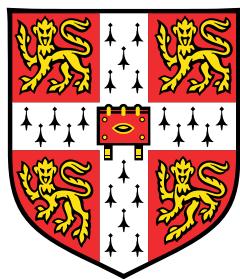


Graphical pangenomics



Erik Peter Garrison

Wellcome Sanger Institute
University of Cambridge

This dissertation is submitted for the degree of
Doctor of Philosophy

Fitzwilliam College

September 2018

for E_2 & E_3

Declaration

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is my own work and contains nothing which is the outcome of work done in collaboration with others, except as specified in the text and Acknowledgements. This dissertation contains fewer than 65,000 words including appendices, bibliography, footnotes, tables and equations and has fewer than 150 figures.

Erik Peter Garrison
September 2018

Graphical pangenomics

Erik Peter Garrison

Completely sequencing genomes is expensive, and to save costs we often analyze new genomic data in the context of a reference genome. This approach distorts our image of the inferred genome, an effect which we describe as reference bias. To mitigate reference bias, I repurpose graphical models previously used in genome assembly and alignment to serve as a reference system in resequencing. To do so I formalize the concept of a *variation graph* to link genomes to a graphical model of their mutual alignment that is capable of representing any kind of genomic variation, both small and large. As this model combines both sequence and variation information in one structure it serves as a natural basis for resequencing. By indexing the topology, sequence space, and haplotype space of these graphs and developing generalizations of sequence alignment suitable to them, I am able to use them as reference systems in the analysis of a wide array of genomic systems, from large vertebrate genomes to microbial pangenomes. To demonstrate the utility of this approach, I use my implementation to solve resequencing and alignment problems in the context of *Homo sapiens* and *Saccharomyces cerevisiae*. I use graph visualization techniques to explore variation graphs built from a variety of sources, including diverged human haplotypes, a gut microbiome, and a freshwater viral metagenome. I find that variation aware read alignment can eliminate reference bias at known variants, and this is of particular importance in the analysis of ancient DNA, where existing approaches result in significant bias towards the reference genome and concomitant distortion of population genetics results. I validate that the variation graph model can be applied to align RNA sequencing data to a splicing graph. Finally, I show that a classical pangenomic inference problem in microbiology can be solved using a resequencing approach based on variation graphs.

Acknowledgements

This work responds to ideas that arose in conversation with Deniz Kural. Our friendship is the first reason that I became a biologist, and his exploration of graphical models for genomes inspired my own. It is thanks to Alexander Wait Zaranek that we had the opportunity to work in George Church’s lab, which pulled us both into biology from our previous fields. There I met Madeline Price Ball, who guided me during an immersive and engaging introduction to biology and genomics.

Deniz introduced me to Gabor Marth, with whom I apprenticed in the art of bioinformatics. Gabor encouraged me to contribute extensively to the 1000 Genomes Project, whose objective captured my imagination and whose participants, in particular the members of the analysis group, taught me many lessons in the way of science. I can thank Hyun Min Kang, Goncalo Abecasis, Adam Auton, Laura Clarke, Gerton Lunter, Mark DePristo, Lisa Brooks, Ryan Poplin, Zamin Iqbal, and Heng Li, for always motivating me, and for helping me to understand and correct the many mistakes I made. Meanwhile, Mengyao Zhao and Wan-Ping Lee gave me my first look inside the alignment algorithms that are such an important part of this thesis.

During those years I had the pleasure of living with Benjamin “Mako” Hill and Mika Matsuzaki, who showed me what it means to work as a scientist for the commons. Not only did I learn from them, but from the many thinkers, dreamers, and travelers who they brought into our life in Somerville. These include Hanna Wallach, who helped me to understand the theory and practice of the learning problems I first encountered in genomics, and Nicolás Della Penna, who continues to shape my understanding of many aspects of the scientific artifice, in particular the fuzzy boundary between the social and the technical. I thank my friends Nathan Trachimowicz and Barbara Eghan, with whom I passed so much time in those years, for not letting me lose touch with the beautiful natural and human world in which this work lives and from which it derives its purpose.

This thesis would cover a considerably narrower set of topics if not for the efforts of the many people with whom I have worked to build the variation graph toolkit, *vg*. These include, but are not limited to: Jouni Sirén, from whom I learned about the world of succinct data structures as I watched him build the graph sequence indexes that brought

`vg` to a level of quality I never could have achieved on my own; Benedict Paten, who applied his unique expertise on genome graphs to help guide the effort of the ever-growing group of project collaborators without a pause in his own stream of contributions; Eric Dawson, whose ready conversation, energy and persistence buoyed me in the early days of `vg`, and who laid the foundation for future work on structural variant calling on graphs; Shilpa Garg, who brought new ideas about assembly and diploid genome inference to our project while helping to establish the long read alignment algorithms in `vg`; Adam Novak, who arrived first and transformed the heart of `vg` from a weak toy model into a foundation suitable for the work of this wide-ranging team, and who continues to carry it forward; Charles Markello, whose precise work on building resequencing pipelines with `vg` has ensured it is and will be widely usable; Emily Kobayashi, who further honed `vg`'s topology index and is pushing improvement in the dynamic graph indexes we use; Wolfgang Beyer, Toshiyuki Yokoyama, Orion Buske, and Ryan Wick, whose work on sequence graph visualization gave me eyes to understand my work; William Jones, whose clear-minded experiments on alignment identity and score comparison provide the basis for so many figures in this work; Hajime Suzuki, whose work on alignment acceleration lies at the core of the next phase of graph based mappers; Jordan Eizenga, with whom I explored the deep complexity of string to graph alignment algorithms; Xian Chang, who is driving the next high-performance paradigm for sequence to graph alignment; Mike Lin, whose experience and tact guided me both in the `vg` project and in extracurricular work for DNAexus; Yohei Rosen, who showed us that old models of haplotype matching could thrive inside our new pangenomic, graphical world; Glenn Hickey, who, in addition to caring for the project, took `vg` full circle and built variant calling into the graph model; Jerven Bolleman, who found a way to link `vg` into the enormous world of semantic biological data; and Toshiaki Katayama, who has explored our work and done so much to bring developers of graphical pangenomic techniques together. These members of the “vgteam” have each shared so much more than I can describe here, and I am deeply grateful to have had the opportunity to work with them. The work that I present here is fundamentally based on the productive collaboration that we have shared, and I could not have completed it without their watchful critique and steadfast exploration of their own research questions.

During my time as a student, I benefited from many long conversations with my fellow Sanger PhD cohort over lunch and coffee at the Genome Campus. In particular I learned from the other students with whom I lived: Ignacio Vázquez García, Martin Fabry, Daniel Bruder, Manuela Carrasquilla, and Girish Nivarti.

Pedro Fernandes gave Tobias Marschall and me the chance to teach a course on pangenomics using *vg*, which motivated many of the applications of variation graphs that I present here. Eppie Jones, Rui Martiniano, and Daniel Wegmann have guided and supported my work on ancient DNA. Remo Sanges and Mariella Ferrante gave me a lab to be part of and a fascinating project to explore in my time in Napoli. My corrections for this thesis drew on Alex Bowe’s excellent series of blog posts on succinct data structures, and I thank him for allowing me to exposit some of his examples here.

The final version of this thesis reflects a long and memorable conversation led by my examiners Aylwyn Scally and Gerton Lunter. I thank them for their clear suggestions for its improvement, and will forever be grateful for the time they devoted to this sprawling work. They focused my time on its roughest and most incomplete aspects. I am proud of the result that their critique has encouraged me to achieve.

Working with Richard Durbin has been a singular pleasure. Richard has an expansive vision for genomics, but he is always ready to dig into the details of a problem. He is a true master of his craft, able to support and guide every aspect of our work. The group he leads is motivated by his wide-ranging interests in biology. I owe its former and current members thanks for their encouragement and imagination.

Without my family, it is unlikely I would have ever begun the meandering trip that has led me to this thesis. My parents, Mark Garrison and Diane Garrison, helped me to be independent, and opened my mind to the world of ideas, which set me out on a wonderful trip. Along the way, my brother Nels and sister Astrid have kept me honest and careful of myself.

Much of this trip has been alongside my partner Enza Colonna. Non so come dire quanto mi ha aiutato, o quanti passi ho fatto in questo viaggio secondo le idee che abbiamo condiviso. I also am grateful to her parents, Donato Colonna e Concetta Tummolo, under whose almond and olive trees I wrote many pages of this work. Mi hanno sollevato dai problemi di vita quotidiana, e con il loro aiuto ho potuto scrivere la prima bozza di questa tesi in un solo mese.

Our daughter, Exa, who always convinces me to play, made sure that I was never too tired to keep going. I look forward to sharing this with her.

Table of contents

List of figures	xii
List of tables	xiv
1 Introduction	1
1.1 Genome inference	4
1.1.1 Reading DNA	4
1.1.1.1 The old school	5
1.1.1.2 “Next generation” sequencing	6
1.1.1.3 Single molecules	7
1.1.2 Genome assembly	9
1.2 Reference genomes	11
1.2.1 Resequencing	12
1.2.2 Sequence alignment	13
1.2.2.1 Compressed full text indexes	15
1.2.3 Variant calling	16
1.2.4 The reference bias problem	17
1.3 Pangenomes	18
1.3.1 On pangenomic models	20
1.3.2 The variation graph	23
1.4 Graphical techniques in sequence analysis	24
1.4.1 (Multiple) sequence alignment	24
1.4.2 Assembly graphs	27
1.4.2.1 Overlap graphs	28
1.4.2.2 De Bruijn graphs	29
1.4.2.3 String graphs	30
1.4.2.4 RNA sequencing graphs	32
1.4.2.5 Genome alignment graphs	32

1.4.3	Pangenomic alignment	34
1.4.3.1	Alignment to unfolded pangenomic references	34
1.4.3.2	Alignment to tiled pangenomic references	35
1.4.3.3	Alignment to graphical assembly models	36
1.4.3.4	Genotyping using a sequence DAG	37
1.4.3.5	Population reference graphs	38
1.4.3.6	Succinct pangenomic sequence indexes	39
1.4.3.7	Mapping to k -mer based pan-genome indexes	42
1.5	Overview and objectives	43
2	Variation graphs	45
2.1	A generic graph embedding for genomics	47
2.1.1	The bidirectional sequence graph	47
2.1.2	Paths with edits	48
2.1.3	Alignments	48
2.1.4	Translations	49
2.1.5	Genotypes	50
2.1.6	Extending the graph	50
2.2	Variation graph construction	51
2.2.1	Progressive alignment	51
2.2.2	Using variants in VCF format	51
2.2.3	From gene models	52
2.2.4	From multiple sequence alignments	53
2.2.5	From overlap assembly and de Bruijn graphs	53
2.2.6	From pairwise alignments	54
2.3	Data interchange	56
2.4	Index structures	57
2.4.1	Dynamic in-memory graph model	58
2.4.2	Graph topology index	58
2.4.3	Graph sequence indexes	61
2.4.3.1	Graph k -mer indexes	62
2.4.3.2	The FM-index and Compressed Suffix Array (CSA)	62
2.4.3.3	BWT-based tree and graph sequence indexes	66
2.4.3.4	The Generalized Compressed Suffix Array	68
2.4.3.5	GCSA2	70
2.4.4	Haplotype indexes	74
2.4.5	Generic disk backed indexes	78

2.4.6	Coverage index	78
2.5	Sequence alignment to the graph	79
2.5.1	MEM finding and alignment seeding	81
2.5.2	Distance estimation	82
2.5.3	Collinear chaining	83
2.5.4	Unfolding	86
2.5.5	DAGification	87
2.5.6	POA and GSSW	87
2.5.7	Banded global alignment and multipath mapping	90
2.5.8	X-drop DP	91
2.5.9	Chunked alignment	93
2.5.10	Alignment surjection	96
2.5.11	Base quality adjusted alignment	97
2.5.12	Mapping qualities	98
2.6	Visualization	99
2.6.1	Hierarchical layout	100
2.6.2	Force directed models	101
2.6.3	Linear time visualization	101
2.7	Graph mutating algorithms	105
2.7.1	Edit	105
2.7.2	Pruning	106
2.7.2.1	k -mer m -edge crossing complexity reduction	106
2.7.2.2	Filling gaps with haplotypes	107
2.7.2.3	High degree filter	107
2.7.3	Graph sorting	108
2.7.4	Graph simplification	108
2.8	Graphs as basis spaces for sequence data	110
2.8.1	Coverage maps	110
2.8.2	Bubbles	110
2.8.3	Variant calling and genotyping	112
3	Applications	114
3.1	Yeast	115
3.1.1	A SNP-based SGRP2 graph	115
3.1.2	Cactus yeast variation graph	117
3.1.3	Constructing diverse <i>cerevisiae</i> variation graphs	121
3.1.4	Using long read mapping to evaluate <i>cerevisiae</i> graphs	124

3.2 Human	126
3.2.1 1000GP graph construction and indexing	126
3.2.2 Simulations based on phased HG002	127
3.2.3 Aligning and analyzing a real genome	127
3.2.4 Whole genome variant calling experiments	129
3.2.5 A graph of structural variation in humans	131
3.2.6 Progressive alignment of human chromosomes	131
3.2.7 Building graphs from the MHC	133
3.2.8 ChIP-Seq	138
3.3 Ancient DNA	139
3.3.1 Evaluating reference bias in aDNA using simulation	140
3.3.2 Aligning ancient samples to the 1000GP pangenome	140
3.4 Neoclassical bacterial pangenomics	145
3.4.1 An <i>E. coli</i> pangenome assembly	146
3.4.2 Evaluating the core and accessory pangenome	146
3.5 Metagenomics	149
3.5.1 Arctic viral metagenome	150
3.5.2 Human gut microbiome	153
3.6 RNA-seq	156
3.6.1 Yeast transcriptome graph	156
4 Conclusions	158
References	161
Appendix Related publications	185

List of figures

1.1	The tree of life, reference genomes, and variation graphs	2
1.2	A variation graph	3
1.3	Computational pangenomics	21
1.4	Pangenomic models	22
2.1	The basic elements of a variation graph	46
2.2	A sketch of the XG index	60
2.3	An example of a suffix tree	63
2.4	Building the BWT and suffix array	64
2.5	Backward search in the BWT and suffix array	65
2.6	The XBW transform	67
2.7	Succinct de Bruijn graph construction	69
2.8	A sequence graph and its de Bruijn transformation	71
2.9	Searching in the GCSA2	73
2.10	The Graph Burrows Wheeler Transform	76
2.11	Alignment of a PacBio read to a yeast pangenome	80
2.12	Finding maximal exact matches (MEMs)	81
2.13	The MEM Chain Model	84
2.14	DAGification	88
2.15	The <i>dozeu</i> X-drop alignment algorithm	92
2.16	The Alignment Chain Model	94
2.17	Hierarchical visualization with Graphviz's <code>dot</code>	102
2.18	Force-directed layout with Graphviz's <code>neato</code>	103
2.19	Force-directed layout with Bandage	103
2.20	Linearized variation graph visualization	104
2.21	Pileup variant calling with <code>vg call</code>	112
2.22	Graph augmentation-based variant calling in <code>vg genotype</code>	113

3.1	Comparing alignment to the linear reference and SGRP2	118
3.2	Cactus yeast variation graph	119
3.3	Cactus yeast simulation	120
3.4	Whole genome alignment graphs for <i>S. cerevisiae</i>	123
3.5	Long read alignment against various <i>S. cerevisiae</i> pangenome graphs . .	125
3.6	Simulated reads from HG002 versus various human pangenome graphs. .	128
3.7	Indel allele balance in HG002	129
3.8	Alignment against the HGSVC graph	132
3.9	Seqwish assembly of the MHC in GRCh38.	135
3.10	<code>vg msga</code> progressive alignment of the MHC in GRCh38.	136
3.11	Path-coincidence dotplots from variation graphs of the MHC in GRCh38.	137
3.12	Resolving reference bias in 36bp ChIP-seq	138
3.13	Comparing <code>bwa aln</code> and <code>vg map</code> using simulated ancient DNA	141
3.14	Downsampling a high-coverage aDNA sample	143
3.15	Allele balance in the Yamnaya sample	143
3.16	<i>D</i> -statistic based ABBA-BABA test of reference bias in aDNA	144
3.17	An <i>E. coli</i> pangenome	147
3.18	Evaluating alignment to the <i>E. coli</i> pangenome.	148
3.19	An arctic freshwater viral metagenome	151
3.20	Comparing <code>vg</code> and <code>bwa</code> alignment to the viral metagenome	152
3.21	A human gut microbiome	154
3.22	Human gut microbiome alignment comparison	155
3.23	Aligning reads against the yeast transcriptome	157

List of tables

3.1	<i>S. cerevisiae</i> variation graphs	121
3.2	1000GP variation graphs	126
3.3	Selected results from the PrecisionFDA Truth Challenge	130