

# **Uniparental disomy and mosaic structural variation in developmental disorders**



**Daniel Alexander King, M.D.**

**Wellcome Trust Sanger Institute**

**University of Cambridge**

**This dissertation is submitted for the degree of**

**Doctor of Philosophy**

**November 2015**



## DECLARATION OF ORIGINALITY

This dissertation describes my research performed between January 2012 and June 2015 under the supervision of Dr. Matthew Hurles. This work is my own and includes nothing which is the outcome of work done in collaboration except where explicitly stated in the text. It has not been previously submitted for any qualification, and it complies with the prescribed word limit set by the Degree Committee for the Faculty of Biology.

## ABSTRACT

Developmental disorders (DDs) are diseases of impaired childhood development and include congenital anomalies, neurodevelopmental disorders, and abnormalities in growth and behaviour. Determining the genetic causes underlying DD is a major goal of contemporary medical research and the recent entrance of exome sequencing data into the rare-disease field has been transformative in uncovering the importance of *de novo* point mutations as a major source of DD-associated mutations. Recent efforts have successfully harnessed exome sequencing data to detect constitutive copy-number variation, a form of large-scale structural abnormality. However, at the inception of my doctoral work, no software tools had yet been developed to identify, from exome sequencing data, uniparental disomy (UPD), a form of copy-neutral variation, nor large-scale ('structural') abnormalities, which have long been implicated as important contributors to DD. The research underlying this work aimed to fill this void.

This dissertation addresses the development of new software tools, UPDio and MrMosaic, which have extended the diagnostic reach of sequencing data to identify UPD and structural mosaicism, and have been made freely available. Simulation analyses show that these tools can detect the large-scale abnormalities identified by karyotyping or microarray in standard clinical testing. Implementation on nearly 5,000 children with undiagnosed diseases demonstrated that UPD and structural mosaicism are enriched in children with developmental disorders compared with healthy children and suggested that most of the detected abnormalities are likely to be pathogenic. Investigation of the clinical impact of the detected events identified several disease-causing mechanisms, including UPD-associated imprinting and recessive diseases, and genomic disorders associated with large mosaic deletions and duplications.

The five chapters of this dissertation are: 1) an introduction, to describe the context of this doctoral work; 2) a description of UPDio, a new method for detecting uniparental disomy from exome trio data; 3) a burden analysis of mosaic structural variation and the clinical consequences of mosaic structural variants found in children with DD; 4) a description of MrMosaic, a new method for the detection of mosaic structural variation using next generation sequence data; and lastly, 5) a discussion that recapitulates the results of these analyses, describes their limitations, and considers future directions.

## ACKNOWLEDGEMENTS

Great thanks are owed to many. In chronological order, my parents deserve credit, overdue and too understated, for teaching me how to walk, talk, and create (thanks mom) and inspiring my passions in science (thanks dad). Josh and Jason provided a sibling rivalry that was a great cure of laziness and I look forward to one day raising our families together.

One of my first patients in medical school was a middle-aged dad, experiencing new, severe headaches, and died from glioblastoma a few months later. Thank you for allowing me to follow your case as the experience led me to pursue medical research at the NIH. There, Dr. Les Biesecker nurtured my interests in programming, advocated for me, and encouraged me to become a physician scientist. Jamie Teer was a patient and forgiving scripting coach and Larry Singh taught me that statistics is really not so boring at all.

I am sincerely grateful to my Ph.D. supervisor, Matt Hurles, with whose intuition this dissertation research has greatly benefited, and with whose guidance the pleasure of research was conspicuous and reinforcing. One of my favourite phrases of yours is ‘it’s a testable hypothesis’, a maxim demonstrating your perpetual loyalty to evidence and methodical experimentation. I am grateful to have pursued with your mentorship and collaboration science as it evolved in front of us, yielding findings often unexpected and fascinating, and hope that such discovery can continue in partnership for many years to come.

Thank you to everyone in Team 29 for putting up with my cheeky birthday surprises and for your friendship. Saeed, Manu and Raheleh, thanks for being good coffee-break companions. Ray, Alejandro, Jeremy, and Tom, I enjoyed our brainstorming sessions. Thank you to the DDD laboratory and informatics teams for performing so many upstream analyses. The work presented here was not possible without the participation of thousands of children and their families; thank you for joining DDD. Wellcome Trust kindly funded this research. Annabel & Christina, and Carol, much appreciation for organising the logistics for all my academic activities.

Thank you Tanya for always being there for me, and for enabling my life inside and outside the lab these Ph.D. years to be adventurous and fulfilling.

# CONTENTS

<b>1 INTRODUCTION.....</b>	<b>1</b>
1.1 STRATEGIES FOR DETECTING STRUCTURAL VARIATION .....	4
1.1.1 Optical cytogenetics.....	4
1.1.2 Molecular cytogenetics .....	6
1.1.3 DNA sequencing.....	10
1.2 STRUCTURAL VARIATION IN DEVELOPMENTAL DISORDERS .....	13
1.2.1 Copy-number variation in DD .....	13
1.2.2 Copy-neutral loss of heterozygosity (uniparental disomy) in DD .....	17
1.2.3 Mosaic structural rearrangements and DD .....	21
1.3 CLINICAL DIAGNOSTIC TESTING OF DEVELOPMENTAL DISORDERS .....	23
1.3.1 Deciphering Developmental Disorders study .....	25
1.4 SUMMARY .....	26
<b>2 UNIPARENTAL DISOMY.....</b>	<b>27</b>
2.1 PUBLICATION NOTE .....	27
2.2 INTRODUCTION.....	27
2.3 METHODS.....	31
2.3.1 Genotype segregation and statistical analysis.....	31
2.3.2 Samples analysed .....	32
2.3.3 Exome processing .....	32
2.3.4 SNP microarray data processing.....	33
2.3.5 Avoiding positions in copy-number variant regions.....	33
2.3.6 Simulation testing .....	34
2.3.7 Assessing pathogenic variation in samples with UPD events .....	35
2.3.8 Using WTCCC data to estimate UPD in the general population.....	36
2.3.9 Computational performance.....	36
2.3.10 Software availability .....	37
2.4 RESULTS.....	38
2.4.1 Simulations .....	39
2.4.2 Comparing UPD detection software tools .....	41
2.4.3 Implementing quality control of UPD detections .....	46
2.4.4 UPD detections .....	52
2.4.5 Investigating UPD frequency.....	53
2.4.6 Investigating pathogenicity in children with UPD events .....	55
2.4.6.1 UPD chromosome is the dominant source of candidate variant(s) .....	55

2.4.6.2 Non-UPD chromosome is the dominant source of candidate variant(s).....	60
2.4.6.3 Variants with uncertain pathogenicity .....	60
2.5 DISCUSSION.....	67
<b>3 MOSAIC STRUCTURAL VARIATION FROM SNP MICROARRAY .....</b>	<b>72</b>
3.1 PUBLICATION NOTE.....	72
3.2 INTRODUCTION .....	72
3.3 MATERIALS & METHODS .....	77
3.3.1 Description of studies.....	77
3.3.2 Mosaic event detection.....	79
3.3.3 Methods of evaluating of clinical significance .....	79
3.3.4 Exome sequencing.....	80
3.4 RESULTS.....	81
3.4.1 Filtering Strategies for MAD output from DDD & SFHS samples .....	81
3.4.1.1 Managing over-segmentation.....	83
3.4.1.2 Managing constitutive homozygosity & unimodal BAF deflection .....	83
3.4.1.3 Managing constitutive CNVs.....	85
3.4.1.4 Inclusion of aberrant standard deviation of BAFs rescues one mosaic event.....	86
3.4.1.5 Filtering strategies for TEDS and ALSPAC .....	86
3.4.2 Assessing the accuracy of filtering strategies .....	87
3.4.3 Mosaicism Frequency in Cases & Controls using MAD .....	89
3.4.4 Additional detections using triPOD .....	92
3.4.5 Validation experiments to explore tissue distribution.....	96
3.4.6 Clinical Interpretation of Probands with Mosaicism.....	96
3.5 DISCUSSION.....	107
<b>4 MOSAIC STRUCTURAL VARIATION FROM TARGETED AND WHOLE-GENOME SEQUENCING.....</b>	<b>110</b>
4.1 PUBLICATION NOTE.....	110
4.2 INTRODUCTION .....	110
4.3 MATERIALS & METHODS .....	113
4.3.1 MrMosaic .....	113
4.3.2 Simulating Mosaicism.....	119
4.3.3 Description of Samples & Sequencing.....	121
4.3.4 Additional filtering implemented in addition to Mscore quality score .....	123
4.3.5 SNP genotyping chip validation.....	124
4.4 RESULTS.....	125

4.4.1 Simulations .....	126
4.4.2 Detections in Exome Data .....	139
4.4.3 Empirical evaluation of detection of mosaicism from WGS data .....	147
4.5 CLINICAL ASSESSMENT.....	149
4.6 DISCUSSION.....	153
<b>5 DISCUSSION .....</b>	<b>160</b>
5.1 SUMMARY OF FINDINGS .....	160
5.2 IMPLICATIONS .....	161
5.3 LIMITATIONS .....	162
5.3.1 Estimates of prevalence .....	162
5.3.2 Algorithmic .....	163
5.3.3 Number of diagnoses .....	165
5.4 FUTURE WORK.....	166
5.5 AND THEN... ..	169
5.5.1 Achieving a higher fidelity genome.....	169
5.5.2 Having achieved a higher fidelity genome .....	169
5.5.3 Challenges further ahead.....	171
<b>6 REFERENCES.....</b>	<b>173</b>



## PUBLICATIONS

King, D.A., Sifrim, A.S., Fitzgerald, T.W. & Hurles, M.E. Detection of structural mosaicism from targeted and whole-genome sequencing data. *In review*.

King, D.A., Jones, W.D., Crow, Y.J., Dominiczak, A.F., *et al.* Mosaic structural variation in children with developmental disorders. *Human molecular genetics* **24**, 2733-2745 (2015).

Carvalho, C.M.B., Pfundt, R., King, D.A., Lindsay, S.J., *et al.* Absence of heterozygosity due to template switching during replicative rearrangements. *The American Journal of Human Genetics* **96**, 555-564 (2015).

Akawi, N., McRae, J., Ansari, M., Balasubramanian, M., *et al.* Discovery of four recessive developmental disorders using probabilistic genotype and phenotype matching among 4,125 families. *Nature genetics* **47**, 1363-1369 (2015).

Wright, C.F., Fitzgerald, T.W., Jones W.D., McRae, J.F., *et al.* Genetic diagnosis of developmental disorders in the DDD study: a scalable analysis of genome-wide research data. *The Lancet* **385**, 1305-1314 (2015).

TDDD Study. Large-scale discovery of novel genetic causes of developmental disorders. *Nature* **12**, 223-228 (2014).

King, D.A., Fitzgerald, T.W., Miller, R., Canham, N., *et al.* A novel method for detecting uniparental disomy from trio genotypes identifies a significant excess in children with developmental disorders. *Genome research* **24**, 673-687 (2014).

## LIST OF TABLES

TABLE 1-1 .....	25
TABLE 2-1 .....	31
TABLE 2-2 .....	41
TABLE 2-3 .....	53
TABLE 2-4 .....	64
TABLE 2-5 .....	66
TABLE 3-1 .....	75
TABLE 3-2 .....	79
TABLE 3-3 .....	98
TABLE 4-1 .....	119
TABLE 4-2 .....	132
TABLE 4-3 .....	143
TABLE 4-4 .....	147
TABLE 4-5 .....	150
TABLE 4-6 .....	157
TABLE 4-7 .....	159

## LIST OF FIGURES

FIGURE 1-1.....	2
FIGURE 1-2.....	4
FIGURE 1-3.....	5
FIGURE 1-4.....	8
FIGURE 1-5.....	10
FIGURE 1-6.....	14
FIGURE 1-7.....	16
FIGURE 1-8.....	17
FIGURE 1-9.....	20
FIGURE 1-10.....	22
FIGURE 2-1.....	38
FIGURE 2-2.....	40
FIGURE 2-3.....	43
FIGURE 2-4.....	43
FIGURE 2-5.....	45
FIGURE 2-6.....	47
FIGURE 2-7.....	50
FIGURE 2-8.....	51
FIGURE 3-1.....	82
FIGURE 3-2.....	82
FIGURE 3-3.....	83
FIGURE 3-4.....	85
FIGURE 3-5.....	87
FIGURE 3-6.....	88
FIGURE 3-7.....	89

FIGURE 3-8.....	90
FIGURE 3-9.....	91
FIGURE 3-10.....	92
FIGURE 3-11.....	95
FIGURE 3-12.....	96
FIGURE 3-13.....	100
FIGURE 3-14.....	105
FIGURE 4-1.....	115
FIGURE 4-2.....	117
FIGURE 4-3.....	118
FIGURE 4-4.....	122
FIGURE 4-5.....	126
FIGURE 4-6.....	127
FIGURE 4-7.....	128
FIGURE 4-8.....	129
FIGURE 4-9.....	132
FIGURE 4-10.....	134
FIGURE 4-11.....	137
FIGURE 4-12.....	139
FIGURE 4-13.....	141
FIGURE 4-14.....	141
FIGURE 4-15.....	144
FIGURE 4-16.....	145
FIGURE 4-17.....	145
FIGURE 4-18.....	146
FIGURE 4-19.....	149

## TABLE OF ABBREVIATIONS AND ACRONYMS

1000G	1000 Genomes study
aCGH	array comparative genomic hybridisation
ALSPAC	Avon Longitudinal Study of Parents and Children
AUC	area under the curve
BAC	bacterial artificial chromosome
BAF	b allele frequency
B <sub>dev</sub>	BAF deviation
C <sub>dev</sub>	copy-number deviation
CNV	copy number variation
DD	developmental disorders
DDD	Deciphering Developmental Disorders
DECIPHER	Database of genomic variation & Phenotype in Humans using Ensembl Resources
FISH	fluorescent <i>in situ</i> hybridisation
GADA	genome alteration detection analysis (software)
GRC	genome reference consortium
GRCh37	genome reference consortium, human genome reference 37
GWAS	genome-wide association studies
HGP	human genome project
HMM	hidden Markov model
HPO	human phenotype ontology
indels	insertions and deletions
LOH	loss of heterozygosity
LRR	log r ratio
MAD	mosaic alteration detection
MAF	minor allele frequency

Mb	megabases
MrMosaic	Mosaic Rearrangements by Merging Of Sequenced Alleles using their Identity and Counts
OMIM	online Mendelian inheritance in man
QC	quality control
RFLP	restriction fragment length polymorphism
ROH	region of homozygosity
SFHS	Scottish Family Health Study
SNP	single nucleotide polymorphism
SV	structural variation
TEDS	Twins Early Development Study
UK10K	United Kingdom 10,000 Genomes Project
UPD	uniparental disomy
VCF	variant call format
WES	whole-exome sequencing
WGS	whole-genome sequencing