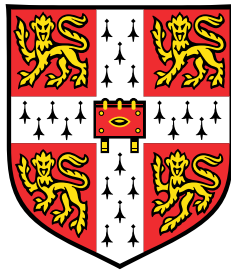


# Lineage tracing of normal human development and childhood cancers



**Tim Coorens**

Wellcome Sanger Institute  
University of Cambridge

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

Clare Hall

September 2020



I would like to dedicate this dissertation to family and friends close by and far away, who have never stopped making me smile and laugh.



## **Declaration**

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is my own work and contains nothing which is the outcome of work done in collaboration with others, except as specified in the text. This dissertation does not exceed 60,000 words in length.

Tim Coorens  
September 2020



# Lineage tracing of normal human development and childhood cancers

Tim Coorens

## Summary

From fertilisation onwards, the cells of the human body continuously experience damage to their genome, either from intrinsic causes or from exposure to mutagens. While the vast majority of DNA damage is repaired and the genome is replicated with extremely high fidelity, cells steadily acquire single nucleotide variants throughout life. Since cells pass these genetic changes on to their descendants, mutations shared between any two cells therefore imply a shared developmental path. In essence, these somatic mutations connect all cells together into one large phylogenetic tree of human development with the zygote at the root.

Reconstructing phylogenies of human development requires readouts of somatic mutations present in single cells. Recently, low-input whole-genome sequencing following laser-capture microdissection has allowed us to reliably call somatic mutations in distinct single-cell derived physiological units, such as colonic crypts and endometrial glands, while retaining spatial information on a microscopic level. In this way, I reconstructed large-scale phylogenies of cells from many different organs of three individuals. These phylogenetic trees recapitulate the early stages of embryonic development and asymmetric cell allocation in the blastocyst, as well as later clonal expansions such as benign prostatic hyperplasia and neoplastic polyp formation.

In a similar way, I also used somatic mutations to investigate the emergence of paediatric cancer, which is thought to be closely linked to aberrations in development. In the context of phylogenetic analyses of tumours, mutations shared between childhood cancers and different normal tissues can shed light on the embryonic lineage of tumours and may reveal the precise juncture at which tumours began to form. Accordingly, I studied the origin of Wilms tumour, the most common childhood cancer of the kidney. I discovered that these tumours often arise from large tissue-resident precursor clones residing in the normal kidney. These embryonal precursors represent an early clonal expansion driven by *H19* hypermethylation.

Lastly, using somatic mutations I discovered that the human placenta is made up of large clonal patches of closely related trophoblast cells. Comparing early embryonic mutations between placental lineages and umbilical cord DNA, which is derived from the inner cell mass, revealed that in approximately half of the cases, a trophectodermal lineage shares no somatic mutations with the umbilical cord. Furthermore, in a quarter of cases, the umbilical

cord is entirely derived from a progenitor later than the zygote. This indicates a natural early segregation between these lineages and a pathway to generate confined placental mosaicism.

This dissertation as a whole provides a new framework to study normal and aberrant human development from whole-genome sequencing. The ability to reconstruct developmental lineages retrospectively can answer fundamental questions about human development and carcinogenesis.



## Acknowledgements

The work described in this thesis would have been impossible without the support of many people. First and foremost, I would like to thank my *de facto* triumvirate of supervisors: Professor Mike Stratton, for his deep insights, endless patience and delightful deliberations on Gibbon's *Decline and Fall* and Procopius' *Secret Histories*; Dr Iñigo Martincorena, for his invaluable analytic mind, steadfast support and imbuing me with a strong affinity for the beta-binomial distribution; Dr Sam Behjati, for his unwavering enthusiasm, endless ideas and the lengthy dialogues about right and wrong in the use of colour in figures.

I would also like to thank the other members of my thesis advisory committee, Dr Peter Campbell, Professor Elizabeth Murchison, Professor Magdalena Zernicka-Goetz, Dr Moritz Gerstung and Dr Nick Goldman, for all the helpful insights, discussions and comments along the road.

I would never have been able to perform this work without the support of the Wellcome Trust Mathematical Genomics and Medicine PhD programme, and I am very grateful to the directors and my fellow students over the years for making this programme so special.

I am especially indebted to the administrative staff, IT team and lab support that make the Sanger Institute run smoothly every day and made this research possible on numerous fronts. A special thanks to Jo Jones, Carl Logan and Wendy McLaughlin, the true unsung heroes of this work, whose impeccable handling of chaotic agendas and enquiries small and large helped me enormously. I am very grateful to all the colleagues I closely worked with over the years, especially Dr Luiza Moore, Dr Taryn Treger, Dr Thomas Oliver, Dr Raheleh Rahbari, Dr Roser Vento-Tormo and Dr Philip Robinson.

I am extremely fortunate to have to thank a small legion of friends, without whom this PhD would not nearly have been so much fun. At the risk of committing the heinous crime of omission, I would like to give a few friends a special thanks for extensively proofreading this dissertation, providing me with much-needed distractions and helping with the crosswords: Jannat Ijaz, Dr Pantelis Nicola, Dr Grace Collord, Dr Alex Cagan, Dr Timothy Butler, Dr Heather Machado, Andrew Lawson, Luke Harvey, Dr Patrick McClanahan, Dr Sophie Allcock, Timo Haber and Natalie Shatashvili.

It almost goes without saying that I want express my deepest gratitude to my family back in the Netherlands, to whom I owe everything and who have supported me every step along the way, especially my parents, Peter and Sandra, my siblings Jodie and Jim, my step-parents, Paul and Angelique, and of course, my grandparents.

And of course I would like to thank you, dear reader, for opening this dissertation and having a look at this work. With all sincerity, I hope you will enjoy this as much as I did.

# Contents

<b>List of Figures</b>	<b>xv</b>
<b>List of Tables</b>	<b>xvii</b>
<b>Nomenclature</b>	<b>xix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 A tree of life in every organism . . . . .	2
1.1.1 A historical perspective . . . . .	2
1.1.2 Early lineage tracing experiments . . . . .	4
1.1.3 Sequencing-based lineage tracing . . . . .	6
1.2 Somatic mutations as natural markers . . . . .	7
1.2.1 Genomic readouts of single cells . . . . .	7
1.2.2 Mutational processes in normal tissues . . . . .	9
1.2.3 Early embryonic mutations and phylogenies . . . . .	11
1.2.4 Driver mutations and cancer precursors . . . . .	11
1.3 Human embryogenesis and its bottlenecks . . . . .	13
1.3.1 Zygote, cleavage, and blastulation . . . . .	13
1.3.2 Symmetry breaking, gastrulation and extraembryonic intercalation .	16
1.3.3 Aneuploidies in blastocysts and trophectoderm . . . . .	17
1.4 Questions and outline of this dissertation . . . . .	18
<b>2 Materials and methods</b>	<b>21</b>
2.1 Samples and sequencing . . . . .	21
2.1.1 Ethics, patient sampling and data availability . . . . .	21
2.1.2 Laser capture microdissection . . . . .	22
2.1.3 Whole-genome sequencing . . . . .	23
2.2 Somatic variant calling and filtering . . . . .	23
2.2.1 Filtering germline variants . . . . .	25

2.2.2	Filtering shared artefacts . . . . .	25
2.2.3	Distinguishing true presence from sequencing noise . . . . .	26
2.3	Binomial mixture models . . . . .	26
2.3.1	Truncated binomial distributions . . . . .	27
2.3.2	Estimating clonality . . . . .	28
2.3.3	Estimating tumour contamination in normal samples . . . . .	29
2.3.4	Timing of copy number gains . . . . .	29
2.4	Phylogenetic tree reconstruction . . . . .	30
2.5	Miscellaneous methods . . . . .	33
2.5.1	Methods pertaining to chapter 3 . . . . .	33
2.5.2	Methods pertaining to chapter 4 . . . . .	35
2.5.3	Methods pertaining to chapter 5 . . . . .	36
<b>3</b>	<b>Extensive phylogenies of normal human development</b>	<b>39</b>
3.1	Introduction . . . . .	39
3.2	Experimental design . . . . .	39
3.3	The clonality of LCM samples . . . . .	41
3.4	Reconstruction of phylogenies . . . . .	43
3.5	Embryonic asymmetries . . . . .	45
3.6	Targeted sequencing and organ-wide mosaic patterns . . . . .	51
3.7	Spatial genomics of mosaicism and embryonic patches . . . . .	54
3.8	Separation of primordial germ cells from somatic lineages . . . . .	59
3.9	Clonal expansions, benign prostatic hyperplasia and polyp formation . . . . .	61
3.10	Recurrent SNVs and the infinite sites model . . . . .	63
3.11	Other types of mutations and recurrent loss of the Y chromosome . . . . .	64
3.12	Conclusion . . . . .	67
<b>4</b>	<b>Embryonal precursors of Wilms tumour</b>	<b>69</b>
4.1	Introduction . . . . .	69
4.2	Detecting early clones in normal kidneys . . . . .	71
4.3	Nephrogenic clones are exclusive to individuals with Wilms tumour . . . . .	74
4.4	The driver of clonal nephrogenesis . . . . .	76
4.5	Timing the early expansion . . . . .	79
4.6	Sufficiency of <i>H19</i> hypermethylation . . . . .	82
4.7	Loss of imprinting versus loss of heterozygosity . . . . .	83
4.8	Conclusion . . . . .	84

---

<b>5</b>	<b>Universal mosaicism of the human placenta</b>	<b>87</b>
5.1	Introduction . . . . .	87
5.2	Placental biopsies contain clonal populations . . . . .	89
5.3	Trophoblast clusters are closely related clonal units . . . . .	92
5.4	Biases in cell allocation to trophoblast and inner cell mass . . . . .	93
5.5	A reversal of trisomy 10 . . . . .	95
5.6	Processes and impact of placental mutagenesis . . . . .	97
5.7	Conclusion . . . . .	99
<b>6</b>	<b>Conclusions and future perspectives</b>	<b>101</b>
6.1	Summary of the main findings . . . . .	101
6.2	Future perspectives and ongoing work . . . . .	102
6.2.1	Other childhood cancers, bilateral tumours and secondary malignancies	103
6.2.2	Further methodological advancements . . . . .	105
6.2.3	Closing remarks . . . . .	106
	<b>Bibliography</b>	<b>107</b>



# List of Figures

1.1	Drawing of a homunculus inside sperm . . . . .	2
1.2	Mutational signatures in normal tissues . . . . .	10
1.3	Overview of early embryogenesis . . . . .	15
3.1	The VAF distribution informs the clonality of LCM samples . . . . .	42
3.2	Phylogenies with mutation burden as branch length . . . . .	44
3.3	Most recent common ancestors of tissues. . . . .	46
3.4	Asymmetric contribution in PD28690 . . . . .	47
3.5	Asymmetric contribution in PD43850 and PD43851 . . . . .	50
3.6	Spatial resolution of embryonic patterns in brain . . . . .	52
3.7	Spatial resolution of embryonic patterns in mesoderm . . . . .	53
3.8	Spatial resolution of embryonic patterns in LCM sections of the small intestine	55
3.9	Embryonic patches in the colon . . . . .	57
3.10	Embryonic deconvolution of polyclonal skin samples . . . . .	58
3.11	Separation of primordial germ cells from somatic lineages . . . . .	60
3.12	Clonal expansions in PD28690 . . . . .	61
3.13	Mitochondrial SNVs and lineage tracing . . . . .	65
3.14	Recurrent loss of the Y chromosome . . . . .	66
4.1	Somatic mutations in PD37272 . . . . .	72
4.2	Distribution of clonal nephrogenesis . . . . .	75
4.3	<i>H19</i> methylation drives clonal nephrogenesis . . . . .	77
4.4	Phylogenies of bilateral and multifocal nephroblastoma . . . . .	80
4.5	PD41750, extensively sampled normal kidney . . . . .	81
4.6	<i>H19</i> loss of heterozygosity and loss of imprinting . . . . .	84
5.1	Overview of sampling of human placentas . . . . .	89
5.2	Burden and clonality of placental bulks . . . . .	90

---

5.3	Clonal architecture of microdissected trophoblast clusters and mesenchymal cores . . . . .	91
5.4	Diagram of embryonic mutations in placenta and umbilical cord . . . . .	93
5.5	Lineage specification in trophoblast phylogenies . . . . .	95
5.6	Trisomic rescue of chromosome 10 in PD45581 . . . . .	96
5.7	Mutational signatures and UPD in placenta . . . . .	98



# List of Tables

- 2.1 Mutation rate estimates in the first two cell generations . . . . . 34
- 2.2 Mutation rate estimates in subsequent cell generations . . . . . 34
  
- 5.1 Overview of different clinical groups in the placenta cohort. . . . . 88



# Nomenclature

## Greek Symbols

$\rho$  The overdispersion parameter in a beta-binomial distribution

## Acronyms / Abbreviations

AIC Akaike information criterion

AML Acute myeloid leukaemia

BIC Bayesian information criterion

ccRCC Clear cell renal cell carcinoma

CMN Congenital mesoblastic nephroma

cnn-LOH Copy number neutral loss of heterozygosity

CNV Copy number variant

DNV Double nucleotide variant

Indel Insertion and deletion

LCM Laser capture microscopy/microdissection

LOH Loss of heterozygosity

LOI Loss of imprinting

MDS Myelodysplastic syndrome

MNV Multinucleotide variant

MRT Malignant rhabdoid tumour

PGC	Primordial germ cell
RNA	Ribonucleic acid
SBS	Single base substitution
SNP	Single nucleotide polymorphism
SNV	Single nucleotide variant
SV	Structural variant
TGS	Targeted genome sequencing
VAF	Variant allele frequency
WES	Whole-exome sequencing
WGS	Whole-genome sequencing