

Using Next-Generation Genomic Datasets In Disease Association

Luke Jostins

King's College
University of Cambridge
September 2012

This dissertation is submitted for
the degree of Doctor of Philosophy

Declaration

This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration except where specifically indicated in the Attributions section or the text. It does not exceed the word limit set out by the Degree Committee for the Faculty of Biology, and is not substantially the same as any work that has been, or is being, submitted to any other university for any degree, diploma or any other qualification.

Luke Jostins
18th September 2012

This is a post-viva dissertation, containing some minor corrections to that submitted on 18th September 2012. The corrections were suggested by Stephen Sawcer and Peter Donnelly.

Luke Jostins
24th January 2012

How are you going to be successful in treatment, if you do not understand the real essence of each disease?

— Galen, On Natural Faculties

Using Next-Generation Genomic Datasets In Disease Association

Luke Jostins
King's College, Cambridge
Ph.D Thesis, September 2012

Abstract

The first generation of genome-wide association studies (GWAS) uncovered thousands of genetic risk factors for hundreds of complex human diseases. However, over the past five years new high-throughput techniques, including next-generation sequencing and low-cost custom genotyping, have allowed us to expand disease association studies into larger sample sizes and across the entire spectrum of human variation. This thesis will explore the potential of these new technologies, and in particular their application to the study of Inflammatory Bowel Disease (IBD) genetics. After reviewing the historical context of complex disease genetics, I introduce the statistical methods and models used in this thesis, and demonstrate how they can be placed into a unified framework of genetic risk models. I then detail three analysis projects that focus on identifying risk variants that the first generation of GWAS was unable to study. The first investigates how genotype imputation, coupled with high-density sequencing reference sets, can aid locus discovery in both European and African populations. The second discusses the use of a custom genotyping chip (the ImmunoChip) to discover risk variants with low effect sizes, by allowing low-cost genotyping of a very large number of samples. The third investigates the use of next-generation sequencing of multiply affected (or “multiplex”) families in order to identify low-frequency, high penetrance risk alleles. Throughout these three projects I describe the discovery of a large number of novel IBD risk loci, and discuss how statistical and biological interrogation of these risk loci can help us to develop and expand biological hypotheses.

Acknowledgements

First and foremost I would like to thank my adviser Jeff Barrett for more things than I could reasonably list here. Needless to say, his dedication to statistical rigour, free sharing of data, clear scientific writing and public communication have provided me with a model scientist that I will always seek to emulate.

I would like express my gratitude to my other advisers, thesis committee members and first year examiners, including Richard Durbin, Julian Parkhill, Chris Tyler-Smith and Ralph McGinnis. I would especially like to thank my external adviser John Todd for the confidence he has shown in me over the last three years.

I would also like to thank the rest of Medical Genomics/ Barrett Group/ Statistical and Computational Genetics/ Team 143 (we'll settle on a name one day): Kate Morley, for teaching me most of what I know, James Morris, for providing the perfect combination of well-planned software development, wry humour and cake, Yang Luo for Chinese sweets and old-school Cambridge amiability (dan shi...), Iris Kolder for her irrepressible sense of mischief and Isabelle Claynen for her infectious passion for biology.

I'd also like to thank all the people that I have got to know during my time at Sanger: Annabel Smith, Christina Hedberg-Delouka and Alex Bateman for running such a tight ship. The Anderson Group, including Carl Anderson, Jimmy Liu and Jamie Floyd, for tea and laughs. Daniel MacArthur for acting as my perennial foil, and Liz Murchison for being a true friend. Sanger's army of pipeline developers past and present, including Thomas Keane, Shane McCarthy, Petr Danecek, Jim Stalker, Josh Randall and Martin Pollard, for

keeping my life (and the lives of hundreds of others) running smoothly, Chloe Noble and Anja Kolb-Kokocinski for somehow keeping Human Genetics from constantly falling over. Anne Wombwell, for taking the pain out of conference travel. All the other people who have made Sanger what is it for me, including (but definitely not limited to) Nicole Soranzo, Aaron Day-Williams, Aarno Palotie, Dan Gaffney, Matt Hurles and Helena Kilpinen.

This thesis would not have been possible without my collaborators, and I would like to thank them all: The clinicians and analysts of the UKIBDGC, for making me feel part of something important, including Miles Parks, Charlie Lees, James Lee, Tim Raine, John Mansfield, Christopher Matthew and Jack Satsangi. The many researchers of the IIBDGC, for showing me just how big “Big Science” can get, including Judy Cho, Andre Franke, Rinse Weersma, Sarah West, Yashoda Sharma, Phil Schumm, Kaida Ning, Severine Vermeire, and Dermot McGovern. All the people that make Boston a powerhouse of genetic analysis, including Stephan Ripke, Johan Essers, Hailiang Huang, Mark Daly, Soumya Raychaudhuri, Xinli Hu and Lizzy Rossin. Closer to home, the statistical know-how of the Diabetes and Inflammation Laboratory, including David Clayton and Chris Wallace. The dedication of the MalariaGEN researchers, Chris Spencer, Gavin Band, Katja Kivinen, Kirk Rockett, Quang Li and Dominic Kwiatkowski, for always doing the right thing, instead of the easy thing. I’d especially like to thank Adam Levine and Tony Segal from UCL for giving me a taste of real biology.

I would like to thank my fellow members of Genomes Unzipped, including Joe Pickrell, Don Conrad, Dan Vorhaus, Caroline Wright and Vincent Plagnol, for fighting all the good fights we so love to fight. And the many people I have met through blogging and twitter for drawing me into their community, and in particular Chris Gunter, Zoe McDougall and Razib Khan for all their support.

I’d also like to thank Daniel Stretch and Stephen Eglan from the Cambridge Computational Biology Institute, and Jules Griffin and Sarah Lummis from King’s College, for letting me loose on five years worth of Cambridge students. And I’d like to thank everyone at King’s, for taking in an overconfident Physics student just out of A-levels and producing an overconfident

geneticist eight years later.

I'd like to thank all of my friends in Cambridge who have been so important to forming how I think about the world, and my parents for putting me on it the first place. I'd also like to thank Hannah "Pixie" Price, for everything.

I would like to thank the research bodies who have funded the research in the thesis, including the Wellcome Trust, the Medical Research Council, the National Association for Colitis and Crohns disease, the Charles Wolfson Charitable Trust and the Gates Foundation, as well as the many international funders who supported the IIBDGC. I'd also like to thank Gil McVean, Anne Pratt, Ali Momin, Helen Thompson and the Wellcome Trust for ensuring that I have somewhere to go next.

Last, but certainly not least, I would like to thank the 104,341 (give or take a few) patients and volunteers who donated their DNA to all the studies that make up this thesis.

Attributions

Many of the projects described in this thesis are of a collaborative nature, and many were performed as part of large, international consortia. **Below** is a summary of the contributions of other scientists to the work described in this thesis.

Chapter 3

The African data used for testing imputation were generated as part of the MalariaGEN Consortial Project 1 (CP1). Samples were collected by partners from 12 centres, and full information on all sample collections is available from <http://www.malariagen.net/projects/cp1>. Genotyping was carried out at the Wellcome Trust Sanger Institute. Pre-imputation data processing, quality control and calling of genotypes was carried out by Kirk Rockett, Katja Kivinen, Gavin Band, Si Quang Le and Chris Spencer.

The list of loss-of-function (LoF) variants was generated by the 1000 Genomes LoF Group (now Functional Integration Group). The list of high quality, polymorphic LoF variants was generated by Daniel MacArthur.

Chapter 4

The samples used in the combined GWAS-Immonchip analysis were collected by the research groups of the International IBD Genetics Consortium (<http://www.ibdgenetics.org/groups.html>). Genotyping was performed across multiple centres, summarised in Table 4.8.

Quality control, genotype imputation and association testing for the

GWAS collections were performed by Stephan Ripke. The optiCall genotyping program used to call ImmunoChip samples was developed by Tejas Shah and Carl Anderson. The tag SNP meta-analysis technique was developed in collaboration with Stephan Ripke and Mark Daly, and implemented by Stephan Ripke. The phenotype likelihood modelling approach was developed and implemented by Jonah Essers.

The DAPPLE analysis was carried out by Stephan Ripke and Lizzie Rossin. The eQTL and cSNP analyses were carried out by Kaida Ning. The immune cell expression enrichment analysis was carried out by Xinli Hu and Soumya Raychaudhuri, and the gene expression network analysis was carried out by Ken Hui and Eric Shadt.

ImmunoChip cluster plots were manually inspected by Mitja Mitrovic, Jeff Barrett, Carl Anderson, Emilie Theatre, Tobias Balschun, Sarah West, Kaida Ning, Zhi Wei, Karin Fransen, Kyle Bailey, Isabelle Cleynen, Suzanne van Sommeren, Philippe Goyette, Sok Meng Evelyn Ng and Martin Ladouceur.

Chapter 5

All research described in Chapter 5 was carried in close collaboration with Adam Levine. There is no single element of this project that I carried out entirely without Adam's input, and virtually all of Section 5.4 was carried out jointly.

However, the following were carried out by Adam with little or no input from me:

- The collection of all samples and phenotype data
- The processing and analysis of the gene expression data
- The calling and quality control of genome-wide genotype data
- The calling of SNPs and indels in the exome data

Gavin Sewell selected the SNPs used to test the load of common IBD risk variants in the family. CytoSNP 12 genotyping was carried out at University College, London, and Sequenom genotyping, whole-genome sequencing

and whole-exome sequencing was carried out at the Wellcome Trust Sanger Institute. QC of the exomes was carried out by Martin Pollard.

Chapter 6

The whole genome sequencing experiment was designed in collaboration with the UK IBD Genetics Consortium, who provided the samples. Sample selection was performed by James Lee. Sequencing was carried out at the Wellcome Trust Sanger Institute, and sequence data was processed and quality controlled by Martin Pollard and Josh Randall. Variant calling, imputation, stringent sample and variant QC and association analysis were carried out by Yang Luo.

Publications

From this Dissertation

- Band, G., Quang, S. L., **Jostins, L.**, Pirinen, M., Kivinen, K., Jallow, M., Sisay-Joof, F., *et al* (2012) Imputation based meta-analysis of severe malaria in three African populations. (Under review).
- **Jostins, L.**, Ripke, S., Weersma, R., Duerr, R. H., McGovern, D. P., Hui, K. Y., Lee, J. C., *et al* (2012) Host-microbe interactions have shaped the genetic architecture of inflammatory bowel disease *Nature* **491**(7422):119-124.
- MacArthur, D. G., Balasubramanian, S., Frankish, A., Huang, N., Morris, J., Walter, K., **Jostins, L.**, *et al* (2012) A systematic survey of loss-of-function variants in human protein-coding genes. *Science* **335**(6070):823-828
- **Jostins, L.** and Barrett, J. C. (2011) Genetic risk prediction in complex disease. *Hum Mol Genet.* **20** (R2):R182-8.
- **Jostins, L.**, Morley, K. I. and Barrett, J. C. (2011) Imputation of low-frequency variants using the HapMap3 benefits from large, diverse reference sets. *Eur J Hum Genet.* **9**(6):662-666.

Arising elsewhere

- **1000 Genomes Project Consortium** (2012) An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**(7422):56-65.

- Liu, J. Z., Almarri, M. A., Gaffney, D. J., Mells, F. G., **Jostins, L.**, Cordell, H. J., Ducker, S. J., *et al* (2012) Dense fine-mapping study identifies new susceptibility loci for primary biliary cirrhosis *Nat Genet.* **44(10):1137-1141.**
- **Jostins, L.** (2012) Dispatches from the functional phase of genome biology *Genome Biol.* **13(6):316.**
- **Jostins, L.**, Pickrell, J. K., MacArthur, D. G. and Barrett J. C. (2012) Misuse of hierarchical linear models overstates the significance of a reported association between OXTR and prosociality. *Proc Natl Acad Sci USA* **109(18):E1048.**
- Sewell, G. W., Rahman, F. Z., Levine, A. P., **Jostins, L.**, Smith, P. J., Walker, A. P., Bloom, S. L., Segal, A. W. and Smith, A. M. (2012) Defective tumor necrosis factor release from Crohn's disease macrophages in response to toll-like receptor activation: Relationship to phenotype and genome-wide association susceptibility loci. *Inflamm Bowel Dis.* **8(11):2120-2127.**
- Franke, A., McGovern, D. P., Barrett, J. C., Wang, K., Radford-Smith, G. L., Ahmad, T., Lees, C. W., *et al* (2010) Genome-wide meta-analysis increases to 71 the number of confirmed Crohn's disease susceptibility loci. *Nat Genet.* **42(12):1118-1125.**
- **1000 Genomes Project Consortium** (2010) A map of human genome variation from population-scale sequencing. *Nature* **467(7319):1061-1073.**
- Krawitz, P., Rödelsperger, C., Jäger, M., **Jostins, L.**, Bauer, S. and Robinson, P.N. (2010) Microindel detection in short-read sequence data. *Bioinformatics* **26 (6):722-729.**

Contents

1	Introduction and Background	1
1.1	Justifications of phenotype and approach	1
1.1.1	Why study inflammatory bowel disease?	2
1.1.2	Why study complex disease genetics?	5
1.2	A brief history of human disease genetics	13
1.2.1	The age of molecular disease: 1940 to 1980	13
1.2.2	The age of linkage for Mendelian traits: 1980-1994	15
1.2.3	The beginning of complex disease genetics: 1994-2005	18
1.2.4	The technological build-up to genome-wide association studies: 1986-2005	23
1.2.5	The age of genome-wide association studies: 2005-Present	26
1.2.6	Technological advances post-GWAS: 2004-Present	29
1.2.7	Next-generation GWAS and post-GWAS studies	31
1.2.8	Conclusions	34
1.3	Outline of this thesis	36
2	Statistical methods and models of genetic risk	41
2.1	Introduction	41
2.1.1	Definition of a genetic risk model	43
2.1.2	Observable parameters of a genetic risk model	44

2.1.3	Genetic risk scores and link functions	46
2.2	From discrete genotypes to continuous risk	47
2.2.1	Properties of a sum of independent variables	48
2.2.2	General covariance for linear functions of allele count .	50
2.2.3	Covariance for non-linear functions of allele count . . .	53
2.3	The log risk model	58
2.3.1	Calculating parameters	59
2.3.2	Case and control distributions	62
2.3.3	Relationship to Risch model	64
2.3.4	Relationship to multiplicative relative risk model and log-linked regression	64
2.3.5	Problem with probabilities greater than 1	66
2.4	The probit risk model	68
2.4.1	Relationship to the liability threshold model	68
2.4.2	Calculating parameters	71
2.4.3	Case and control distributions	71
2.4.4	Relationship to probit regression and latent variable modelling	74
2.5	The logit risk model	76
2.5.1	Calculating parameters	77
2.5.2	Fitting the logit risk model numerically	79
2.5.3	Case and Control Distributions	80
2.5.4	Relationship to the multiplicative odds ratio model . .	80
2.5.5	Relationship to logistic regression	87
2.6	Comparing models of risk	88
2.6.1	Comparing disease probability distributions in cases . .	88
2.6.2	Comparing relative recurrence risk	91
2.6.3	Comparing ROC curves for risk prediction	93
2.7	Conclusion	96
2.7.1	Summary of models	96
2.7.2	Limitations of this approach	96
2.7.3	Problems generated by model ambiguity	99

3	Investigating new reference and target sets in genotype imputation	101
3.1	Introduction	101
3.1.1	Overview of imputation software and methods	102
3.1.2	New reference and target sets in imputation	104
3.2	The impact of reference set diversity in Europeans	107
3.2.1	Performing and Scoring Imputation	108
3.2.2	Reference Set Quality	110
3.2.3	Reference Set Size	111
3.2.4	Reference Set Diversity	115
3.2.5	Discussion	118
3.3	Imputation in African populations	121
3.3.1	HapMap-based imputation in a GWAS meta-analyses	122
3.3.2	1000 Genomes-based imputation in a single, diverse population	127
3.4	Using imputation to explore the impact of loss-of-function variants on complex disease	134
3.4.1	Loss-of-function variants and the 1000 Genomes project	134
3.4.2	Performing imputation and association analysis	135
3.4.3	Results	136
3.5	Concluding remarks	138
4	Investigating IBD genetics using the ImmunoChip	139
4.1	Introduction	139
4.1.1	Overview of this chapter	140
4.2	An overview of the ImmunoChip	142
4.2.1	The economics of the ImmunoChip	142
4.2.2	The content of the ImmunoChip	146
4.2.3	The biology of the ImmunoChip	153
4.3	QC and association analysis of the IIBDGC ImmunoChip dataset	158
4.3.1	The IIBDGC ImmunoChip dataset	158
4.3.2	Genotyping, imputation and quality control	161
4.3.3	Association analyses	166
4.3.4	GWAS and ImmunoChip analyses	166

4.3.5	Deep replication meta-analysis	168
4.3.6	Combining signals into loci	170
4.3.7	Crohn's disease/Ulcerative colitis likelihood modelling .	171
4.3.8	Comparison of this locus list to previous CD and UC lists	173
4.4	Biological and bioinformatic interpretation of 163 IBD loci . .	185
4.4.1	Global patterns in the "IBD genome"	186
4.4.2	IBD genetics in the context of autoimmunity and in- fection	187
4.4.3	Prioritising candidate genes in IBD loci	190
4.4.4	Testing for enrichment of functional terms within IBD loci	192
4.4.5	Natural selection in IBD loci	201
4.4.6	Gene expression analyses of IBD loci	203
4.4.7	Take home messages about the biology of IBD	205
4.5	IBD and Y haplogroups	208
4.5.1	Calling Y SNPs and assigning haplogroups	208
4.5.2	Association analyses and controlling for stratification .	209
4.5.3	Identifying candidate causal variants	211
4.6	Fine-mapping the <i>NOD2</i> locus	213
4.6.1	Characterising coding mutations in <i>NOD2</i>	213
4.6.2	Characterising a common regulatory signal at <i>NOD2</i> .	215
4.7	Conclusions	218
5	High-throughput genomic studies of multiplex families	221
5.1	Introduction	221
5.2	A history of multiplex family studies in complex disease	224
5.3	Modelling and controlling polygenic risk in multiplex families .	228
5.3.1	A combined polygenic/penetrant model of multiplex families	230
5.3.2	Risk prediction in multiplex families	241
5.4	Linkage and sequence analysis of a multiplex IBD family . . .	253
5.4.1	Description of the family	253

5.4.2	Segregation analysis	255
5.4.3	Known IBD risk variants in the family	256
5.4.4	Linkage and haplotype analysis of the family	259
5.4.5	Whole-genome sequencing in the family	265
5.4.6	Whole-exome sequencing in the family	271
5.4.7	Identifying candidate variants in the family	272
5.5	Follow-up of candidate causal variants	282
5.5.1	Technical validation of causal variants	283
5.5.2	Independent replication of causal variants	286
5.6	Conclusions	293
6	Conclusions	295
6.1	Connections and themes	295
6.2	A next-generation GWAS using low-coverage sequencing	299
6.3	Towards the ideal locus discovery experiment	303
6.4	Beyond locus discovery	307

