.6	References . . . . .	18
1.7	Appendices . . . . .	35
1.7.1	Definitions . . . . .	35
<b>2</b>	<b>Deconvoluting simulated metagenomes: the performance of hard- and soft- clustering algorithms applied to metagenomic chromosome conformation capture (3C)</b>	<b>38</b>
2.1	Authorship Declaration . . . . .	38
2.2	Abstract . . . . .	39
2.3	Introduction . . . . .	40
2.4	Materials and Methods . . . . .	43
2.4.1	Representation . . . . .	43
2.4.2	Clustering . . . . .	43
2.4.3	Appropriate Validation Measures . . . . .	44
2.4.4	Clustering Algorithm Selection . . . . .	45
2.4.5	Gold Standard . . . . .	47
2.4.6	Graph Generation . . . . .	47
2.4.7	Validation . . . . .	48
2.4.8	Simulating Hi-C/3C read-pairs . . . . .	49
2.4.9	Pipeline Design . . . . .	49
2.4.10	Simulation . . . . .	51
2.4.11	Parameter Sweep . . . . .	52
2.4.12	Assembly Entropy . . . . .	54
2.4.13	Graph Complexity . . . . .	54
2.5	Results . . . . .	55
2.5.1	Experimental Design . . . . .	55
2.5.2	Assembly Complexity . . . . .	56
2.5.3	Graph Complexity . . . . .	57
2.5.4	Clustering Validation . . . . .	58
2.6	Discussion . . . . .	62
2.6.1	Limitations and Future Work . . . . .	63
2.7	Conclusion . . . . .	66
2.8	Additional Information and Declarations . . . . .	67
2.8.1	Competing Interests . . . . .	67
2.8.2	Author Contributions . . . . .	67

2.8.3	Data Availability . . . . .	67
2.8.4	Funding . . . . .	67
2.9	List of abbreviations . . . . .	67
2.10	Acknowledgments . . . . .	68
2.11	References . . . . .	69
2.12	Appendices . . . . .	78
<b>3</b>	<b>sim3C: simulation of Hi-C and Meta3C proximity ligation sequencing technologies</b>	<b>79</b>
3.1	Authorship Declaration . . . . .	79
3.2	Abstract . . . . .	80
3.3	Findings . . . . .	81
3.3.1	Software testing . . . . .	81
3.3.2	3C sequencing . . . . .	82
3.3.3	Experiment scenarios . . . . .	83
3.3.4	Error Modelling . . . . .	84
3.3.5	Simulation modes . . . . .	85
3.3.6	Structurally related interactions . . . . .	89
3.3.7	Example scenarios . . . . .	92
3.3.8	Bacterial . . . . .	92
3.3.9	Eukaryotic . . . . .	95
3.3.10	Metagenomic . . . . .	96
3.3.11	Limitations and future work . . . . .	100
3.4	Methods . . . . .	100
3.4.1	Reference Data . . . . .	100
3.4.2	Read Generation . . . . .	101
3.4.3	Contact Maps . . . . .	103
3.5	Availability of data and materials . . . . .	103
3.6	Availability of supporting source code and requirements . . . . .	104
3.7	List of abbreviations . . . . .	104
3.8	Declarations . . . . .	105
3.8.1	Funding . . . . .	105
3.8.2	Authors contributions . . . . .	105
3.9	Acknowledgements . . . . .	105
3.10	References . . . . .	106

<b>4 bin3C : exploiting Hi-C sequencing data to accurately resolve metagenome-assembled genomes (MAGs)</b>	112
4.1 Authorship Declaration . . . . .	112
4.2 Abstract . . . . .	113
4.3 Background . . . . .	114
4.4 Method . . . . .	117
4.4.1 Simulated Community . . . . .	117
4.4.2 Read-set generation . . . . .	118
4.4.3 Ground Truth Inference . . . . .	119
4.4.4 Performance Metrics . . . . .	119
4.4.5 Real Microbiome . . . . .	120
4.4.6 Initial Processing . . . . .	120
4.4.7 Hi-C Read Mapping . . . . .	121
4.4.8 Contact Map Generation . . . . .	121
4.4.9 Bias Removal . . . . .	122
4.4.10 Genome binning . . . . .	123
4.5 Results . . . . .	124
4.5.1 Simulated Community Analysis . . . . .	124
4.5.2 Library Recommendations . . . . .	126
4.5.3 Real Microbiome Analysis . . . . .	129
4.5.4 Comparison to previous work . . . . .	131
4.6 Discussion . . . . .	134
4.6.1 Limitations and future work . . . . .	134
4.7 List of abbreviations . . . . .	136
4.8 Declarations . . . . .	136
4.8.1 Author contributions . . . . .	136
4.8.2 Competing interests . . . . .	137
4.8.3 Consent for publication . . . . .	137
4.8.4 Ethics approval and consent to participate . . . . .	137
4.8.5 Funding . . . . .	137
4.8.6 Availability of data and materials . . . . .	137
4.8.7 Supporting tools . . . . .	138
4.9 Acknowledgements . . . . .	138
4.10 References . . . . .	139
4.11 Appendices . . . . .	148

# List of Figures

1.1	Main steps of the Hi-C library protocol . . . . .	9
2.1	Parametric sweep pipeline . . . . .	50
2.2	Two simple phylogenetic models . . . . .	52
2.3	Branch length scale factor . . . . .	53
2.4	Impact of evolutionary divergence on assembly . . . . .	59
2.5	First two principal components of sweep . . . . .	60
2.6	External validation over parametric sweep . . . . .	61
3.1	Logical schema of Hi-C and meta3C modes . . . . .	88
3.2	Statistical modelling of proximity ligations . . . . .	91
3.3	Prokaryotic contact maps, simulated vs real . . . . .	94
3.4	Eukaryotic contact maps, simulated vs real . . . . .	96
3.5	Metagenomic contact maps, simulated vs real . . . . .	99
4.1	Synthetic human gut microbiome composition . . . . .	118
4.2	Validation of bin3C genome binning solutions . . . . .	127
4.3	Simulated community extracted MAGs by rank . . . . .	128
4.4	The effect of shotgun depth on solution quality . . . . .	130
4.5	Single vs dual enzyme digestion . . . . .	131
4.6	Extracted MAGs by rank for real microbiome . . . . .	132
4.7	Comparison of bin3C to MaxBin and ProxiMeta . . . . .	133
4.S1	Modelled abundance and returned cluster medians . . . . .	148
4.S2	Clustered contact map . . . . .	153
4.S3	Completeness and contamination of simulated microbiome . . . . .	154
4.S4	Extracted MAGs by rank at half shotgun depth . . . . .	155

# List of Tables

2.1	Parametric sweep variable ranges . . . . .	55
2.2	Parametric sweep combinatorial tally . . . . .	56
2.S1	Scale factor vs phylogenetic distance . . . . .	78
3.1	Real datasets used for comparison . . . . .	101
3.2	Composition of synthetic Hi-C community . . . . .	102
3.3	Composition of synthetic meta3C community . . . . .	102
3.4	Simulation runtime parameters . . . . .	103
4.1	Simple proposed MAG quality standard . . . . .	115
4.2	Genome Standards Consortium quality standard (MIMAG) . . . . .	115
4.3	Assembly statistics for the synthetic and real metagenomes . . . . .	121
4.S1	Reference genomes used in simulated community . . . . .	149
4.S2	Simulated community CheckM validation at full depth . . . . .	150
4.S3	Real microbiome cluster report and CheckM validation . . . . .	151
4.S4	Genome binning statistics for simple MAG ranks . . . . .	152
4.S5	Genome binning statistics for GSC MIMAG ranks . . . . .	152

# Abstract

Much of our understanding of the microbial world has been obtained using culture-based methodologies, a paradigm that has stood since the 19th century. And yet it has long been known that most of the Earth's microbial species are resistant to laboratory culture. It is reasonable to expect; therefore, that applying equal scrutiny to all microbial life will lead to significant discoveries. Motivated by this, metagenomics eliminates the culturing dependency by directly sampling DNA from an environment; successfully shedding light on the once unseen majority.

The technical limitations of present-day sequencing technologies have meant, however, that in achieving culture-independence, traditional shotgun metagenomic sequencing experiments make a considerable sacrifice. That sacrifice comes in the form of information loss where, in preparing DNA for sequencing, much of the "same-cell" and "same-chromosome" information is destroyed; information which is essential when reconstructing the individual genomes. Purely computational solutions to overcoming this sacrifice have proved insufficient; surpassed instead by strategies which employ changes in the experimental design aimed at reducing the information loss.

A recent strategy is the inclusion of a new form of sequencing data, provided by the Hi-C sequencing technique. Originally conceived to study the three-dimensional structure of chromatin, the Hi-C sequencing technique captures *in vivo* proximity interactions between DNA loci in an all-vs-all manner. When applied to direct metagenomic sampling, the physical structure of the microbial community (chromosome, cell and community) strongly influences the probability of observing proximity interactions between loci, and this pronounced modulation can be exploited to recover the information lost during shotgun sequencing.

This thesis details the effective integration of Hi-C into metagenomic sequencing studies to accurately reconstruct individual genomes, thereby deconvoluting the metagenome. To accomplish this, first an *in silico* investigation of the effectiveness of graph clustering as a

means of metagenome deconvolution was conducted; where Hi-C proximity interactions defined the edges and assembly contigs defined the nodes. A parametric sweep of experimental and community composition parameters was carried out, exploring how the degree of evolutionary divergence (from species to strains) affected the quality of deconvolution. For each iterate in the sweep, a ground-truth was constructed and quality assessed using a novel external validation measure supporting overlapping clusters and variable object weights.

This work led to the design and implementation of the first metagenomic Hi-C read-pair simulator, `sim3C`, capable of simulating complex community definitions and simple three-dimensional structural elements. While in pursuit of the final objective of metagenome deconvolution, `sim3C` enabled an externally validated development process.

Lastly, as the outcome of the final objective, `bin3C` is demonstrated; an open-source solution to Hi-C driven metagenome deconvolution. In an unsupervised manner, `bin3C` reconstructs individual genomes from metagenomic data. Using external validation of simulated data, `bin3C` is shown to have high precision and good recall. When a real human microbiome was analysed, `bin3C` achieved leading performance, resolving 20 more nearly-complete MAGs (57% gain) than its closest competitor.

Dissertation directed by Associate Professor Aaron E. Darling  
i3 institute, Faculty of Science  
University of Technology Sydney