



UNIVERSITY OF  
CAMBRIDGE

# Primary sclerosing cholangitis: from genetic risk to disease biology

Elizabeth Claire Goode



Clare College

This dissertation is submitted for the degree of Doctor of Philosophy, May 2020



# Declaration

This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration except as declared in the Preface and specified in the text. It is not substantially the same as any that I have submitted, or am concurrently submitting, for a degree or diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. I further state that no substantial part of my dissertation has already been submitted, or is being concurrently submitted, for any such degree, diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. This dissertation does not exceed the prescribed limit of 60 000 words.

Elizabeth Claire Goode  
May, 2020



# Abstract

## **Primary sclerosing cholangitis: from genetic risk to disease biology**

**Elizabeth Claire Goode**

One in 10,000 people in the Western world lives with Primary Sclerosing Cholangitis (PSC), an immune-mediated, inflammatory disease of the bile ducts that is highly comorbid with inflammatory bowel disease (IBD). PSC confers risk of serious disease sequelae including hepatobiliary malignancy and progression to end-stage liver failure, for which the only treatment option is liver transplantation. The absence of effective medical therapies for PSC reflects our current limited understanding of the disease's aetiology and pathogenesis.

Our DNA, laid down at conception, gives us an unrivalled opportunity to understand the underlying causal biology of disease. This is because the genetic variants associated with disease susceptibility perturb genes and biological pathways that contribute to disease causality. Twenty-two regions of the genome, outside of the HLA, have been associated with PSC risk. These loci offer the potential for huge insight into the causal biology of PSC, if only we can robustly identify the true causal variants driving these loci and the genes they perturb. However, this is complicated by several scientific challenges. Firstly, the majority of disease-associated risk loci occur within non-coding regions of the genome. Secondly, patterns of correlation between variants within a risk locus means that the true causal variant driving the signal could be any of those highly correlated with the variant with the smallest p-value.

In this thesis, I present analyses aimed at identifying the true genes and causal variants underlying each of the twenty-two PSC risk loci. Many non-coding risk variants associated with complex disease exert a quantitative effect upon gene expression i.e. are expression quantitative trait loci (eQTLs). Colocalisation assesses the evidence that a single shared causal variant is responsible for driving PSC risk and gene expression via an eQTL. In order to assign dysregulated genes to PSC risk loci, I perform colocalisation with eQTLs mapped in multiple cell-types and tissues mechanistically relevant to PSC. Because PSC is rare, eQTLs have not previously been mapped in all cell-types most relevant to this disease. In addition, I therefore map eQTLs in six peripheral blood T-cell subsets (including the rare CCR9+ gut-homing T-cells) from ~80 patients with PSC and IBD. With colocalisation, I assign causal genes to five PSC risk loci, and assign other epigenetic regulatory features

including methylation or histone modification, to six risk loci. Statistical fine-mapping of each risk locus in both the GWAS and eQTL data enables me to resolve three PSC risk loci to a single causal variant and nine loci to 95% credible sets containing ten or fewer variants.

The results presented in this thesis identify three genes (*PRKD2*, *ETS2* and *UBASH3A*), causal in the pathogenesis of PSC, which are currently the target of existing or experimental therapeutic agents. Firstly, reduced expression of *PRKD2* causes excessive cell-autonomous T-follicular helper cell development and B-cell activation, and is associated with increased risk of PSC. Several studies are investigating the therapeutic effects of increasing the kinase activity of PRKD2. *ETS2* is involved in the induction of pro-inflammatory cytokine release from macrophages and IL-2 regulation in Th to Th0 transition. ETS2 inhibitors are currently the subject of early therapeutic trials. Finally, *UBASH3A* attenuates the NF- $\kappa$ B/I-KK $\beta$  pathway, an inflammatory pathway that is already targeted by proteasome inhibitors and acetylsalicylic acid, both of which could be potentially therapeutic in PSC.

PSC is a debilitating disease with serious disease sequelae, for which new therapeutic options are urgently needed. In this thesis, I elucidate multiple genes with a causal role in PSC pathogenesis, several of which are potential candidates for future therapeutic target.

# Acknowledgements

As I began to write these final, but perhaps most important, few paragraphs I realise how fortunate I have been to be assisted by so many wonderful individuals along the road to completing this thesis. There a great number of people to whom I owe thanks and would like to make special mention to just some of them below.

In May of 2016, I began the work presented in this thesis under the supervision of Dr Carl Anderson. I wish to express my sincere gratitude to Carl for welcoming me to his research group and for his support and guidance in the years since. It is only with his enthusiasm, encouragement, scientific brilliance and humour that the work in this thesis has been made possible. I would also like to express my sincere appreciation to my secondary supervisors, Professor Nicole Soranzo, Dr Tim Raine and Dr Simon Rushbrook who have provided additional scientific clarity, wisdom and support. I would like to pay special regards to Simon for introducing me to this field of research back in 2012, and his mentorship ever since.

There are three talented scientists, whose assistance has been integral to the work presented in thesis and all of whom have shown me extensive kindness and patience. Firstly, I would like to thank Dr Loukas Moutsianas for teaching a naïve medic to code and for his guidance with fine-mapping and colocalisation. Secondly, I am grateful to Dr Laura Fachal who guided me through the functional annotation analyses and was always there to answer every genetics question, big or small. Thirdly, I would like to express my gratitude to Dr Nikos Panousis, whose direction was invaluable in the RNA-seq and eQTL mapping component of this thesis. To all three of you, I am indebted.

For the past four years I have been fortunate to be supported by a fellowship from the University of Cambridge Wellcome Trust Clinical PhD program. The Wellcome Sanger Institute has been an incredible source of support throughout my fellowship and a wonderful place to spend the last few years. I would also like to thank each and every member of the Anderson Team at the Wellcome Sanger Institute. This team has evolved and expanded during my time at Sanger and I have been lucky to spend time with a group of such kind, talented and good-humoured individuals. I would like to pay special thanks to Dr Rebecca McIntyre who has been a source of personal support and who also performed all of the DNA extraction from my study samples. I would also like to thank Mr

Ben Bai for his guidance with mashR (*'Nobody really understands how mashR works...'*), Dr Alex Sazonovs for his help with Latex (*'It's us against Latex!'*), Dr Carla Jones for her patient explanations of T-cell biology and Mr Sigurgeir Olafsson, Alex and Ben for helping solve my innumerable coding issues!

I would also like to extend my gratitude to Dr Norihito Kawasaki and Dr Alex Wittmann at the Quadram Institute, for their guidance in the preparation of samples for my T-cell eQTL study, alongside the many patients from the Norfolk and Norwich University Hospital who donated their time (and blood!) to this study. In addition I would like to thank the patient support charity *PSC Support*, who generously funded the sample acquisition for the T-cell study.

To my esteemed examiners, Dr Richard Sandford and Professor Heather Cordell, I would like to express my gratitude for their invaluable comments and for a thoroughly enjoyable viva (albeit via Zoom)!

Undoubtedly, my proudest achievement during my PhD has been the arrival of my daughter, Alexandra. Thank you for being a source of such joy, for reminding me of the most important things in life, and for sleeping well so that Mummy could write her thesis! I would also like to express my sincere gratitude to Dr Pat Tate whose unwavering positivity encouraged me to put confidence in myself, and has kept me (mostly) sane. My final and most special debt of gratitude goes to my wonderful husband, Grégory for his great love and support. Je tiens à exprimer ma profonde gratitude pour ton amour et ton soutien et pour m'avoir aidée à prendre confiance en moi. To him, I am wholeheartedly grateful.

Finally, to my parents and the many family, friends and colleagues who I have not named individually, but who have also given me their support. I hope that along the way I have shown you my appreciation, and to you all I give a heart-felt thank you.



# Contents

<b>1</b>	<b>Introduction</b>	<b>19</b>
1.1	What is Primary Sclerosing Cholangitis? . . . . .	19
1.2	PSC is a complex disease . . . . .	21
1.3	Genome-wide association studies . . . . .	21
1.4	Genetic associations within the human leucocyte antigen . . . . .	23
1.5	Genetic associations outside of the HLA . . . . .	24
1.5.1	PSC risk loci in coding regions of the genome . . . . .	25
1.5.2	PSC risk loci in non-coding regions of the genome . . . . .	26
1.5.3	Current genetic understanding of PSC subtypes . . . . .	27
1.6	Current hypotheses of disease pathogenesis in PSC . . . . .	29
1.6.1	The ‘gut-homing T-cell’ hypothesis . . . . .	29
1.6.2	The ‘toxic bile’ hypothesis . . . . .	32
1.6.3	The ‘leaky gut’ hypothesis . . . . .	33
1.7	Challenges in deciphering PSC risk loci . . . . .	34
1.7.1	Expression quantitative trait loci . . . . .	35
1.7.2	Histone modification . . . . .	37
1.7.3	DNA methylation . . . . .	37
1.8	Translating genetic risk loci into biological drug targets . . . . .	38
1.9	Outline of this thesis . . . . .	39
<b>2</b>	<b>Fine-mapping of disease-associated risk loci in Primary Sclerosing Cholan- gitis</b>	<b>41</b>
2.1	Introduction . . . . .	41
2.2	Chapter overview . . . . .	44
2.3	Methods . . . . .	44
2.3.1	Fine-mapping . . . . .	44
2.3.2	Functional annotation . . . . .	46
2.4	Results . . . . .	48
2.4.1	Loci mapped to a single causal variant . . . . .	53
2.4.2	Variants with a greater than 50% posterior probability of causality .	55

2.4.3	Variants with a greater than 20% posterior probability of causality .	58
2.4.4	Loci not well-resolved with fine-mapping . . . . .	60
2.5	Discussion . . . . .	63

**3 Statistical colocalisation of Primary Sclerosing Cholangitis risk loci with functional quantitative trait loci 67**

3.1	Introduction . . . . .	67
3.2	Chapter overview . . . . .	69
3.3	Methods . . . . .	70
3.3.1	Colocalisation analysis . . . . .	70
3.3.2	Functional QTL data . . . . .	72
3.3.3	Fine-mapping of functional QTL loci . . . . .	74
3.4	Results . . . . .	76
3.4.1	The <i>PRKD2</i> locus . . . . .	80
3.4.2	The <i>ETS2</i> locus . . . . .	83
3.4.3	The <i>UBASH3A</i> locus . . . . .	85
3.4.4	The <i>SH2B3</i> locus . . . . .	88
3.4.5	The Chr18:67543688 locus . . . . .	88
3.5	Discussion . . . . .	89

**4 T-cell expression quantitative trait loci maps in Primary sclerosing cholangitis 93**

4.1	Introduction . . . . .	93
4.2	Chapter Overview . . . . .	94
4.3	Methods . . . . .	95
4.3.1	Sample type and Patient recruitment . . . . .	95
4.3.2	Sample preparation . . . . .	96
4.3.3	RNA extraction, library preparation and sequencing . . . . .	98
4.3.4	Read alignment, counts and quality control . . . . .	99
4.3.5	Differential gene expression . . . . .	104
4.3.6	Genotype QC and imputation . . . . .	109
4.3.7	eQTL mapping . . . . .	112
4.3.7.1	Identifying sample mismatches and amplification bias . . .	113
4.3.7.2	Identifying <i>cis</i> -eQTLs . . . . .	115
4.3.8	Identifying shared and tissue-specific eQTL . . . . .	117
4.3.9	Colocalisation . . . . .	118
4.4	Results . . . . .	119
4.4.1	Differential gene expression . . . . .	119
4.4.2	eQTL mapping . . . . .	123

4.4.3	Shared and tissue-specific eQTLs . . . . .	126
4.4.4	Colocalisation of disease-risk loci with eQTL . . . . .	126
4.5	Discussion . . . . .	138
<b>5</b>	<b>Conclusions</b>	<b>143</b>
	<b>Bibliography</b>	<b>149</b>



# List of Figures

1.1	Twenty of the twenty-two non-HLA PSC risk loci plotted according to their effect size (OR) and MAF in Ji <i>et al</i> 's GWAS data [42]. . . . .	25
1.2	Figure taken from Ji <i>et al</i> demonstrating odds ratios (and their 95% confidence intervals) for PSC, UC and CD across the 6 PSC associated SNPs demonstrating strong evidence for a shared causal variant (maximum posterior probability >0.8) [42]. . . . .	28
1.3	The 'gut-homing' T-cell hypothesis of PSC pathogenesis. . . . .	31
2.1	Power (y axis) to identify the causal variant in a correlated pair increases with the significance of the association (x axis), and therefore with sample size and effect size (vertical dashed line shows genome-wide significance level). Figure taken from Huang, Fang, Jostins <i>et al</i> [56]. . . . .	43
2.2	Summary of fine-mapping the PSC risk loci. . . . .	50
2.3	Regional association plots for PSC risk loci mapped to single variants. . . .	54
2.4	Regional association plots for PSC risk loci mapped to casual variants with >50% posterior probability of causality. . . . .	56
2.5	Regional association plots for PSC risk loci mapped to casual variants with >20% posterior probability of causality. . . . .	59
2.6	Regional association plots for PSC risk loci not well resolved with fine-mapping. . . . .	62
2.7	Regional association plots for PSC risk loci not well resolved with fine-mapping. . . . .	63
3.1	Schematic diagram of the GWAS fine-mapping - colocalisation - functional-trait fine-mapping pipeline to resolve the causal variants driving PSC risk loci, and the genes they perturb. . . . .	75
3.2	Colocalisation between seven PSC risk loci with UC and the evidence for PP4 and PP3 with varying $p^{12}$ . . . . .	76
3.3	Chr19:47205707 regional association plots for most probable fine-mapped SNP, rs313839, in PSC GWAS data and colocalising eQTL data for <i>PRKD2</i> in monocytes. . . . .	82

3.4	Chr21:40466744 regional association plot showing the most probable fine-mapped SNP for PSC GWAS (rs4817987) and colocalising eQTL data for <i>ETS2</i> in monocytes (fine-mapped to rs4817987) and for a H3K27ac histQTL in monocytes (fine-mapped to rs2836878). . . . .	84
3.5	Chr21:43855067 regional association plots for fine-mapped SNP, rs1893593, in PSC GWAS and colocalising eQTL data for <i>UBASH3A</i> and spliceQTL data for <i>UBASH3A</i> . . . . .	87
4.1	Sample preparation pipeline. . . . .	96
4.2	Gating strategy used for FACS separation of CD4+CCR9-, CD4+CCR9+, CD8+CCR9- and CD8+CCR9+ central effector T-cells from peripheral blood mononuclear cells. . . . .	98
4.3	Proportion of reads mapped to exons for a subset of 96 of the total 456 samples, highlighting an experimental outlier which was subsequently excluded due to a low proportion of reads mapped to exons compared to the mean. .	100
4.4	Principal component analysis of the top 500 most variably expressed genes, identifying two experimental outliers which did not cluster with their expected cell types. . . . .	102
4.5	PCA analysis of the top 500 most variably expressed genes, identifying four experimental outliers from two patients. . . . .	103
4.6	Expression of marker genes across all cell types. . . . .	104
4.7	Schematic representation of the <i>DESeq2</i> method of normalisation. . . . .	106
4.8	MA plots with and without shrinkage applied. Points are coloured red where the adjusted p-value is less than 0.05, and plotted as open triangles pointing either up or down if they fall outside of the window. . . . .	108
4.9	PCA of study samples compared to 1000 Genomes samples of known ethnicity using a pruned set of 62,805 independent variants with an $r^2 < 0.2$ and $MAF > 0.01$ . . . . .	110
4.10	Outline of pre-imputation QC of genotype data. . . . .	111
4.11	Concordance at heterozygous genotypes (x-axis) versus concordance at homozygous genotypes (y-axis), for each individual genotype sample (black dots). A match between genotype (box at top) and gene expression data (plot title) is coloured red (two left hand examples). A mismatch or amplification bias is coloured black (right hand example). . . . .	114
4.12	Concordance at heterozygous genotypes (X-axis) versus concordance at homozygous genotypes (Y-axis), for each individual genotype sample (black dots). An sample mismatch is shown by a match between a different genotype (in box at top) and gene expression data (plot title) in all four examples. . . . .	115

4.13	Gene ontology pathway analysis for DEGs in T-memory cells of UC compared to HC. Figure generated using g:profiler [236], 20/12/2019. . . . .	122
4.14	Number of significant eQTLs (y-axis) mapped for each individual cell type at 5% (blue line) and 10% FDR (red line), using covariate models with different numbers of gene-expression derived PCs from zero to fifty (x-axis). . . . .	124
4.15	Distance from transcription start site (TSS) for each significant eQTL (coloured red for those less than 5% FDR) per cell type. . . . .	125
4.16	Number of cell-type specific and shared QTLs. . . . .	126
4.17	Regional association plot for the Chromosome 21 rs1893592 risk locus in PSC GWAS data. . . . .	128
4.18	Regional association plots for colocalisation between PSC GWAS and eQTLs for <i>UBASH3A</i> in T-cells at Chromosome 21 rs1893592 risk locus, using <i>mashR</i> eQTL data. . . . .	129
4.19	Expression of <i>UBASH3A</i> according to Chromosome 21 rs1893592 genotype in T-memory cells. . . . .	130
4.20	Colocalisation between PSC GWAS and <i>AP003774.1</i> eQTL data from the individual cell-type analysis, at the chromosome 11 rs663743 PSC risk locus. . . . .	131
4.21	Colocalisation between PSC GWAS and <i>AP003774.1</i> eQTL data from the <i>mashR</i> analysis, at the chromosome 11 rs663743 PSC risk locus. . . . .	132
4.22	Expression of <i>AP003774.1</i> according to Chromosome 11 rs663743 genotype in T-regulatory cells. . . . .	133
4.23	Expression of <i>AP003774.4</i> across multiple human tissues (figure generated by GTEx portal, 25/02/20 [176]). . . . .	135
4.24	Expression of <i>AP003774.4</i> across multiple immune cell types (figure generated by the Database of immune cell eQTL expression [261], 26/02/2020). . . . .	136







