

Identifying and modelling genes that are associated with rare developmental disorders



Keren Jacqueline Carss
Wellcome Trust Sanger Institute
Queens' College
University of Cambridge

This dissertation is submitted for the degree of
Doctor of Philosophy
August 2014

Declaration

I declare that this dissertation describes my own, original work. It only includes work done in collaboration where clearly and specifically indicated in the text. No part of this dissertation has been submitted for any other degree, diploma, or qualification at any university or institution. This dissertation does not exceed 60,000 words.

Keren Jacqueline Carss

August 2014

Abstract

While individually rare, collectively developmental disorders are common, affecting around 3% of live births in the UK. The aetiology of these disorders often includes a genetic component. Next generation sequencing provides a powerful tool with which to identify variants that cause rare developmental disorders. This dissertation describes three distinct projects in which next generation sequencing was used for this purpose, along with statistical or functional follow-up approaches.

A cohort of 30 fetuses with a diverse range of structural abnormalities, along with their parents, was exome sequenced. I analysed these data to identify rare, high quality, coding variants consistent with a *de novo* or recessive inheritance model. I investigated several methods of variant interpretation, including manual and computational methods, and found causative variants for 10% of the cohort. These results suggest that next generation sequencing is a promising method for prenatal genetic diagnostics.

As part of the UK10K project, 996 patients with moderate to severe intellectual disability (ID) underwent sequencing of 565 known or candidate ID-associated genes. I developed and implemented a pipeline to identify likely causative loss of function (LOF) variants through extensive quality filtering. From these data, causative variants were identified for ~14% of the cohort, and the novel ID-associated gene *SETD5* was identified. Next, I performed a series of case-control enrichment analyses to evaluate the contribution of different classes of possibly pathogenic variants. Patients with ID had a significant enrichment of both LOF and missense variants in known ID-associated genes, compared to controls with non-syndromic congenital heart defects.

One strategy to investigate the consequences of a potentially pathogenic variant is to inhibit expression of the gene in an appropriate animal model, and assess the extent to which aspects of the human phenotype are recapitulated. I applied this technique to two genes identified from the UK10K project as likely to be associated with dystroglycanopathy, a subtype of muscular dystrophy. I inhibited the expression of both genes, *B3GALNT2* and *GMPPB*, in zebrafish embryos using morpholino oligonucleotides. The phenotype of both models mimicked several aspects of the human phenotype including morphological defects such as microphthalmia and hydrocephalus, structural defects of the tissue such as disordered muscle fibres, and the precise molecular defect, which is hypoglycosylation of α -dystroglycan.

Acknowledgements

Throughout my PhD, I have been involved in several large collaborative projects, so there are a great many people who I would like to thank. My first and biggest thank you is to my primary supervisor, Dr Matthew Hurlles. He provided me with the opportunity to work on several exciting and fruitful projects. He also provided the perfect balance of support and freedom, along with a constant stream of outstanding ideas, and encouragement when I needed it. Thank you also to my secondary supervisor, Dr Derek Stemple. He introduced me to the zebrafish, and gave me valuable support throughout my time working in the zebrafish laboratory. Thank you to the additional members of my thesis committee: Dr Helen Firth and Dr David Adams, for guidance and advice throughout the last four years.

Thank you to the following clinical collaborators at the University of Birmingham: Prof Eamonn Maher, Prof Mark Kilby, Dr Dominic McMullan, and Dr Sarah Hillman. I very much enjoyed working on the abnormal fetal development project, and I learned a lot from it. Thank you to Dr Vijaya Parthiban, Dr Alejandro Sifrim and Dr Damian Smedley for running the CoNVex, eXtasy, and PhenoDigm programs respectively during this project.

I am grateful to the UK10K consortium, particularly the rare disease group, for the opportunity to be involved in two exciting UK10K projects. Thank you to Dr Lucy Raymond for the opportunity to work on the UK10K intellectual disability cohort, in an exciting and fruitful collaboration. I also thank the many other people involved in this project, most importantly Dr Detelina Grozeva, Dr Olivera Spasic-Boskovic, and Dr James Floyd.

Thank you to Prof Francesco Muntoni for giving me the opportunity to work on the UK10K dystroglycanopathy project. Also to other members of this project including Dr Elizabeth Stevens, Dr Silvia Torelli, Dr Sebahattin Cirak, and Dr Reghan Foley. I owe a particularly large thank you to Dr Yung-Yao Lin, who provided a huge amount of support, supervision, and training to me during this project. Thank you also to Dr Sebastian Gerety for help and support with designing zebrafish experiments.

Thank you to all patients who participate in research studies, without whose generosity none of this work would have been possible.

I am very grateful to the Wellcome Trust for generously funding my PhD. I have been fortunate to do my PhD at the Wellcome Trust Sanger Institute, which provides not only

a very stimulating environment, but also unparalleled facilities. I therefore thank the many people who work in the various pipelines and services teams at the Wellcome Trust Sanger Institute, including the sample reception and sequencing teams, the GAPI team, Human Genome Informatics, and the many other people who keep the laboratory and computational facilities running. I am grateful to Dr Alex Bateman, Dr Julian Rayner, Dr Annabel Smith, Christina Hedberg-Delouka, and Carol Dunbar for practical and administrative support.

I thank all past and present members of teams 29 and 31 for help with experiments or coding, and for many fun and enlightening conversations over the years. I would like to especially thank Dr Saeed Al Turki, Dan King, and Dr Vijaya Parthiban for their great generosity in sharing their scripts, skills, and time with me. I am also particularly grateful to Dr Ana Cvejic and Dr Jovana Serbanovic-Canic for teaching me basic experimental techniques in zebrafish when I first started.

On a more personal note, I thank the many additional scientists who have encouraged, mentored, and inspired me throughout my career, both directly and by example. These include Prof Graham Anderson, Dr Kevin O'Shaughnessy, Dr Annie Mercier Zuber, Dr Darren Logan, Dr Inês Barroso, Dr Eleanor Wheeler, Dr Eleftheria Zeggini, Prof Jane Worthington, and Prof Eamonn Maher. Last but not least I thank my friends and family, especially the fabulous PhD10, the wonderful Norman and Carss families, and my endlessly supportive husband.

Publications

Publications arising from work associated with this thesis:

- Mackie FL, **Carss KJ**, Hillman SC, Hurles ME, & Kilby MD. Exome sequencing in fetuses with structural malformations. [Review article] *Journal of Clinical Medicine*. 2014 July, 3(3), 747-762.
- Grozeva D, **Carss KJ**, Spasic-Broskovic O, Parker M, Archer H, Firth HV, *et al*. *De novo* mutations in *SETD5* cause intellectual disability and associated features of 3p25 microdeletion syndrome. *American Journal of Human Genetics*. 2014 April, 3;94(4):618-24.
- **Carss KJ**, Hillman SC, Parthiban V, McMullan DJ, Maher ER, Kilby MD & Hurles ME. Exome sequencing improves genetic diagnosis of structural fetal abnormalities revealed by ultrasound. *Human Molecular Genetics*. 2014 June, 15,23(12):3269-77.
- **Carss KJ***, Stevens E*, Foley AR, Cirak S, Riemersma M, Torelli S, *et al*. Mutations in *GDP-Mannose pyrophosphorylase B* cause congenital and limb-girdle muscular dystrophies associated with hypoglycosylation of α -dystroglycan. *American Journal of Human Genetics*. 2013 July, 11;93(1):29-41.
- Stevens E*, **Carss KJ***, Cirak S, Foley AR, Torelli S, Willer T, *et al*. Mutations in *B3GALNT2* cause congenital muscular dystrophy and hypoglycosylation of α -dystroglycan. *American Journal of Human Genetics*. 2013 March, 7;92(3):354-65.

* Jointly contributing authors.

Table of Contents

1	Introduction.....	1
2	Exome sequencing improves genetic diagnosis of structural fetal abnormalities	5
2.1	Introduction	5
2.1.1	The impact and causes of fetal structural abnormalities	5
2.1.2	Current techniques for prenatal genetic diagnosis.....	6
2.1.3	Next generation sequencing.....	9
2.1.4	Variant prioritisation strategies	10
2.1.5	Prenatal next generation sequencing: proof of concept.....	12
2.1.6	Aims, context, and colleagues.....	12
2.2	Methods	14
2.2.1	Cohort.....	14
2.2.2	Exome sequencing	15
2.2.3	VCF file merging, annotation, and quality control	15
2.2.4	Identification of <i>de novo</i> SNVs and indels	16
2.2.5	Identification of inherited recessive and X-linked SNVs and indels.....	17
2.2.6	Identification of CNVs	17
2.2.7	Sanger sequencing	18
2.2.8	Interpretation of variants	18
2.3	Results.....	21
2.3.1	The exome sequencing data are of high quality	21
2.3.2	There is a mean of 1.13 validated <i>de novo</i> SNVs or indels per fetus.....	26
2.3.3	There are three candidate <i>de novo</i> or X-linked copy number variants.....	28
2.3.4	There is a mean of 13 candidate genes with inherited recessive or X-linked variants per fetus, in the preliminary round of analysis	30
2.3.5	<i>De novo</i> SNVs in <i>FGFR3</i> and <i>COL2A1</i> are highly likely to be causal.....	30
2.3.6	<i>De novo</i> SNVs in <i>NF1</i> and <i>SMARCC2</i> are possibly causal.....	32
2.3.7	Two unrelated fetuses with no clear clinical overlap have <i>de novo</i> SNVs in <i>PARD3B</i>	33
2.3.8	A <i>de novo</i> deletion that overlaps with <i>OFD1</i> is highly likely to be causal ..	34
2.3.9	Inherited recessive or X-linked SNVs in five fetuses are possibly causal, in the preliminary round of analysis.....	34
2.3.10	The variant prioritisation program eXtasy identifies 36 possibly causal variants, with an enrichment of <i>de novo</i> mutations	36

2.3.11	The variant prioritisation program PhenoDigm identifies possibly causal variants in 18 genes	39
2.3.12	There is a degree of overlap between the variants identified as possibly causal by the three different prioritisation methods.....	41
2.3.13	The continuing need for manual curation	42
2.3.14	Inherited recessive or X-linked SNVs in five fetuses are possibly causal, in the final round of analysis	42
2.3.15	The estimated diagnostic yield of this study is 10%.....	44
2.4	Discussion.....	47
2.4.1	Summary	47
2.4.2	The diagnostic yield in context	47
2.4.3	Comparison of variant interpretation methods.....	48
2.4.4	The ethics of next generation sequencing for prenatal genetic diagnosis	50
2.4.5	Next generation sequencing is the future of prenatal genetic diagnostics	52
3	Case-control analysis of 565 known and candidate intellectual disability-associated genes.....	55
3.1	Introduction.....	55
3.1.1	The impact of intellectual disability	55
3.1.2	Discovery of intellectual disability-associated genes	57
3.1.3	Biology of intellectual disability-associated genes	59
3.1.4	Case-control enrichment analysis of rare variants.....	62
3.1.5	Aims, context, and colleagues.....	65
3.2	Methods.....	67
3.2.1	Samples, sequencing, and quality control	67
3.2.2	Annotation and filtering pipeline	67
3.2.3	Principal component analysis.....	68
3.2.4	Cohort allelic sums test.....	69
3.3	Results.....	72
3.3.1	Targeted resequencing of 565 intellectual disability-associated genes in cases and controls	72
3.3.2	The sequencing data are of good quality	72
3.3.3	There is no substantial difference in population structure between the intellectual disability and congenital heart disease cohorts.....	74
3.3.4	14% of intellectual disability patients have a likely causative variant in a sequenced intellectual disability-associated gene	76
3.3.5	<i>SETD5</i> is a novel intellectual disability-associated gene	77
3.3.6	Individuals with intellectual disability have an enrichment of loss of function variants in sequenced ID-associated genes, compared to controls	78

3.3.7	In known ID-associated genes on the X chromosome, unique missense variants tend to be more damaging in ID patients than controls.	81
3.3.8	Evidence for an enrichment of unique, predicted damaging, missense variants in sequenced ID-associated genes in the ID cohort	83
3.4	Discussion.....	86
3.4.1	Summary	86
3.4.2	Loss-of-function of the histone methyltransferase gene <i>SETD5</i> is probably responsible for the cardinal features of 3p25 microdeletion syndrome	86
3.4.3	Insights from case-control enrichment analyses.....	89
3.4.4	Limitations of this study.....	91
3.4.5	Further work.....	92
4	Modelling dystroglycanopathy in zebrafish embryos by knockdown of <i>B3GALNT2</i> and <i>GMPPB</i>	94
4.1	Introduction.....	94
4.1.1	The phenotypic spectrum of dystroglycanopathy.....	94
4.1.2	Dystroglycan structure, function and glycosylation	95
4.1.3	Known dystroglycanopathy-associated genes.....	96
4.1.4	Frequency of variants, and genotype-phenotype correlations	101
4.1.5	Zebrafish models of genetic disease	102
4.1.6	Animal models of dystroglycanopathy	104
4.1.7	Aims, context, and colleagues.....	105
4.2	Materials and methods	107
4.2.1	Sequencing of clones.....	107
4.2.2	Reverse transcription polymerase chain reaction	107
4.2.3	Design and injection of morpholino oligonucleotides	107
4.2.4	Generation of green fluorescent protein-tagged RNA.....	108
4.2.5	Immunofluorescence staining.....	109
4.2.6	Evans blue dye assay	109
4.2.7	Immunoblotting	109
4.3	Results.....	110
4.3.1	<i>B3GALNT2</i> and <i>GMPPB</i> are conserved with their zebrafish orthologues 110	
4.3.2	Expression of <i>b3galnt2</i> and <i>gmppb</i> throughout early zebrafish development	114
4.3.3	Finding the optimal morpholino oligonucleotide and dose for <i>b3galnt2</i> and <i>gmppb</i>	115
4.3.4	Morpholino oligonucleotides reduce the expression of <i>b3galnt2</i> and <i>gmppb</i> 116	

4.3.5	<i>b3galnt2</i> morphants have gross morphological defects including hydrocephalus and impaired motility	117
4.3.6	<i>b3galnt2</i> morphants have muscle defects including gaps in the myosepta and lesions between fibres	119
4.3.7	<i>b3galnt2</i> morphants have hypoglycosylated α -dystroglycan.....	120
4.3.8	Coinjection with wildtype RNA fails to rescue the <i>b3galnt2</i> morphant phenotype	120
4.3.9	<i>gmppb</i> morphants have gross morphological defects including microphthalmia and impaired motility	122
4.3.10	<i>gmppb</i> morphants have muscle defects including disordered fibres, incomplete myosepta, and interfibre spaces.....	124
4.3.11	<i>gmppb</i> morphants have hypoglycosylated α -dystroglycan	124
4.4	Discussion.....	126
4.4.1	Summary	126
4.4.2	Phenotypic rescue	126
4.4.3	The function of B3GALNT2	127
4.4.4	The function of GMPPB	128
4.4.5	Zebrafish phenotypes in context	130
4.4.6	Technical limitations of this study	133
4.4.7	Future research.....	133
5	Discussion	136
6	References	140
7	Appendices	179

List of Figures

Figure 2 - 1: Matrix showing categories of phenotypes in the cohort of fetuses with structural abnormalities	15
Figure 2 - 2: Decision tree for classifying candidate genes into three categories.....	20
Figure 2 - 3: Target coverage of exome sequencing reads by sample.	21
Figure 2 - 4: Quality control metrics for single nucleotide variants.....	24
Figure 2 - 5: Quality control metrics for indels.	25
Figure 2 - 6: Log ₂ ratios of candidate CNVs in fetuses with structural abnormalities. ..	29
Figure 2 - 7: Pedigree of trio 23, showing Sanger sequencing of <i>de novo</i> mutation in <i>FGFR3</i>	31
Figure 2 - 8: Pedigree of trio 20, showing Sanger sequencing of <i>de novo</i> mutation in <i>COL2A1</i>	32
Figure 2 - 9: Venn diagram showing overlap between the genes prioritised by each of the three methods.	41
Figure 3 - 1: Quality control metrics for the UK10K targeted resequencing study.....	73
Figure 3 - 2: Principal component analysis.....	75
Figure 3 - 3: Classes of variant identified through the R filtering pipeline.	76
Figure 3 - 4: Facial appearance of individuals with <i>SETD5</i> mutations.	78
Figure 3 - 5: Patients with intellectual disability have an enrichment of loss of function variants in sequenced intellectual disability-associated genes compared to controls.	79
Figure 3 - 6: In known ID-associated genes on the X chromosome, unique missense variants are predicted to be more damaging in ID patients than controls.....	83
Figure 3 - 7: In known ID-associated genes on the autosomes, unique missense variants are not predicted to be more damaging in ID patients than controls.....	83
Figure 4 - 1: Model of α -DG interactions.	96
Figure 4 - 2: The function of GMPPB in glycosylation pathways.	106
Figure 4 - 3: Protein alignment showing conservation of B3GALNT2.	112
Figure 4 - 4: Protein alignment showing conservation of GMPPB.	113
Figure 4 - 5: Reverse transcription PCR shows expression of <i>b3galnt2</i> and <i>gmppb</i> throughout early zebrafish development.....	114
Figure 4 - 6: Coinjection of <i>p53</i> MO rescues the neurodegeneration induced by <i>b3galnt2</i> MO.	115
Figure 4 - 7: The <i>b3galnt2</i> MO inhibits expression of recombinant GFP-tagged <i>b3galnt2</i> RNA.....	116
Figure 4 - 8: The <i>gmppb</i> splice blocking MO disrupts RNA splicing.	117

Figure 4 - 9: <i>b3galnt2</i> knockdown zebrafish embryos have muscle defects and hypoglycosylated α -DG at 48 hpf.	118
Figure 4 - 10: Coinjection with wildtype human <i>B3GALNT2</i> RNA fails to rescue the <i>b3galnt2</i> morphant phenotype.....	121
Figure 4 - 11: Coinjection with wildtype zebrafish <i>b3galnt2</i> RNA fails to rescue the <i>b3galnt2</i> morphant phenotype.....	122
Figure 4 - 12: <i>gmppb</i> knockdown zebrafish embryos have morphological defects, damaged muscle, and hypoglycosylated α -DG at 48 hpf.....	123
Figure 4 - 13: α -DG of <i>gmppb</i> morphants is hypoglycosylated relative to wildtype embryos.	125
Figure 4 - 14: B3GALNT2 catalyses the synthesis of the trisaccharide GalNAc β 1-3-GlcNAc- β 1,4-Man, which is required for laminin binding.	128

List of Tables

Table 2 - 1: Exome sequencing coverage and quality control metrics.	23
Table 2 - 2: Validated <i>de novo</i> SNVs in fetuses with structural abnormalities.	27
Table 2 - 3: Candidate CNVs in fetuses with structural abnormalities.	28
Table 2 - 4: Candidate genes identified as possible causal by eXtasy.	38
Table 2 - 5: Candidate genes identified as possibly causal by PhenoDigm.	40
Table 2 - 6: Summary of all candidate genes identified in 30 fetuses with structural abnormalities.....	46
Table 3 - 1: Functional classes of ID-associated genes.	59
Table 3 - 2: List of 204 sequenced intellectual disability-associated genes that are known.	70
Table 3 - 3: List of 361 sequenced intellectual disability-associated genes that are candidates.....	71
Table 3 - 4: Candidate genes with the highest number of LOF variants.	77
Table 3 - 5: Enrichment of unique LOF variants in the ID cohort, split by category.....	80
Table 3 - 6: Enrichment of unique, predicted damaging, missense variants in the ID cohort, split by category.	84
Table 4 - 1: Primers used to analyse <i>b3galnt2</i> and <i>gmppb</i> expression in early zebrafish development, and splicing disruption in <i>gmppb</i> morphants.	107
Table 4 - 2: Morpholino oligonucleotide sequences and predicted effects.	108
Table 4 - 3: Percentage identity of B3GALNT2 orthologues of five diverse eukaryotic species with human B3GALNT2.....	111
Table 4 - 4: Percentage identity of GMPPB orthologues of five diverse eukaryotic species with human GMPPB.	111
Table 4 - 5: <i>b3galnt2</i> morphants are significantly more likely to have hydrocephalus than wildtype embryos.	119
Table 4 - 6: Phenotypic comparison of zebrafish dystroglycanopathy models.....	132

List of Appendices

Appendix 1: Primer sequences for Sanger sequencing of variants that passed validation.....	180
Appendix 2: Inherited recessive and X-linked SNPs and indels that pass filters, in fetuses with structural abnormalities (preliminary round of analysis).	188
Appendix 3: High-quality, rare, coding, inherited recessive and X-linked SNPs and indels (final round of analysis).	192

List of Abbreviations

aCGH	Array comparative genomic hybridisation
α -DG	α -dystroglycan
ASD	Autism spectrum disorder
β -DG	β -dystroglycan
BioGPS	Biology Gene Portal System
bp	Base pair
CAST	Cohort allelic sums test
CDG	Congenital disorder of glycosylation
cfDNA	Cell-free DNA
CHD	Congenital heart disease
CK	Creatine kinase
CMD	Congenital muscular dystrophy
CNS	Central nervous system
CNV	Copy number variant
CRISPR	Clustered regularly interspaced short palindromic repeat
DDD	Deciphering developmental disorders
DDG2P	Developmental Disorder Gene2Phenotype
DGC	Dystrophin-glycoprotein complex
DNA	Deoxyribonucleic Acid
DoI-P-Man	Dolichol phosphate mannose
EBD	Evans blue dye
ECM	Extracellular matrix
ENU	<i>N</i> -ethyl- <i>N</i> -nitrosourea
ER	Endoplasmic reticulum
ESP	Exome Sequencing Project
FCMD	Fukuyama-type CMD
FISH	Fluorescence in situ hybridisation
FORGE	Finding of rare disease genes
GalNAc	N-acetylgalactosamine
Gb	Gigabases
GDP	Guanosine diphosphate
GFP	Green fluorescent protein
GlcNAc	N-acetylglucosamine
GPI	Glycosylphosphatidylinositol
GWAS	Genome-wide association study
HDL-C	High-density lipoprotein cholesterol
HMQ	High mapping quality
Hpf	Hours post fertilisation
HPO	Human phenotype ontology
IC	Information Content
ID	Intellectual disability

IEM	Inborn error of metabolism
IGV	Integrative Genomics Viewer
IKMC	International knockout mouse consortium
Indel	Insertion deletion
IQ	Intelligence quotient
Kb	Kilobase
LD	Linkage disequilibrium
LGMD	Limb girdle muscular dystrophy
LOF	Loss of function
Mb	Megabase
MEB	Muscle-eye-brain disease
MO	Morpholino oligonucleotide
MPO	Mammalian phenotype ontology
MRI	Magnetic resonance imaging
NGS	Next generation sequencing
NHGRI	National Human Genome Research Institute
NHLBI EVS	National Institute of Health: Heart, Lung, and Blood Institute exome variant server
NIPT	Non-invasive prenatal testing
OMIM	Online Mendelian Inheritance in Man
PAGE	Prenatal Assessment of Genomes and Exomes
PCA	Principal component analysis
PCR	Polymerase chain reaction
PKU	Phenylketonuria
PSD	Postsynaptic density
PTR	Primary target region
QF-PCR	Quantitative fluorescent polymerase chain reaction
RD	Retinal dystrophy
RNA	Ribonucleic Acid
RT-PCR	Reverse transcription polymerase chain reaction
SB	Splice blocking
SimJ	Jaccard Index
SKAT	Sequence kernel association test
SNP	Single nucleotide polymorphism
SNV	Single nucleotide variant
TALEN	Transcription activator-like effector nuclease
TB	Translation blocking
UTR	Untranslated region
VCF	Variant call format
VEP	Variant effect predictor
VOUS	Variant of unknown significance
WTSI	Wellcome Trust Sanger Institute
WWS	Walker-Warburg Syndrome
ZFIN	Zebrafish information network
ZFN	Zinc finger nuclease