

**Whole-genome Sequencing-based
Association Studies of Cardiovascular Biomarkers**



Jie Huang

This dissertation is submitted to the University of Cambridge
Faculty of Biology for the degree of Doctor of Philosophy

Darwin College

University of Cambridge

February 2015

PREFACE

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 text.

This dissertation is not substantially the same as any that I have submitted, or, is being concurrently submitted 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 of similar institution except as declared in the Preface and specified in the text.

The dissertation does not exceed the page limit of 300 specified by the Biology Degree Committee

ABSTRACT

Background: Genome-wide association studies (GWAS) have significantly advanced the genetic study of complex human traits. With the advent of whole-genome sequencing (WGS) technologies and the increased capacity to identify rare variants, GWAS that use WGS data are expected to provide further opportunities for the discovery of variants that have larger and even causal effects. The UK10K project is one of the largest studies that use WGS to investigate the contribution of low frequency and rare genetic variants to medical traits.

Research aims: My research aims to address the utility of WGS-based imputation and associations for identifying the genetic determinants of a select quantitative traits that are associated with cardiovascular risks. Under the UK10K project framework, I study a suite of circulating biomarkers that have been reported for association with CVD. Specifically, I seek to evaluate the following three broad aspects: 1. what are the characteristics of phasing and imputation with WGS data? 2. what novel analytic methods could be applied to a large scale WGS based association study on a rich of phenotypes? 3. can I identify novel and potentially stronger effect genetic variants that are associated with the chosen CVD traits?

Methods: My study leverages existing WGS data from the UK10K project ($N = \sim 4,000$) and further uses it as a reference to impute more samples ($N > 10,000$) that have genome-wide SNP array data. In doing so, I first evaluate the quality of the WGS data and its utility for imputation, by comparing it to WGS data from the 1000 Genomes Project. Then, I examine the associations between genotypes and phenotypes for 13 quantitative traits, first in samples having WGS and then in samples having imputed data. The 13 CVD related biomarkers include four lipid traits (high-density lipoprotein cholesterol (HDL), low-density lipoprotein cholesterol (LDL), total cholesterol (TC), triglycerides (TG)), one inflammatory biomarker (C-reactive protein (CRP)), and eight haematological traits (hemoglobin (HGB), mean corpuscular hemoglobin (MCH), mean corpuscular hemoglobin concentration (MCHC), mean corpuscular volume (MCV), packed cell volume (PCV), platelet counts (PLT), red blood cell counts (RBC), white blood cell counts (WBC)).

To my dear parents: YuanYu Huang & Youquan Deng

To my wife: Weilin Chen

To my daughter Valerie and my son Jimmy

ACKNOWLEDGEMENTS

From the three years of my PhD study there are a lot of people I want to mention and thank for giving me help and advice. The first and most important person is Dr. Nicole Soranzo, my PhD thesis supervisor who gave me an opportunity to come to Cambridge and Sanger Institute to pursue my PhD study and work on a most important and frontier project on WGS. Her dedication to science and attentive guidance to trainees inspires me to become a good researcher and scientific leader myself.

I would like to thank the other faculty members who served on my thesis committee and provided strategic guidance for my PhD project: Dr. Richard Durbin, Dr. Carl Anderson, Dr. Eleftheria Zeggini from Sanger Institute, and Dr. Adam Butterworth from University of Cambridge.

I would like to thank our wonderful teammates (team 151): Lu Chen, Klaudia Walter, Louella Vasquez, Valentina Iotchkova, Massimiliano Cocca, Matthias Geihs, Yasin Memari, So-Youn Shin. Over the past three years, we sit close together and worked even closely with each other. Whenever there is a question, I feel that I could simply turn around my chair and get insightful feedback and help instantly.

I would also like to thank the larger community of the UK10K project and outsider collaborators who contributed data for our large meta-analysis, particularly to Drs. Nic Timpson and Josine Min from Bristol University.

Finally, I would like to thank my former supervisor, Dr. Chris O'Donnell at Framingham Heart Study, for offering me research experience in cardiovascular genetics and continuous support for my PhD study and career growth.

Table of Contents

1	Introduction	21
1.1	The burden of cardiovascular disease in modern society.....	21
1.2	Established and emerging risk factors for CVD	22
1.3	The allelic architecture of complex traits.....	26
1.4	Genome-wide association studies (GWAS).....	28
1.5	GWAS studies of CVD events and cardiovascular biomarkers.....	30
1.6	Rare variants and the motivation for whole genome sequencing (WGS).....	32
1.7	The UK10K Project	35
1.8	This thesis.....	39
2	Methods	41
2.1	Introduction	42
2.2	Study samples	43
2.2.1	UK10K WGS cohorts	43
2.2.2	UK10K GWA cohorts	43
2.2.3	Expanded discovery cohorts	44
2.3	Genetic data.....	49
2.3.1	UK10K WGS data.....	49
2.3.2	Imputation using WGS reference panel.....	51
2.4	Phenotype harmonization.....	51
2.5	Statistical methods for association studies	55
2.5.1	Power estimation.....	55
2.5.2	Single-variant based association studies	58
2.5.3	Loci selection for single marker results	59
2.5.4	Rare variants aggregation analysis.....	62
2.5.5	Loci selection for rare variant aggregation results.....	63
2.5.6	Other statistical methods	64
2.6	Conclusion & Discussion.....	67
3	Imputation.....	71
3.1	Introduction	72
3.1.1	How imputation works.....	72
3.1.2	Use of imputation in GWAS.....	72
3.1.3	Imputation with WGS reference panels	72
3.1.4	Aims of this study	73
3.2	Methods.....	75

3.2.1	WGS Reference Haplotypes	75
3.2.2	Test GWAS datasets	77
3.2.3	Running imputation	78
3.3	Results	80
3.3.1	Characteristics of UK10K WGS panel	80
3.3.2	Imputation evaluation on UK10K vs. 1000GP reference panels	83
3.3.3	Evaluation of metrics for choosing reference haplotypes	84
3.3.4	Evaluation of combining two reference panels	88
3.4	Conclusion & Discussion.....	90
4	Lipids	93
4.1	An introduction to lipids.	93
4.1.1	Biology and physiology circulating lipids.....	93
4.1.2	Lipids as risk factors for CVD	94
4.1.3	Genetic determinants of lipids levels.....	97
4.1.4	Aims of this study	103
4.2	Methods.....	104
4.2.1	Cohorts & phenotype measurements.....	104
4.2.2	Single marker based discovery and follow-up	108
4.2.3	Rare variant aggregation based discovery and follow-up	109
4.2.4	Fine-mapping of known loci	110
4.4	Results	111
4.4.1	Novel loci and novel variants from single marker analysis.....	111
4.3.2	Fine mapping of known and novel loci.....	123
4.3.3	Novel loci based on rare variants aggregation test	126
4.4	Conclusion & Discussion.....	130
4.4.1	Summary of main findings	130
4.4.2	Interpretation of results	130
4.4.3	Future direction.....	133
5	Full Blood Counts	135
5.1	An introduction to full blood counts	135
5.1.1	Biology and physiology of FBC.....	135
5.1.2	FBC traits as risk factors for CVD.....	136
5.1.3	Genetic determinants of FBC.....	137
5.1.4	Aims of this study	140
5.2	Methods.....	141
5.2.1	Cohorts & phenotype measurements.....	141

5.2.2	Single marker based discovery and follow-up	143
5.2.3	Rare variant aggregation based discovery and follow-up	143
5.2.4	Fine-mapping of known loci	144
5.3	Results	146
5.3.1	Novel loci and novel variants from single marker analysis.....	146
5.3.2	Fine mapping of known and novel loci.....	160
5.3.3	Novel loci based on rare variants aggregation test	162
5.3.4	Host-response eQTL.....	166
5.4	Conclusion & Discussion.....	168
5.4.1	Summary of main findings	168
5.4.2	Interpretation of results	168
5.4.3	Future direction.....	169
6	CRP	173
6.1	An introduction on CRP	173
6.1.1	Biology and physiology of circulating CRP	173
6.1.2	CRP as risk factors for CVD.....	174
6.1.3	Genetic determinants of CRP.....	175
6.1.4	Aims of this study	178
6.2	Methods.....	178
6.2.1	Cohorts & phenotype measurements.....	178
6.2.2	Single marker based discovery and follow-up	179
6.2.3	Rare variant aggregation based discovery and follow-up	180
6.2.4	Fine-mapping of known loci	180
6.3	Results	182
6.3.1	Novel loci and novel variants from single marker analysis.....	182
6.3.2	Fine mapping of known and novel loci.....	192
6.3.3	Novel loci based on rare variants aggregation test	193
6.4	Conclusion & Discussion.....	195
6.4.1	Summary of main findings	195
6.4.2	Interpretation of results	195
6.4.3	Future direction.....	196
Chapter 7.	Summary & Discussion.....	199
7.1	This thesis.....	199
7.2	Implication of findings for genetics of complex traits	199
7.3	Strength and limitations of the current study.....	203
7.4	Recommendations for future research in the field.....	205

7.4.1	Larger sample size with increased power.....	205
7.4.2	High genotyping accuracy through high-depth WGS	206
7.4.3	Better methods for rare variants aggregation test and replication	207
7.4.4	System biology approach that integrates various functional data	207
7.4.5	Pleiotropy analysis	208
7.4.6	Thinking genetics in the context of the trend of metabolic syndrome.	209
References.....		211
Appendix.....		238
Appendix 1 Manhattan plots of individual GWA.....		238

LIST OF TABLES

Table 1.1 List of traits in UK10K-Cohorts	37
Table 3.1 Sequence quality and variation metrics for UK10K Cohorts	81
Table 3.2 Descriptive for imputation reference panels	82
Table 4.1 Gene discovery in monogenic dyslipidemias	100
Table 4.2 GWAS studies of lipids	102
Table 4.3 NGS studies on lipids	103
Table 4.4 Characteristics of participating cohorts	106
Table 4.5 Phenotype harmonization protocol for lipids traits.....	107
Table 4.6 Putative novel variants of low or rare frequency from UK10K WGS.....	114
Table 4.7 Replication results of WGS top hits	115
Table 4.8 SKAT results for single point test top hits.....	116
Table 4.9 Expanded discovery(14-way meta-analysis) top hits	120
Table 4.10 Cohort specific results for four top variants based on 14-way meta-analysis	121
Table 4.11 Predictive causal variants based on fine mapping	125
Table 5.1 GWAS studies on FBC traits	140
Table 5.2 Phenotype harmonization protocol for FBC traits.....	142
Table 5.3 Characteristics of participating cohorts	145
Table 5.4 Putative novel variants of low or rare frequency from UK10K WGS.....	149
Table 5.5 Novel FBC variants based on expanded discovery (12-way meta-analysis).....	153
Table 5.6 Cohort specific results of top hits from expanded discovery analysis.....	154
Table 5.7 Top hits from a further expanded discovery (18-way meta-analysis)	156
Table 5.8 LD of three putative novel variants in known locus.....	157
Table 5.9 Putative causal variants based on fine mapping	161
Table 5.10 Rare variants aggregation tests based top hits for FBC traits	164
Table 6.1 GWAS studies of CRP.....	177
Table 6.2 Characteristics of participating cohorts	181
Table 6.3 Novel associations of CRP from expanded discovery meta-analysis.....	187

Table 6.4 Cohort specific results of novel associations from expanded discovery	188
Table 6.5 LD between novel and known variants in <i>HIST1H3G</i>	191
Table 6.6 Putative causal variants based on fine mapping	193

LIST OF FIGURES

Figure 1.1 Established and new/emerging risk factors for CVD.....	24
Figure 1.2 The cardiovascular disease continuum.....	25
Figure 1.3 The allelic spectrum of human disease predisposition.....	34
Figure 2.1 UK10K WGS samples data production.....	50
Figure 2.2 Evaluation of batch effects and trait distribution.....	53
Figure 2.3 Phenotype harmonization protocol.....	54
Figure 2.4 Power calculation in the UK10K cohorts.....	57
Figure 2.5 Flow of step-wise conditional analysis.....	61
Figure 3.1 imputation evaluation workflow.....	79
Figure 3.2 Imputation performance for different reference panels and strategies.....	86
Figure 3.3 Illustration of reference states (haplotypes) copied by IMPUTE2.....	87
Figure 3.4 Performance of combining UK10K and 1000GP panels.....	89
Figure 4.1 Lipids loci overlap between candidate gene studies and GWAS.....	101
Figure 4.2. Single point association results of lipids on WGS samples.....	113
Figure 4.3 Association results of 14-way meta-analysis of the four main lipid traits.....	119
Figure 4.4 Regional plots of two loci with replicated novel associations.....	122
Figure 4.5 Number of putative causal variants within fine-mapped loci.....	123
Figure 4.6 QQ plots of SKAT tests for lipids.....	127
Figure 4.7 Rare variants aggregation test results for lipids.....	128
Figure 4.8 Regional plot of SKAT-O locus <i>EGF-ELOVL6</i>	129
Figure 4.9 Statistical power and novel variants from single marker analysis.....	132
Figure 5.1 Association results for WGS based samples for FBC traits.....	148
Figure 5.2 Results for 12-way meta-analysis.....	152
Figure 5.3 Regional plots of two known loci with putative novel variants.....	158
Figure 5.4 Regional plots of top hits from 18-way meta-analysis.....	159
Figure 5.5 Rare variants aggregation test results for FBC traits.....	163
Figure 5.6 Regional plots of <i>RHBDL2</i>	165
Figure 5.7 eSNPs associated with host response to TB and Malaria.....	167

Figure 5.8 Statistical power and novel variants from single marker analysis	171
Figure 6.1 Association Results of CRP based on WGS samples.....	183
Figure 6.2 Single marker association results of CRP from expanded meta-analysis	186
Figure 6.3 Regional plots of two novel associations of CRP	190
Figure 6.4 Rare variants aggregation test results for CRP.....	194
Figure 6.5 Statistical power and novel variants from single marker analysis	197
Figure 7.1 Allelic spectrum for single marker association results in UK10K.....	201
Figure 7.2 QQ plot of association tests for 31 UK10K core traits.....	202

LIST OF ABBREVIATIONS

1000GP	1000 Genomes Project
ADH	Autosomal dominant hypercholesterolemia
ALSPAC	Avon Longitudinal Study of Parents and Children
Apo-A1	apolipoprotein A-I
Apo-B	apolipoprotein B
Apo-E	apolipoprotein E
AMD	age-related macular degeneration
BF	Bayes' factor
BGI	Beijing Genomics Institute
BP	blood pressure
CAD	coronary artery disease
CBR	Cambridge BioResource
CHD	coronary heart disease
CNV	copy number variation
CKD	chronic kidney disease
CRP	C-reactive protein
CVD	cardiovascular disease
DALYs	disability-adjusted life years
DHS	DNaseI hypersensitive sites
EAF	effect allele frequency
EMR	electronic medical records
ERFC	Emerging Risk Factors Collaboration
FHS	Framingham Heart Study
FVG	Friuli Venezia Giulia
GWAS	genome-wide association studies
HDL	high-density lipoprotein
HGB	hemoglobin
HELIC	HELlenic Isolated Cohorts study
HMM	hidden markov model
HWE	hardy-weinberg equilibrium
IBD	identify by descent

IBS	identify by state
InDel	insertion/deletion polymorphism
INGI	Italian Network of Genetic Isolates
LD	linkage disequilibrium
LDL	low-density lipoprotein
LMT	lipid modification therapies
LoF	loss of function
LOLIPOP	London Life Sciences Population study
LURIC	Ludwigshafen Risk and Cardiovascular Health
MAF	minor allele frequency
HGB	haemoglobin
MCH	mean corpuscular hemoglobin
MCHC	mean corpuscular hemoglobin concentration
MCV	mean cell volume
MDS	multidimensional scaling
MI	myocardial infarction
MR	mendelian randomisation
OR	odds ratio
PCA	principle component analysis
PCV	packed cell volume
PLT	platelet count
PROCARDIS	Precocious Coronary Artery Disease Study
QC	quality control
RBC	red blood cell
RCT	reverse cholesterol transport
SKAT	sequence kernel association test
SKAT-O	sequence kernel association test - optimized
SNP	single nucleotide polymorphism
SNV	single nucleotide variation
TC	total cholesterol
TFBS	transcription factor binding sites
TG	triglycerides
TSS	transcription start site

TwinsUK	UK Adult Twin Registry
UK10K	10,000 UK genome sequencing project
UTR	untranslated regions
VB	Val Borbera
WBC	white blood cell
WGS	Whole Genome Sequencing
WTCCC	Wellcome Trust Case Control Consortium
WTSI	Wellcome Trust Sanger Institute

PUBLICATIONS ARISING FROM THIS DISSERTATION

- * *Co-first author*
 - *For papers with more than 10 authors, my name is listed together with the first 3 and the last 3 authors. When there are more than 3 co-starred first-authors, all of them are listed.*
1. Gormley P*, Downes K*, **Huang J***, Kettunen J, Aki S, ..., Palotie A, Ripatti S, Soranzo N. A polygenic panel of platelet-associated SNPs is associated with risk of incident ischaemic stroke. (*submitted*)
 2. Walter K*, Min M*, **Huang J***, Lucy Crooks*, ..., Timpson NJ, Durbin R, Soranzo N. The UK10K project: rare variants in health and disease. (*under revision*)
 3. **Huang J**, Howie B, Memari M, ..., Timpson NJ, Marchini J, Soranzo N, UK10K Project. A reference panel of 3,781 genomes from the UK10K Project increases imputation performance over the 1000 Genomes Project. *Nature Communications*. (*accepted*)
 4. Taylor P, Porcu E, Chew S, ... **Huang J**, ..., Soranzo N, Timpson NJ, Wilson S, the UK10K Consortium. Whole genome sequence based analysis of thyroid function. *Nature Communications*. 2015 Mar 6;6:5681
 5. Timpson NJ, Walter K, Min JL, ..., **Huang J**, ..., Humphries SE, Zeggini E, Soranzo N; UK10K consortium members. A novel low-frequency variant near APOC3 is associated with plasma triglyceride and VLDL levels in Europeans. *Nature Communications*. 2014 Sep 16;5:4871
 6. O'Connell J, Gurdasani D, Delaneau O, ..., **Huang J**, ..., Soranzo N, Sandhu MS, Marchini J. A general approach for haplotype phasing across the full spectrum of relatedness. *PLOS Genetics* 2014 Apr 17;10(4):e1004234

PUBLICATIONS ARISING ELSEWHERE (from 2012-01 to 2015-01)

- * *Co-first author*
 - *For papers with more than 10 authors, my name is listed together with the first 3 and the last 3 authors. When there are more than 3 co-starred first-authors, all of them are listed.*
1. Baumert J*, **Huang J***, McKnight B*, Sabater-Lleal M*, Steri M*, ..., Strachan DP, Peters A, Smith NL. No evidence for genome-wide interactions on plasma fibrinogen by smoking, alcohol consumption and body mass index: results from meta-analyses of 80,607 subjects. *PLoS One*. December 31, 2014 DOI: 10.1371
 2. Shin SY, Fauman EB, Petersen AK, ..., **Huang J**, ..., Kastenmüller G, Spector TD, Soranzo N. An atlas of genetic influences on human metabolism. *Nature Genetics* 2014 Jun;46(6):543-50
 3. Han B, Luo H, Raelson J, **Huang J**, Li Y, Tremblay J, Hu B, Qi S, Wu J. TGFBI (BIG-H3) is a diabetes risk gene based on mouse and human genetic studies. *Hum Mol Genet*. 2014 Apr 11
 4. **Huang J**, Huffman JE, Yamkauchi M, ..., Lowenstein CJ, Strachan DP, O'Donnell CJ; CHARGE Consortium Hemostatic Factor Working Group. Genome-wide association study for circulating tissue plasminogen activator levels and functional follow-up implicates endothelial STXBP5 and STX2. *Arteriosclerosis, Thrombosis, and Vascular Biology*. 2014 Feb 27
 5. Sabater-Lleal M*, **Huang J***, Chasman D*, Naitza S*, Dehghan A*, ..., Strachan DP, Hamsten A, O'Donnell CJ. Multiethnic meta-analysis of genome-wide association studies in >100 000 subjects identifies 23 fibrinogen-associated Loci but no strong evidence of a causal association between circulating fibrinogen and cardiovascular disease. *Circulation*. 2013 Sep 17;128(12):1310-24.
 6. **Huang J**, Liu Y, Welch R, Willer C, Hindorff LA, Li Y. WikiGWA: an open platform for collecting and using genome-wide association (GWA) results. *European Journal of Human Genetics*. 2013 Apr;21(4):471-3
 7. O'Seaghda CM, Wu H, Yang Q, ..., **Huang J**, ..., Bonny O, Fox CS, Bochud M. Meta-analysis of genome-wide association studies identifies six new loci for serum calcium concentrations. *PLoS Genetics*. 2013 Sep;9(9):e1003796
 8. Kleber ME, Seppälä I, Pilz S, ..., **Huang J**, ..., Lehtimäki T, März W, Meitner A. Genome-wide association study identifies three genomic loci significantly associated with

serum levels of homoarginine – The AtheroRemo Consortium. *Circulation Cardiovascular Genetics*. 2013 Sep 18

9. McGrath LM, Cornelis MC, ..., **Huang J**, ..., Sullivan P, Perlis RH, Smoller JW. Genetic predictors of risk and resilience in psychiatric disorders: a cross-disorder genomewide association study of functional impairment in major depressive disorder, bipolar disorder, and schizophrenia. *American Journal of Medical Genetics Part B: Neuropsychiatric Genetics*. 2013 Sep 13
10. ALSGEN Consortium, Ahmeti KB, Ajroud-Driss S, ..., **Huang J**, ..., Veldink JH, Yang Y, Zheng JG. Age of onset of amyotrophic lateral sclerosis is modulated by a locus on 1p34.1. *Neurobiology of Aging*. 2013 Jan;34(1):357.e7-19
11. Chen M-H*, **Huang J***, Chen W-M, Larson MG, Fox CS, Vasani RS, Seshadri S, O'Donnell CJ, Yang Q. Using family-based imputation in genome-wide association studies with large complex pedigrees: the Framingham Heart Study. *PLoS ONE*. 2012;7(12):e51589. (*co-first author)
12. **Huang J**, Sabater-Lleal M, Asselbergs FW, ..., Liu Y, O'Donnell CJ, Hamsten A. Genome-wide association study for circulating levels of plasminogen activator inhibitor-1 (PAI-1) provides novel insights into the regulation of PAI-1. *Blood*. 2012 Dec 6;120(24):4873-81
13. **Huang J**, Ellinghaus D, Franke A, Howie B, Li Y. 1000 Genomes-based imputation identifies novel and refined associations for the Wellcome Trust Case Control Consortium phase 1 data. *European Journal of Human Genetics*. 2012 Jul;20(7):801-5
14. Bis JC, DeCarli C, Smith AV, ..., **Huang J**, ..., Launer LJ, Ikram MA, Seshadri S; Cohorts for Heart and Aging Research in Genomic Epidemiology Consortium. Common variants at 12q14 and 12q24 are associated with hippocampal volume. *Nature Genetics*. 2012 Apr 15;44(5):545-51.
15. Willour VL1, Seifuddin F, Mahon PB, ..., **Huang J**, ..., Gurling H, Purcell S, Smoller JW, Craddock N, DePaulo JR Jr, Schulze TG, McMahon FJ, Zandi PP, Potash JB. A genome-wide association study of attempted suicide. *Mol Psychiatry*. 2012 Apr;17(4):433-44.

