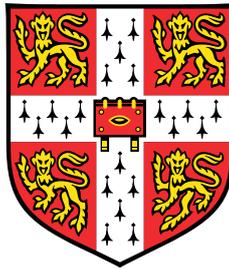


# The Intra-Tumour Heterogeneity Landscape of Human Cancers



**Stefan Christiaan Dentro**

Christ's College  
University of Cambridge

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

December 2017



## **Declaration**

This thesis 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, 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 thesis 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. It does not exceed the prescribed word limit for the relevant Degree Committee

Stefan Christiaan Dentre  
December 2017



# **Abstract**

## **The Intra-Tumour Heterogeneity Landscape of Human Cancers**

**Stefan Christiaan Dentre**

Tumours accumulate many somatic mutations in their lifetime. Some of these mutations, drivers, convey a selective advantage and can induce clonal expansions. Incomplete clonal expansions give rise to intra-tumour heterogeneity. Somatic mutations can be measured through massively parallel sequencing, where mutations that are supporting incomplete expansions will appear as subclonal. These mutations can be used as a marker of the existence of the expansion and allow for a window into the clonal and subclonal architecture of the tumour at diagnosis.

During my Ph.D. I have developed computational methods to infer intra-tumour heterogeneity from massively parallel sequencing data and applied these to the 2,778 tumour whole genome sequences in the International Cancer Genome Consortium Pan-Cancer Analysis of Whole Genomes initiative to paint the pan-cancer landscape of intra-tumour heterogeneity.

I will first introduce the methods; a method to call somatic copy number alterations (Battenberg) and a method to infer subclones from single nucleotide variants (DPClust). Both are extensively validated on simulated and on real data, and I describe a rigorous quality control procedure. The methods are then applied to a single sample to showcase what can be learned about the life history of a cancer, before introducing additional computational methods for a pan-cancer study of heterogeneity. Finally, I describe the findings.

I find that nearly all cancers, for which there is sufficient power, contain at least one subclone (96.7% of 1,801 primary tumours). The subclones contain driver mutations that are under positive selection, and known cancer genes contain subclonal driver mutations in low proportions. 9.5% of tumours contain only subclonal drivers that are clinically actionable, suggesting that heterogeneity could inform treatment choices. Finally, the analysis reveals that activity of smoking and UV-light associated mutational signatures goes down as the tumour evolves, while activity of the APOBEC associated signatures goes up.



## Acknowledgements

I would like to thank my supervisors, Peter Van Loo, David Wedge and David Adams. I'm very grateful to have had the opportunity to work with you and on this very exciting project! Thank you very much for your guidance, thoughts, ideas and your patience. It has been an enormous pleasure to learn from you and to work with you! Thank you for the great discussions and for continuously pushing me to be a better scientist in a friendly and positive environment.

During the PCAWG project I've been in the fortunate position to closely collaborate with members of the Van Loo lab. Thank you Maxime Tarabichi, Kerstin Haase, Clem Jolly, Jonas Demeulemeester and Matt Fittall for the fantastic collaboration and the many in-depth discussions on tumour evolution and heterogeneity. I'd like to also thank the Van Loo lab for being my 'home' for the last few years. Annelien, Clem, Jonas, Kerstin, Lilly, Matt, Maxime and Peter: it's been an enormous pleasure!

Most of the work in this thesis was done as part of the international PCAWG collaboration, I'd like to thank everyone of the PCAWG collaborators and in particular Quaid Morris, Moritz Gerstung, Jeff Wintersinger, Ignaty Leshchiner, Amit Deshwar, Yulia Rubanova and Peter Campbell.

I would like to extend my thanks to Jason J. Pitt for the fruitful collaboration on the West-African breast cancer study and to members of the Adams lab, the Wedge lab and of the Cancer Genome Project for stimulating discussions and continuous helpful feedback.

This thesis could not have happened without the love and support from my parents and sister. Thank you for supporting me and believing in me!

Finally, I would like to thank the Wellcome Trust for generously funding this Ph.D.



## Preface

During my Ph.D. I have been in the very fortunate position to heavily collaborate with colleagues close by and far away. Nearly all of the work reported in this thesis was performed as part of an international collaboration project to jointly analyse the cancer whole genome sequencing samples that are part of the International Cancer Genome Consortium (ICGC) The Pan-Cancer Analysis of Whole Genomes (PCAWG) initiative. As this data set is referenced throughout this thesis I have included below a high level description of the project and data set.

Nearly all of the work reported in this thesis is performed in collaboration with, or building on top of the work done by others. Throughout this thesis I therefore systematically refer to work that is solely to my credit with "I", and work that was done by others (with my involvement) as "we". There is also the occasional reference to "a collaborator", where work was done that I had no involvement with.

To help the story line I have spread the introduction across Chapters 1 and 2. Chapter 1 contains a brief review of the relevant literature, but descriptions of algorithms which were already in advanced development when I started are described in Chapter 2. I have worked on and with these algorithms throughout my Ph.D. and have made numerous improvements and adjustments (of which improvements on computational resource requirements are not reported as they are not of direct scientific interest). I felt it would make this thesis more easily readable when the algorithms and updates are described in one chapter.

This setup was chosen to paint a comprehensive story that can hopefully be understood from this thesis alone.

### **The ICGC Pan-Cancer Analysis of Whole Genomes initiative**

The International Cancer Genome Consortium (ICGC) was created to coordinate cancer genome sequencing projects spanning 50 different types of cancer, with the aim to sequence over 25,000 cancer genomes (ICGC Consortium, 2010). ICGC is organised as a series of projects based in countries spread all over the world that focus on analysis of a single cancer

type. Over 17,000 cancers have now been sequenced, of which the majority through whole exome sequencing.

The Pan-Cancer Analysis of Whole Genomes (PCAWG) project was launched to comprehensively characterise those samples for which whole genome sequencing (WGS) is available as a single data set (Campbell et al., 2017). The advantage of focussing on samples for which WGS is available is that the full genome can be interrogated, including the full array of single nucleotide variants (SNV), indels (short insertions and deletions) and structural variants (SVs).

The project consists of 16 working groups with each their own distinct theme. I have been a member of the working group that focusses on tumour evolution and heterogeneity, which is a collaboration of about 60 scientists representing 12 different laboratories.

The tumours that are part of ICGC PCAWG had to meet a series of criteria to be included: a minimal set of clinical annotations should be available, both tumour and normal samples have to be paired-end sequenced from an Illumina machine to a coverage of at least 30x and 25x respectively (Campbell et al., 2017). The data set consists of primarily treatment-naive primary tumours and nearly all matched normals are generated from blood samples.

Data from 2,834 donors was selected to be included in ICGC PCAWG, of which data from 2,658 donors passed quality control procedures (Whalley et al., 2017). The analysed data set consists of 2,778 tumours, of which 2,605 are primary tumours and 173 from a metastasis or relapse case, and spreads 39 histologically distinct types of cancer. Each sequencing sample was processed using a standardised set of primary analysis pipelines, that includes alignment of the sequencing reads and variant calling and filtering from pipelines provided by the Sanger, Broad and EMBL/DKFZ (Yung et al., 2017). These pipelines were extended by one additional SNV and one additional indel caller to further increase the reliable detection of low allele frequency variants (Campbell et al., 2017). Clinical data was systematically collected and standardised (Campbell et al., 2017).

Output from the three primary variant calling pipelines was combined into a high quality set of somatic consensus SNVs (Campbell et al., 2017), indels (Campbell et al., 2017), SVs (Campbell et al., 2017) and copy number alterations (CNAs) (Dentro et al., 2017, manuscript in preparation), of which the latter is described in this thesis. The optimal strategy to find consensus SNVs and indels was found by first running 19 different variant callers across a selected set of 64 tumours. 250,000 calls were selected for validation through deep targeted capture sequencing such that every combination of variant callers is represented and were stratified by allele frequency, after which consensus strategy that maximises precision and recall was then developed to generate the final PCAWG calls (Campbell et al., 2017).

These steps have created the largest set of whole genome cancer sequences to date, spreading a broad range of cancer types. The data set is uniformly processed and the variant calling pipelines have been extensively validated. It therefore provides a unique opportunity for a high quality, in-depth study of tumour heterogeneity.



# Table of contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Cancer . . . . .	1
1.2	Hallmarks . . . . .	1
1.3	Genetics . . . . .	2
1.3.1	Early observations suggest a role for the genome . . . . .	2
1.3.2	Confirmation of a role in cancer . . . . .	3
1.4	Evolution . . . . .	3
1.4.1	Estimates of the number of mutations to form a cancer . . . . .	3
1.4.2	An integrated theory of tumour evolution . . . . .	4
1.5	Drivers . . . . .	4
1.5.1	Oncogenes . . . . .	5
1.5.2	Tumour suppressor genes . . . . .	5
1.5.3	Drivers and passenger mutations . . . . .	5
1.6	High throughput technology . . . . .	5
1.6.1	Array based technology . . . . .	6
1.6.2	Sequencing technology . . . . .	6
1.6.3	The emergence of sequencing consortia . . . . .	6
1.7	Copy number . . . . .	7
1.7.1	Confirming classic knowledge . . . . .	7
1.7.2	Pan-cancer overview of CNAs . . . . .	7
1.8	Massively parallel sequencing of cancer genomes . . . . .	8
1.8.1	Early findings from exome sequencing . . . . .	8
1.8.2	The full compendium of somatic alterations . . . . .	9
1.8.3	Whole genome sequencing reveals the extent of somatic alterations in cancer genomes . . . . .	9
1.9	Mutational processes . . . . .	10
1.9.1	Automated extraction of mutational signatures . . . . .	10

1.9.2	Linking signatures to mutational processes . . . . .	10
1.10	Heterogeneity . . . . .	11
1.10.1	Detecting heterogeneity from sequencing data . . . . .	11
1.10.2	Cancer type specific studies highlight evolutionary properties . . . . .	11
1.10.3	Pan-cancer studies reveal widespread ITH across cancer types . . . . .	12
1.11	Subclonal inference . . . . .	13
1.11.1	Clustering of mutations . . . . .	13
1.11.2	Statistical and computational strategies for subclonal reconstruction . . . . .	14
1.12	Copy number calling . . . . .	14
1.13	Tumour micro-environment . . . . .	16
1.13.1	Carcinoma-associated fibroblasts . . . . .	16
1.13.2	Tumour-associated macrophages . . . . .	16
1.13.3	Tumour-infiltrating lymphocytes . . . . .	17
1.13.4	Tumour-associated neutrophils . . . . .	17
1.13.5	Other cell types . . . . .	17
1.13.6	Immune evasion . . . . .	17
1.14	Clinical implications of heterogeneity . . . . .	18
1.15	Summary . . . . .	19
<b>2</b>	<b>Methods</b>	<b>21</b>
2.1	Principles of subclonal reconstruction . . . . .	21
2.2	The Battenberg algorithm . . . . .	21
2.2.1	Pre-processing . . . . .	22
2.2.2	Reconstructing haplotype blocks . . . . .	22
2.2.3	Fitting a global copy number profile . . . . .	25
2.2.4	Testing whether a segment is clonal . . . . .	25
2.2.5	Fitting subclonal copy number . . . . .	26
2.2.6	Extensions to segmentation . . . . .	27
2.2.7	GC content correction . . . . .	28
2.3	Subclonal architecture inference with DPCLust . . . . .	30
2.3.1	Estimating cancer cell fractions . . . . .	30
2.3.2	Filtering . . . . .	34
2.3.3	Algorithm . . . . .	35
2.3.4	Post-processing . . . . .	38
2.3.5	Extension to multi-sample cases . . . . .	39
2.3.6	Co-clustering of indels and CNAs . . . . .	40
2.3.7	Alternative post-processing steps . . . . .	43

---

2.3.8	A downsampling strategy . . . . .	44
2.4	Automated post-hoc tree building . . . . .	45
2.4.1	Cluster-pair classification . . . . .	45
2.4.2	Tree building . . . . .	47
<b>3</b>	<b>Validation of methods</b>	<b>49</b>
3.1	Introduction . . . . .	49
3.2	Simulating subclonality with SimClone . . . . .	49
3.2.1	Introduction . . . . .	49
3.2.2	Assumptions . . . . .	50
3.2.3	Simulating a tree . . . . .	50
3.2.4	Determining cluster sizes . . . . .	51
3.2.5	Simulating mutations . . . . .	51
3.2.6	Extension to simulating multi-sample cases . . . . .	53
3.2.7	Simulating copy number . . . . .	53
3.3	SimClone1000, a validation data set for PCAWG . . . . .	59
3.4	Metrics to evaluate a subclonal reconstruction . . . . .	61
3.5	A lower bound generated by RandomClone . . . . .	61
3.5.1	Introduction . . . . .	61
3.5.2	RC – Stick breaking . . . . .	62
3.5.3	RC – Informed . . . . .	62
3.5.4	RC – Uniform . . . . .	62
3.5.5	RC – Single cluster . . . . .	63
3.6	Validation of multiplicity calls . . . . .	63
3.7	Assesment of a subclonal architecture through resimulations . . . . .	63
3.8	Validation of DPclust . . . . .	66
3.9	Validation of assignments of gained mutations . . . . .	69
3.10	Validation of Battenberg . . . . .	71
<b>4</b>	<b>A rigorous quality control procedure</b>	<b>77</b>
4.1	Introduction . . . . .	77
4.2	Quality control metrics . . . . .	78
4.2.1	Large homozygous deletions . . . . .	79
4.2.2	50% subclone in copy number . . . . .	83
4.2.3	Empty odd numbered copy number state . . . . .	83
4.2.4	No clonal copy number alteration . . . . .	87
4.2.5	Shifted clonal mutation cluster . . . . .	89

4.2.6	Mutation cluster at 50% of tumour cells . . . . .	89
4.2.7	Empty mutation copy number state . . . . .	92
4.3	Resolving a quality control failure . . . . .	95
4.3.1	Automatic correction . . . . .	95
4.3.2	Manual correction . . . . .	95
4.4	Inventory of metric triggers in the PCAWG data set . . . . .	96
<b>5</b>	<b>The subclonal architecture and life history of a single cancer</b>	<b>101</b>
5.1	Introduction . . . . .	101
5.2	Background . . . . .	102
5.3	Methods . . . . .	102
5.4	Subclonal architecture . . . . .	104
5.5	Mutational Signatures . . . . .	106
5.6	Kataegis . . . . .	107
5.7	The life history of N010985 . . . . .	109
<b>6</b>	<b>Methods for a pan-cancer study of tumour heterogeneity</b>	<b>113</b>
6.1	Introduction . . . . .	113
6.2	Consensus copy number . . . . .	114
6.2.1	Assumptions behind different copy number callers . . . . .	114
6.2.2	Determining consensus segment breakpoints . . . . .	117
6.2.3	Constructing consensus copy number . . . . .	120
6.2.4	Rounding subclonal copy number . . . . .	121
6.2.5	Chromosomes X and Y . . . . .	122
6.2.6	Panel of experts review . . . . .	123
6.2.7	Determining consensus purity . . . . .	124
6.2.8	Filtering . . . . .	126
6.3	Consensus subclonal architecture . . . . .	127
6.3.1	Three consensus approaches . . . . .	128
6.3.2	Performance comparison . . . . .	130
6.4	Purity, ploidy and sequencing coverage determine ability to detect subclones	135
6.5	Correcting for the winner's curse . . . . .	136
<b>7</b>	<b>A pan-cancer overview of tumour heterogeneity</b>	<b>139</b>
7.1	Introduction . . . . .	139
7.2	Methods . . . . .	140
7.3	Nearly all primary tumours contain detectable subclones . . . . .	141

---

7.4	Metastatic melanomas are often clonal . . . . .	142
7.5	Subclonal driver mutations in known cancer genes . . . . .	143
7.6	14% of mutations are undetected . . . . .	147
7.7	Clear signs of positive selection in subclonal mutations . . . . .	148
7.8	Subclonal clinically actionable events . . . . .	150
7.9	Evidence of additional heterogeneity . . . . .	152
7.10	Cancer types follow individual evolutionary narratives . . . . .	154
<b>8</b>	<b>Discussion</b>	<b>159</b>
8.1	Overall summary . . . . .	159
8.2	Future directions . . . . .	161
8.2.1	A more in-depth view of intra-tumour heterogeneity . . . . .	161
8.2.2	Tumour evolution . . . . .	162
8.2.3	Towards clinical application . . . . .	162
8.2.4	Methods . . . . .	163
	<b>References</b>	<b>167</b>
	<b>List of figures</b>	<b>193</b>
	<b>List of tables</b>	<b>195</b>
	<b>Glossary</b>	<b>197</b>
	<b>Appendix A The evolutionary history of 2,658 cancers</b>	<b>201</b>
	<b>Appendix B The evolutionary history of breast adenocarcinoma</b>	<b>235</b>
	<b>Appendix C Published works</b>	<b>237</b>

