

Development of new methods for accurate estimation of tumour heterogeneity

by

Sebastian Hollizeck

[ORCID: 0000-0002-9504-3497](#)

A thesis submitted in total fulfillment for the
degree of Doctor of Philosophy

in the
Sir Peter MacCallum department of Oncology
Melbourne School of Awesome
THE UNIVERSITY OF MELBOURNE

November 24, 2021

THE UNIVERSITY OF MELBOURNE

Abstract

Sir Peter MacCallum department of Oncology
Melbourne School of Awesome

Doctor of Philosophy

by [Sebastian Hollizeck](#)

ORCID: [0000-0002-9504-3497](#)

Intra-patient tumour heterogeneity is a widely accepted cause of resistance to therapy [1, 2], but the possibility to study this phenomenon is so far underexplored as the acquisition of multi region data sets is costly and ethically challenging [3]. With circulating tumour DNA (ctDNA) as a proxy it is possible to analyze a snapshot of the unified heterogeneity, but there is still an unmet need for new analysis methods to optimize the analysis of these very valuable data and drive new treatment targets [4]. In this work we will develop new methods to study genetic heterogeneity from next generation sequencing (NGS) of tumour tissue as well as ctDNA to elucidate the role of tumour heterogeneity on treatment resistance.

Declaration of Authorship

I, AUTHOR NAME, declare that this thesis titled, 'THESIS TITLE' and the work presented in it are my own. I confirm that:

- The thesis comprises only my original work towards the NAME OF AWARD except where indicated in the preface;
- due acknowledgement has been made in the text to all other material used; and
- the thesis is fewer than the maximum word limit in length, exclusive of tables, maps, bibliographies and appendices as approved by the Research Higher Degrees Committee.

Signed:

Date:

Preface

This preface includes a summary of all chapters in this work as well as a comprehensive summary of my contributions and everyone else's contribution. This is a thesis *with* publications and each publication included in a chapter is shown here.

Hollizeck S., Wong S.Q., Solomon B., Chandrananda D.¹, Dawson S-J.¹ "Custom workflows to improve joint variant calling from multiple related tumour samples: FreeBayesSomatic and Strelka2Pass" *Bioinformatics*. 2021. DOI: [10.1093/bioinformatics/btab606](https://doi.org/10.1093/bioinformatics/btab606)

Chapter 1: Introduction is an original work providing background and overview relevant to understanding the thesis and its relevance to the field. It includes an introduction to DNA, ctDNA, DNA sequencing, somatic variant calling and lung cancer.

¹These authors contributed equally and are considered shared last.

Chapter 2: Joint somatic variant calling is an original work describing two workflows for the joint analysis of multiple related tumour samples and has been published in *Bioinformatics* as "Custom workflows to improve joint variant calling from multiple related tumour samples: FreeBayesSomatic and Strelka2Pass" on 21st September 2021. In addition to the published analysis, I have added longitudinal analysis and its evaluation.

Contributions for this chapter:

- I conceptualised the work
- I implemented the workflows and containerised all required tools
- I performed the data simulation
- I performed the analysis presented in the publication
- I wrote the draft of the manuscript and performed revisions
- D.C. and S-J.D. provided advice in planning and writing the manuscript
- D.C. provided guidance for method development
- S-J.D. provided guidance for method evaluation
- S.W. performed the targeted amplicon validation
- S.W. and B.S. read the draft manuscript and provided feedback
- B.S. provided clinical expertise for human data

Chapter 3:

summary plus contributions

Chapter 4:

Chapter 5:

Other publications These publications i have contributed to in my candidature, but they are not presented in this work

Burr M.L., Sparbier C.E., Chan K.L., Chan Y-C., Kersbergen A., Lam E.Y.N., Azidis-Yates E., Vassiliadis D., Bell C.C., Gilan O., Jackson S., Tan L., Wong S.Q., **Hollizeck S.**, Michalak E.M., Siddle H.V., McCabe M.T., Prinjha R.K., Guerra G.R., Solomon B.J., Sandhu S., Dawson S-J., Beavis P.A., Tothill R.W.,

Cullinane C., Lehner P.J., Sutherland K.D., Dawson M.A. "An evolutionarily conserved function of polycomb silences the MHC class I antigen presentation pathway and enables immune evasion in cancer" *Cancer cell.* 2019. DOI: [10.1016/j.ccr.2019.08.008](https://doi.org/10.1016/j.ccr.2019.08.008)

Solomon B.^{J²}, Tan L.², Lin J.J.², Wong S.Q.², Hollizeck S.², Ebata K., Tuch B.B., Yoda S., Gainor J.F., Lecia V. Sequist L.V., Oxnard G.R., Gautschi O., Drilon A., Subbiah V., Khoo C., Zhu E.Y., Nguyen M., Henry D., Condroski K.R., Kolakowski G.R., Gomez E., Ballard J., Metcalf A.T., Blake J.F., Dawson S-J., Blosser W., Stancato L.F., Brandhuber B.J., Andrews S., Robinson B.G., Rothenberg S.M "RET Solvent Front Mutations Mediate Acquired Resistance to Selective RET Inhibition in RET-Driven Malignancies" *Journal of Thoracic Oncology.* 2020. DOI: [10.1016/j.jtho.2020.01.006](https://doi.org/10.1016/j.jtho.2020.01.006)

Add Katie and Danes paper

Funding:

All necessary funding goes here

Instructions: Where applicable, the following information must be included in a preface:

- a description of work towards the thesis that was carried out in collaboration with others, indicating the nature and proportion of the contribution of others and in general terms the portions of the work which the student claims as original;
- a description of work towards the thesis that has been submitted for other qualifications;
- a description of work towards the thesis that was carried out prior to enrolment in the degree;
- whether any third party editorial assistance was provided in preparation of the thesis and whether the persons providing this assistance are knowledgeable in the academic discipline of the thesis;
- the contributions of all persons involved in any multi-authored publications or articles in preparation included in the thesis;
- the publication status of all chapters presented in article format using the descriptors below;
 - Unpublished material not submitted for publication
 - Submitted for publication to [publication name] on [date]
 - In revision following peer review by [publication name]

²These authors contributed equally and are considered shared first.

- Accepted for publication by [publication name] on [date]
 - Published by [publication name] on [date]
- an acknowledgement of all sources of funding, including grant identification numbers where applicable and Australian Government Research Training Program Scholarships, including fee offset scholarships.

Acknowledgements

The acknowledgements and the people to thank go here, don't forget to include your project advisor...

Lots of figures in the introductory chapter 1 were created with the help of [BioRender.com](#)

think of where to put the package citations; Probably at the end as appendix

0.1 Software and packages

This section is dedicated to all the software that usually gets uncited because they are "standard" or backbone

Most analysis in a prototype state was done on a linux cluster running Centos 7 [5] with Bash [6] and due to the high amount of data, parallel [7] was used of the multi-cpu architecture of HPCs.

0.1.1 R

In depth data analysis and visualisation was done with R [8] with the help of packages listed below.

Most of the parallelisation in R was performed with BiocParallel [9], which is available through BiocManager [10].

Colour schemes and manipulation was performed with colorspace [11, 12].

Copynumber analysis was performed with sequenza [13], FACETS [14, 15] and PURPLE [16]. Some analysis was also directly performed with copynumber [17, 18].

Variant effect prediction was performed with VEP [19].

Table manipulation was performed with data.table [20].

Violin plots were generated with vioplot [21].

Heatmaps and UpSet plots were generated with ComplexHeatmap [22]

Phylogenetic analysis was performed with both ape [23] and phangorn [24] followed by dendextend [25].

Google sheets and its built in scripts were used to collect stats on docker pull requests and the data was then read in R through googlesheets4 [26].

Additional libraries, which were used for a multitude of things are listed in no particular order below: Rsamtools [27], GenomicRanges [28], optparse [29], VariantAnnotation [30], MultiAssayExperiment [31], circlize [32], BioQC [33], Biostrings [34], deconstructSigs [35], BSgenome [36], QDNAseq [37], RColorBrewer [38], pheatmap [39], ensemblVEP [40], stringdist [41], Rsubread [42], svglite [43], grImport [44], XML [45], kableExtra [46], lsa [47], irlba [48], ggplot2 [49]

maybe itemize over just a blob

o.1.2 python

Analysis for [chapter 4](#) was mostly done through python [50] with the help of many different packages, which are listed here in no particular order: numpy [51], ncls [52], pysam [53, 54, 55], zarr [56], pandas [57, 58], quadprog [59] as well as scipy [60].

Contents

Abstract	iii
Declaration of Authorship	iv
Preface	v
Acknowledgements	ix
o.1 Software and packages	ix
o.1.1 R	ix
o.1.2 python	x
List of Figures	xiii
List of Tables	xv
Abbreviations	xvii
Constants	xx
Symbols	xxiii
1 Introduction	1
1.1 DNA	1
1.1.1 Ploidy	3
1.1.2 Mutations	4
1.2 cfDNA	5
1.3 DNA sequencing	6
1.3.1 Library preparation	7
1.3.2 Next generation sequencing	8
1.3.3 Long read sequencing	8
1.4 DNA analysis	8
1.4.1 Mapping	10
1.4.2 Variant calling	10
1.4.3 Germline	11
1.4.4 Somatic	11
1.5 Cancer	12
1.5.1 Lungcancer	14
1.6 Overview	15

2 Joint somatic variant calling	17
2.1 Introduction	17
2.2 Publication	18
2.2.1 Summary	19
2.2.2 FreeBayesSomatic workflow	19
2.2.3 Strelka2Pass workflow	20
2.2.4 Validation	21
2.3 Effects on downstream analysis	25
2.3.1 Phylogenetic reconstruction	25
2.4 Longitudinal analysis	28
2.4.1 Clonal deconvolution	30
2.4.2 Longitudinal enriched phylogeny	32
2.5 Usage	33
3 CASCADE	35
3.1 Introduction	35
3.2 Publication	35
3.3 Cohort analysis	35
3.4 Mitochondrial phylogenetic reconstruction	35
3.5 Outlook	35
4 Mismatchfinder	37
4.1 Introduction	37
5 Conclusion	39
A Strelka2Pass and FreeBayesSomatic publication	41
A.1 Introduction	42
A.2 Materials and methods	43
A.2.1 FreeBayesSomatic workflow	43
A.2.2 Strelka2Pass workflow	44
A.3 Validation	45
A.3.1 Simulated data	46
A.3.2 Clinical data	47
A.4 Discussion	48
A.5 Supplementary methods	61
A.5.1 Alignment of clinical data	61
A.5.2 Validation of clinical data	61
A.5.3 Purity estimation with sequenza	62
A.5.4 Performance of individual steps in Strelka2Pass	62
A.5.5 Ensemble workflows – user suggestions	62
Bibliography	63

List of Figures

1.1	Overview DNA structure	2
1.2	Overview Chromosome structure	3
1.3	Overview DNA structure	5
1.4	Library preparation for NGS	7
1.5	Sequencing by synthesis (Illumina)	9
1.6	Original hallmarks of cancer	13
1.7	New hallmarks of cancer	14
2.1	Comparison of joint multi-sample variant calling and single tumour-normal paired calling methods	22
2.2	Reconstructed phylogenies of joint samples	26
2.3	Tanglegram of the reconstructed phylogenies	27
2.4	Improved somatic variant calling in longitudinal data	29
2.5	Longitudinal data informs diagnostic variant calling	30
2.6	Reconstructed clonal trees for joint and pairwise variant calling	31
2.7	Reconstructed phylogeny with longitudinal ctDNA samples	32
2.8	Usage statistics joint workflows	33
A.1	Comparison of joint multi-sample variant calling and single tumour-normal paired calling methods	45
A.2	Characteristics of simulated data	51
A.3	Performance of workflows using simulated data	51
A.4	Variant allele frequencies (VAF) of variants detected by joint sample analysis	52
A.5	Performance of individual steps in the Strelka2pass workflow using the simulated data	53
A.6	Summary of variant filters assigned by Mutect2	54
A.7	Assessing the performance of different workflows using tumour samples with different evolutionary relationships in the simulated data	55
A.8	Correlation of variant allele frequencies in validation	56
A.9	Performance of the different workflows using clinical samples from eight cancer patients	57
A.10	Correlation between cellularity and proportion of variants found only with joint calling using FreeBayesSomatic	58
A.11	Improvement in recall using FreeBayesSomatic and Strelka2pass over Mutect2 in the clinical samples	58
A.12	Performance of ensemble variant calling strategies	59

List of Tables

A.1	Sample name mapping	60
A.2	Runtime of different workflows on simulated data	61

Abbreviations

DNA	D eoxyribo N ucleic A cid
RNA	R ibo N ucleic A cid
cfDNA	c ell f ree D NA
ctDNA	c irculating t umour D NA
bp	b ase p air
ChIP	C hromatin I mmuno P recipitation
WGS	W hole G enome S equencing
WES	W hole E xome S equencing
SCLC	S mall C ell L ung C ancer
NSCLC	N on- S mall C ell L ung C ancer
RAID	R edundant A rray of I ndependent D isks
SNP	S ingle N ucleotide P olymorphism
InDel	I nsertion or D eletion
SV	S tructural V ariant
PON	P anel O f N ormals
GATK	G enome A nalysis T ool K it
NJ	N eighbour J oining
UPGMA	U nweighted P air G roup M ethod with A rithmetic mean
WPGMA	W eighted P air G roup M ethod with A rithmetic mean
F81	F elsenstein 19 81 model
HKY85	H asegawa, K ishino and Y ano 19 85 model
HPC	H igh P erformance C omputing
MRCA	M ost R ecent C ommon A ncestor

Sort alphabetically

Constants

Speed of Light c = $2.997\,924\,58 \times 10^8$ ms⁻¹ (exact)

Symbols

a	distance	m
P	power	W (Js ⁻¹)
ω	angular frequency	rads ⁻¹

“Begin at the beginning,” the King said, very gravely, “and go on till you come to the end: then stop.”

— Lewis Carroll, *Alice in Wonderland*



Introduction

This first introduction chapter contains all the necessary background information as well as an overview for the work discussed in this thesis. It summarised basic biological properties of DNA and cell biology as well as the respective technologies to read, analyse and measure these biological concepts and then how to evaluate the output of these methods. [Section 1.1](#) delineates the role DNA plays for the cell and then [section 1.2](#) shows how these standards are changed in the tumour and cell free context. [Section 1.3](#) introduces the current technologies used to measure and detect DNA and its variations. With [section 1.4](#) covering the computational analysis methods to read out changes in the DNA. Then [section 1.5.1](#) relates how these changes lead to cancer and what we can learn from them. The introduction concludes with [section 1.6](#) as an overview over the thesis aims and my work in addressing them in the following chapters.

1.1 DNA as a information storage unit

It is a widely accepted fact, that Deoxyribonucleic acid (DNA) serves as the long term information storage molecule of our cells. This information is protected and allows correction of simple errors through its double helix structure [61, 62]. The nucleotides, which consist of a deoxyribose sugar (hence the name), a phosphate group and the nitrogenous base, are joined together by phosphate groups. Even though there are six common naturally occurring nitrogenous bases: Adenine (A), Thymine (T), Guanine (G), Cytosine (C), Uracil (U) and nicotinamide, only the first four are used to encode the genetic information into DNA. Each of the strands mirrors the other, so that an adenine will be paired up with a thymine forming two hydrogen bonds. Similarly cytosine will pair with guanine forming an even stronger bond with three hydrogen bonds. While other pairings which do

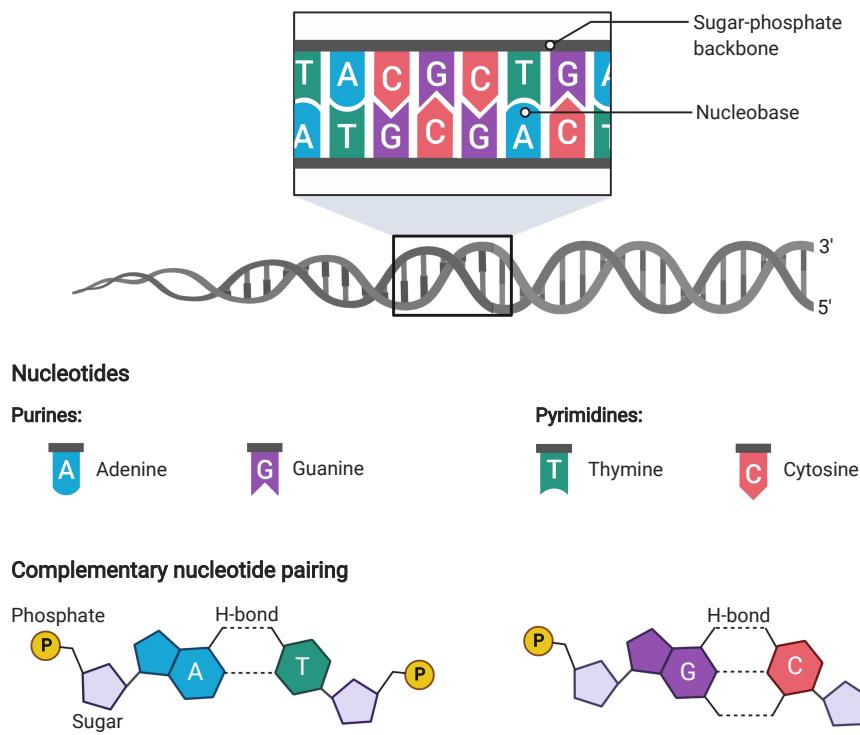


Figure 1.1: Overview of DNA structure and the nucleobases, which form DNA strands. Nucleotides are split into Purines and Pyrimidines by the structure of the nitrogen ring; complementary pairing of bases is shown as shapes of the bases as well as with 2D structures; Hyrdogen (H) bonds are shown as dotted lines; Phosphates are shown as P; 3' and 5' ends are defined by the internal number of the carbon atom of the sugar which is exposed; Adapted from “DNA structure” by BioRender.com (2021) Retrieved from <https://app.biorender.com/biorender-templates>

not follow those rules are chemically possible, they are mostly observed in ribonucleic acid (RNA) [63]. These very strict bonding rules enable the DNA to be similar to a hard drive with backup on a computer. And as only one strand contains all the information, the DNA polymerase enzyme does only need access to one strand, which allows parallel replication during cell division, but also error corrections, by proof reading the newly synthesised strand with the template. In order to be able to distinguish the two strands, they were assigned the names 3' and 5' depending on the numbering of the carbon atom in the sugar, which is exposed (Figure 1.1).

The entirety of the DNA encoding the organism is commonly called “the genome” with all genes, which consist of introns and exons are called exome. Unicellular organisms usually only have a very small amount of introns, which to current knowledge only provide limited information and are only responsible for structure. In vertebrates introns as well as intergenic DNA (the DNA between genes) contribute most of the DNA in the genome. For example in humans, only 1% of the genetic

material is considered to be exonic, whereas introns contribute $\approx 24\%$ and the rest is intergenic ($\approx 75\%$)[64].

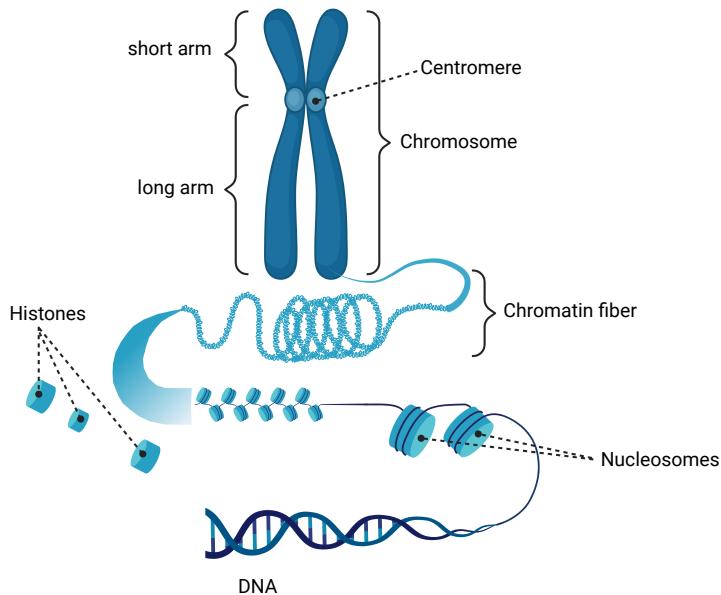


Figure 1.2: Structural overview of the metaphase condensed chromosome: DNA is first wrapped around Histones to form nucleosome, which then associate with each other to form the chromatin fiber, which in the metaphase of the cell cycle is condensed even more into the X-shaped chromosome

The DNA in eukaryotes however is not free floating around in the nucleus of a cell, but rather in most eukaryotic organisms, it is highly condensed and structured, first wrapped around nucleosomes like thread on a spool, then organised around histones, into either open (accessible) or closed chromatin, which then can be even further condensed into chromosomes, which have a X-like shape, with two shorter and two longer arms (Figure 1.2). This allows some of the DNA to be accessible where the use of other areas can be restricted[65]. Through this restriction, the availability of certain genes, which are the sections of the DNA, which encode for short term storage molecules like RNA. This restriction plays an important role in cell fate and cell viability. Ultimately all information stored to create a new highly complex organism is stored in just the DNA of one cell. Whichever parts are used and how they are used decides the function and the identity of the cell[66].

1.1.1 Ploidy - its good to have a backup, if you do it right

Similar to the already discussed RAID-like setup of the DNA in two strands, another concept of data security, a spatial different storage is also implemented. Most eukaryotic organisms have at least

two of each chromosome (diploid) with some species reaching up to septaploid [67]. However, this concept is not the only reason for the ploidy of somatic cells. For sexually reproducing organisms, at least a diploid set of chromosomes is necessary to enable information to be joined from both parents. Germline cells (sperm and egg) are generally monoploid, such that the resulting cell will be diploid, but the ploidy of the somatic cells is not as uniform within a species, where it can vary between organisms based on gender or rank [68]. In most organisms, a change in ploidy is fatal [69] and only partial ploidy changes like extra copies of chromosome 17 [70], chromosome 18 [71] and chromosome 21 [72] are tolerated. These syndromes can occur when there is an uneven split of chromosomes during cell division. The additional advantage, apart from sexual reproduction, is that a second almost identical copy of a chromosome allows repair of DNA, even when both strands are damaged, for example in a double strand break. In this case, the information from the sister chromosome will be used, by first cutting the double strand break ends to have overhang (resection). This overhang will then merge with the sister chromosome's mirrored strand. In this state, the two chromosomes are fused together in a Holliday junction, which allows the missing part from the resection and the double strand break to be synthesised [73]. During this process, which is part of the homology directed repair (HDR) machinery, the sister chromosomes exchange parts of their DNA, when resolving the Holliday junction. As these stretches of DNA do not need to be 100% identical, this plays an important role in evolution and diversity [74, 75].

Even though this X-like structure is the most commonly used and known structure, the DNAs 3D structure is usually very different and only takes this shape for the very short time of the cell cycle. Most of the time, the chromosomes are unravelled into something resembling a ball of yarn, where the “open” chromatin regions are on the outside and the “closed” regions are “hidden” in the inside and each chromosome establishes its own “territory” inside the nucleus (Figure 1.3). This structure allows another DNA cross over with non-sister chromosomes, which is called a chiasma.

1.1.2 Phantastical mutations and where to find them

However even though the DNA is highly stable and error correction methods are constantly working to not introduce any changes in the DNA, the source of evolution and adaptation of species is sourced in a steady mutation rate [76, 77]. These changes in normal tissue are mostly irrelevant to the organism as a whole and will not be passed on to the next generation. These changes are known as somatic mutations. This type of mutation accumulates in a cell linearly over the course of the lifespan of the cell and is not bound to just cell divisions [78, 79]. In contrast, if one of those

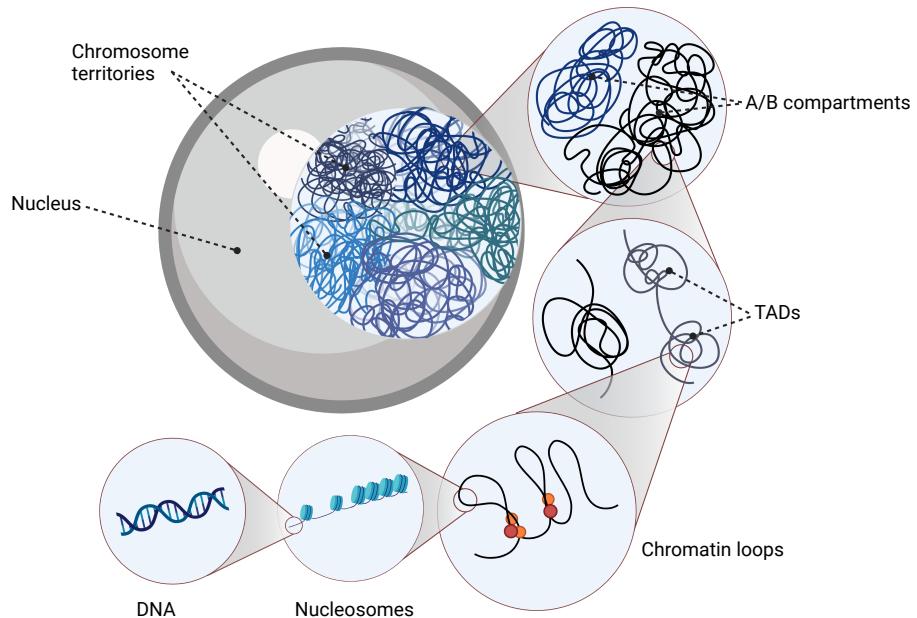


Figure 1.3: Individual chromosomes occupy a subspace in the nucleus called chromosome territories. Chromosome territories can be further partitioned to distinct A and B compartments, which are enriched for active and repressed chromatin, respectively. Genomic regions within topologically associating domains (TADs) display increased interactions, while their interactions with neighbouring regions outside of the TADs are rather limited.

mutations occurs in the germline cell, eg. sperm or egg producing cells, these mutations will be propagated to all offspring and be present in all cells of that organism and in term all its offspring. These mutations are called germline mutations. These mutations are also called germline variants, as they establish in the population and represent a variation of the organism. Mutations can also be classified depending on either their size ranging from single nucleotide polymorphisms (SNPs) over small insertions or deletions (InDels) to large structural changes, like the deletion of parts of or even a whole chromosome arm. like previously described with ploidy changes, usually smaller changes have less impact on the overall fitness of the organism, however even SNPs can lead to changes which are not compatible with life[80, 81].

1.2 Cell free DNA is more than just bits and pieces

When a cell from a multicellular organism dies, through which ever method, there will be many different enzymes involved, which clear the debris and recycle material. This means that proteases

digest proteins into amino acids, which will later be used for either building new proteins or possibly even digested further for energy production. The same happens with the DNA in the cell. However as discussed in the previous section 1.1 the DNA is wrapped around histones and organised in structures called nucleosomes. These protect the DNA from being cut by DNAases by hindering the access to the DNA, similar to how they stopped the access for transcription into RNA. This then in turn leads to the DNA being cut into pieces mainly in the length of 167 base pairs (bp). These DNA fragments, which are called cell free DNA (cfDNA), can then be detected in bodily fluids, like blood or even stool. By analysing these fragments, non invasive tests for prenatal care have been possible, as the DNA of the foetus is detectable in the mothers blood [82, 83]. Similar to the process, a cancer also sheds DNA, titled circulating tumour DNA (ctDNA), when its cells die, either through intervention of the immune system or through other forcefull processes. These ctDNA fragments can also be analysed and molecular properties measured, without even knowing the exact location of the tumour. As a blood test can be routinely performed in the clinic or even a general practitioner, the monitoring of cancer progression is significantly easier and safer than through other measures. Of course it is, similar to the prenatal test, only a proxy for the cells which are still alive, as these have not shed their DNA. Additionally the amount of shedded DNA is highly variable between tumours, with a general higher amount for later stages, so that sometimes there is almost no ctDNA present, even though the cancer is fairly advanced [84, 85].

1.3 DNA sequencing - when is next generation sequencing the current generation?

As we know the building blocks, that make DNA as well as the process and the enzymes responsible, we can synthesise DNA in vitro. By chemically modifying the nucleotides supplied to the synthesis process, the sequence of the copied strand can be analysed. The first method to make use of this used the lambda phage to fuse known ends for the primers needed for the reaction to the piece of DNA and supplied labelled nucleotides [86].This method was then superseded by "Sanger sequencing" after Frederick Sanger who with colleagues published this method in 1977, by adding dideoxynucleotides in a low concentration, the polymerase chain reaction would terminate trying to integrate these nucleotides and by labelling them radioactively or flourecently, a gel can be used to determine the sequence of a piece of DNA [87, 88], which made the method better suited for larger scale projects.

However this method has multiple issues for modern research questions. Mostly, that it is fairly labour and time consuming to analyse multiple pieces of DNA at the same time and it is very challenging to sequence all the DNA of an organism. The human genome project, which was started in 1990 used machines which automated the Sanger sequencing procedure and it still took hundreds of researchers 13 years to complete the DNA sequence of just one human [89, 64]. Even though this was a very long project, it laid the ground work for the usage of the current sequencing technologies.

1.3.1 Library preparation - what we learned from using phages

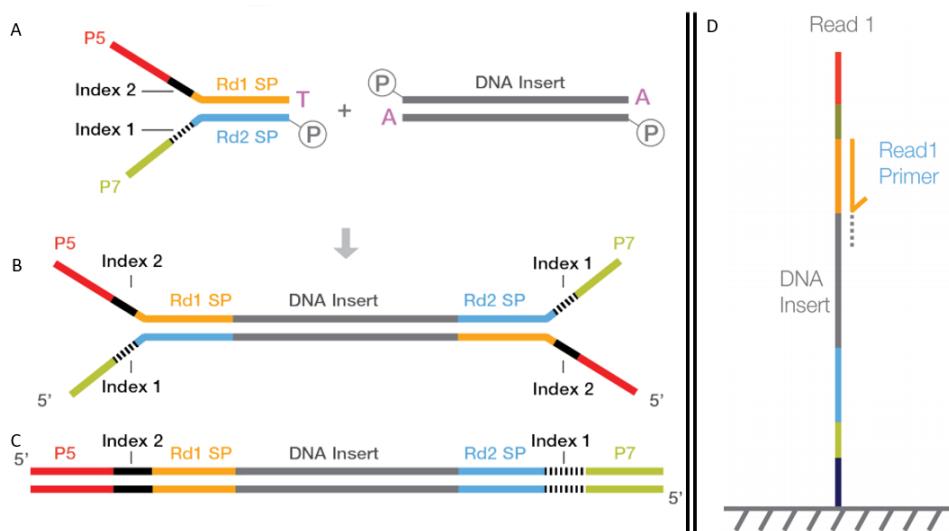


Figure 1.4: Adapter ligation during library preparation. The adapters are added to the DNA insert during library preparation. A. The DNA insert is prepared by adding an A-tail and phosphorylation. B. The adapter complex which includes the P5/P7 flow cell binding adapter is added to the DNA insert. C. The DNA insert is ready for sequencing. D. The DNA insert binds to the flow cell for sequencing. Primers bind to the DNA insert to generate reads;
Figure adapted from "How short inserts affect sequencing performance" [90]

Library preparation is the name of the preprocessing step, which is done before it is sequenced with the current technologies. The first step to sequence DNA is to obtain the DNA, which is done by lysing the cells of interest, which disrupts the cell membrane and therefore spills all its contents. The then spilled DNA is fragmented into smaller pieces, by either restriction enzymes or sonication, to have a size of about between 200-800bp. These steps are not necessary when preparing sequencing of ctDNA, as discussed in [Section 1.2](#), the DNA is unbound and already digested into short fragments. Once the DNA is ready, it is both phosphorelated as well as an A-tail is added, before the adapter complex is ligated. This enabled the DNA to bind to the flow cell which is covered with the reverse complement of the adapter ([Figure 1.4](#)).

1.3.2 Next generation sequencing

Next generation sequencing (NGS) is the coined term for basically any standard high-throughput sequencing performed, which includes exome, genome, transcriptome, protein-DNA interactions (ChIP) and other epigenome studies. The term NGS is still widely used, even though it has been more than 10 years since the first NGS approach was commercially available. While in the beginning of next generation sequencing there were multiple approaches, the current lion share (80% of sequencing data) of protocols use the Illumina short read sequencing by synthesis approach ([Figure 1.5](#))[\[91, 92\]](#), which is based on the concept of alternating integration of fluorescently labelled nucleotides and imaging with a microscope ([Figure 1.5](#)) as well as multiplexing, where a DNA fragment is ligated to an index, which allows the sequencing of multiple samples at once [\[93, 94\]](#) as it is shown in [Figure 1.4](#). This method allows highly accurate determination of the sequence of a DNA fragment and depending on the flow cell and sequencing machine allows to sequence a whole genome in just 24h.

1.3.3 Long read sequencing - the "third" generation sequencing

By now, multiple methods which broke free of the size limitations of NGS exist, which are commonly referred to as long read sequencing. Most of the current methods trade the very high accuracy of the second generation NGS methods for the capability of sequencing huge continuous strands of DNA (current record 2.3 Million bp [\[95\]](#)) with normal library preparation ranging between 10-30 Kbp. These methods are expected to revolutionise our understanding of the highly repetitive elements that exist in the genome, such as the centromeres of chromosomes. Methods such as the direct molecule sequencing approach by Oxford Nanopore are even able to distinguish post transcriptional modifications on RNA[\[96\]](#). So far, these methods however are still very expensive and as this work is dealing with ctDNA, which is highly fragmented, these methods offer only limited advantages over the short read sequencing, while being much more expensive.

1.4 DNA analysis- what to do with the sequence

The types of analysis that can be done with the output from the sequencing machine stretches far, however, all methods need to first infer the location in the genome, the sequenced piece of DNA

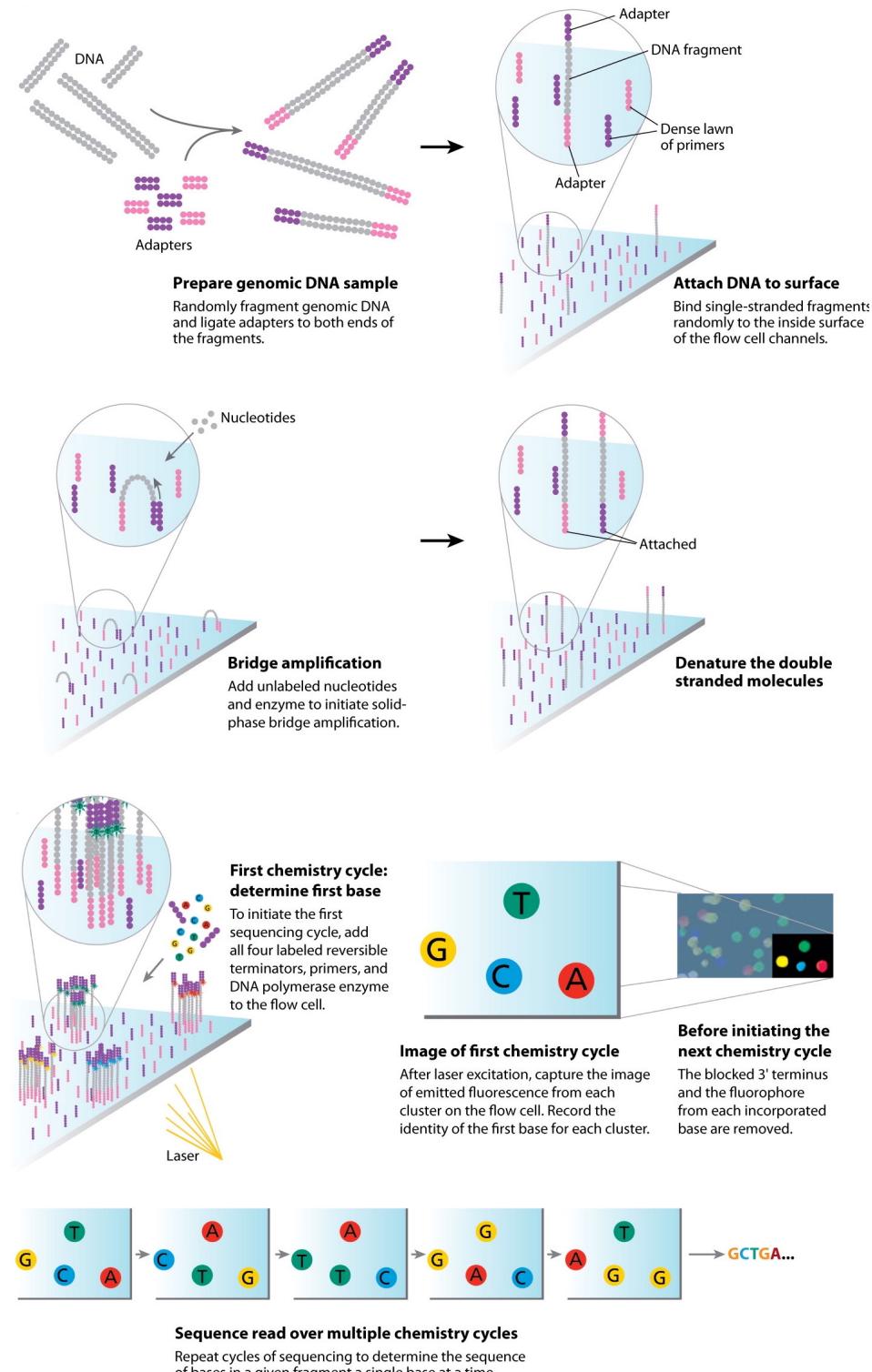


Figure 1.5: The Illumina sequencing-by-synthesis approach. Cluster strands created by bridge amplification are primed and all four fluorescently labeled, 3'-OH blocked nucleotides are added to the flow cell with DNA polymerase. The cluster strands are extended by one nucleotide. Following the incorporation step, the unused nucleotides and DNA polymerase molecules are washed away, a scan buffer is added to the flow cell, and the optics system scans each lane of the flow cell by imaging units called tiles. Once imaging is completed, chemicals that effect cleavage of the fluorescent labels and the 3'-OH blocking groups are added to the flow cell, which prepares the cluster strands for another round of fluorescent nucleotide incorporation; Figure adapted from Mardis[91]

originated from. As the current methods randomly fragment the DNA ([Section 1.3.1](#)), the genomic location information is completely lost. This process is referred to as mapping.

1.4.1 Mapping - Ey man, where is my genomic location?

In this process, the fragments of DNA, which were sequenced, are assigned a genomic coordinate on the reference genome. This is only possible, due to the fact, that we have a resolved genome sequences ([Section 1.3](#)) for a high number of species. The location a sequenced piece of DNA fits to the reference genome might be unique, but it could also fit to multiple locations, due to highly repetitive regions or due to the existence of pseudo genes with almost 100% identity. In addition to this, the reference genome might not accurately reflect the genome of the organism that has been sequenced. Each mapping position is therefore assigned a quality score, which reflects how likely it is the actual position of the sequence. As Illumina sequencers have the ability to sequence both ends of the DNA fragment, the position of the ends (read 1 and read 2) to each other can also be used to infer the quality, as they should be within a reasonable distance to each other ([Figure 1.4](#))

As this process is time consuming and the exact location of the fragment might not be as important, there exists a subset of tools called pseudo-mapper, which are based on k -mers, which are predefined DNA sequences of length k , which help to identify certain regions of interest. These tools are especially common for RNAseq, where the exact location of a read doesn't matter, only that the read is within a gene [[97](#), [98](#)], but also for methods that estimate similarity between sequences (DNA, RNA or protein) [[99](#), [100](#)].

For this work however, the exact position of reads is important, so only real mapping methods like BWA [[101](#)] or Bowtie 2 [[102](#)], which are optimised for short reads from Illumina systems, provide the necessary functions.

add things about alternative contigs and reference genome?

1.4.2 Variant calling - spot the difference

As intra-species genetic variation is intended for adaptation and evolution, there will be places where the DNA sequence of the subject will differ from the sequence of the reference (see [Section 1.1.2](#)). These variants give insight into medical background as well as treatment options for patients and can even be used to guide family planning. Depending on the type of variation that is

of interest, a different set of computational methods are needed, as germline and somatic variants have different properties.

1.4.3 Germline variant calling - the cards you have been dealt at birth

The most common source of DNA used for germline variant analysis is the mono nuclear layer from the blood of the subject, but really almost any cell can be used for this process, as all cells in the organism will share all germline variants ([Section 1.1.2](#)). The only important input on top of the DNA sequence from the sequencer are the reference genome of the organism as all variant nomenclature is based on the reference and the ploidy of the organism ([Section 1.1.1](#)). The ploidy is important to infer, at which ranges of allele frequency a variant can biologically occur. For example in a human diploid genome, germline variants can occur either in one or both chromosomes, which mean we assume reads should show an allele frequency of around 50% and 100%, where the hexaploid commercial wheat [103] allele frequency for variants would be 16%, 33%, 50%, 66%, 0.83% and 100%. Due to the random sampling and possible sequencing errors, however the observed allele frequencies will differ. Most state of the art germline variant calling method will also use haplotype reconstructions through de-Brujin graphs, which features a remapping of reads in relation to each other [104, 105, 106, 107, 108] where the original mapping location assigned by the aligner ([Section 1.4.1](#)) is only used as a guideline. This allows to resolve even complex haplotypes of the sample by not restricting the method to the linear setup of the reference genome.

1.4.4 Somatic variant calling - life is ever changing

In contrast to germline variant calling, somatic variant calling methods cannot rely on allele frequency, as not all cells sequenced are expected to have the change in nucleotide. The allele frequency is instead a measure of the sub clonal size. A subclone is here defined as the set of cells, which were derived from the cell, which originally acquired the somatic mutation. Depending on the selective advantage, just random drift and also the time point when the variant was introduced, these clones can be very variable in size and therefore their contribution to the DNA in the sequencing. As not all cells have the variants, the selection of the tissue for library preparation is very important, unlike for germline calling. The main use of somatic variant calling is the genetic diagnosis and research of cancer samples, where the main question is, which changes are present in the tumour, which lead to the disease.

The ideal scenario for tumour somatic variant calling is when a biopsy of the tumour as well as a normal sample of the patient is available. In most clinical cases, this will be the diagnostic biopsy as well as the mono nuclear layer from blood, just like for germline calling ([Section 1.4.3](#)). This needs to be adjusted depending on the type of malignancy, because if the tumour is a leukemia, the mono nuclear layer of the blood might contain tumour cells, but for solid tumours, the blood is a routine, minimally invasive option. These two samples are then analysed together and only changes that are only in the somatic tumour sample and not in the normal sample are reported. Even though this concept sounds simple, there are some pitfalls [[109](#)]. First of all, there might be some tumour contamination in the normal sample, which needs to be adjusted for [[106, 110](#)]. Second, there might be normal “contamination” in the tumour sample, this means that not all cells in the tumour sample are actually tumour. This means that the signal of the tumour changes is reduced and harder to find.

All of these issues are amplified in the case, when there is no “normal” sample available, either because the patient didn't consent, due to other medical issues, or because for diagnostic tests there usually is no need for a germline sample. In this case, there is the option for “tumour only” variant calling, which requires a database of germline variants in the population, to distinguish between somatic and germline variants, as the variant calling is very similar to just germline variant calling ([Section 1.4.3](#)) without the restriction of the ploidy. However, even with an extensive database like gnomAD [[111](#)] it is unlikely that it is possible to remove all germline variants from the analysis and as there is no direct comparison, the precision of the “tumour only” method is significantly lower [[112](#)].

1.5 Cancer

While cancer is a massively heterogenous disease, as it can arise through a multitude of ways in almost any tissue, there are some fundamental defining features, which most, if not all malignancies acquire, before they are truly cancers. The original characteristics comprise 1. Sustaining proliferative signaling 2. Evading growth suppressors 3. Activating invasion and metastasis 4. Enabling replicative immortality 5. Inducing angiogenesis 6. Resisting cell death ([Figure 1.6](#)). These hallmarks were long considered the core of tumour development and the authors themselves hypothesised, that these core mechanics allow us to condense the complexity that cancer displays, both in the clinic as well as in labs, with a small set of rules, which all cancers have to obey [[113](#)].

In their exact words: “We foresee cancer research developing into a logical science, where the complexities of the disease, described in the laboratory and clinic, will become understandable in terms of a small number of underlying principles”

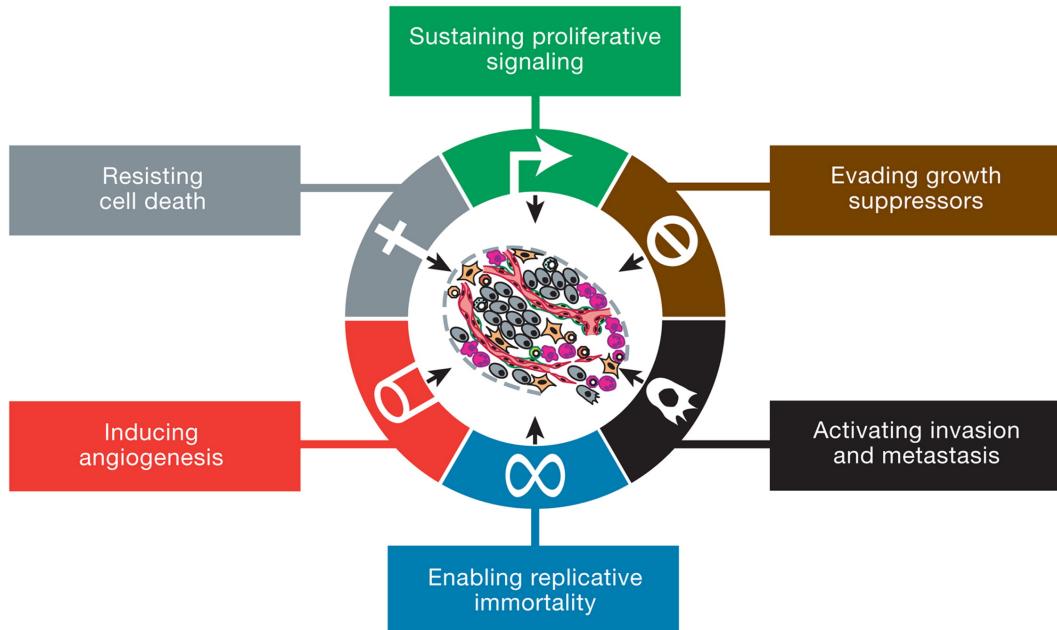


Figure 1.6: Acquired capabilities of cancer; Functional capabilities acquired by most cancers during their development; Figure adapted from Hanahan et al.[113]

However, with 11 years of additional research into the topic, more hallmarks have been found and the original list was revised by the authors to contain additional characteristics, namely 1. Avoiding immune destruction 2. Tumour-promoting inflammation 3. Genome instability and mutation 4. Deregulating cellular energetics (Figure 1.7) [114]. And even then a few years later, even more hallmarks e.g. metabolic rewiring are now considered a part of the characteristics of cancer [115].

So while the original set of the hallmarks was not sufficient or complete, it offered a good attempt at abstraction of biological concepts to describe cancers. In the following pages, I will outline each of those hallmarks and how it influences my research.

describe the hallmarks

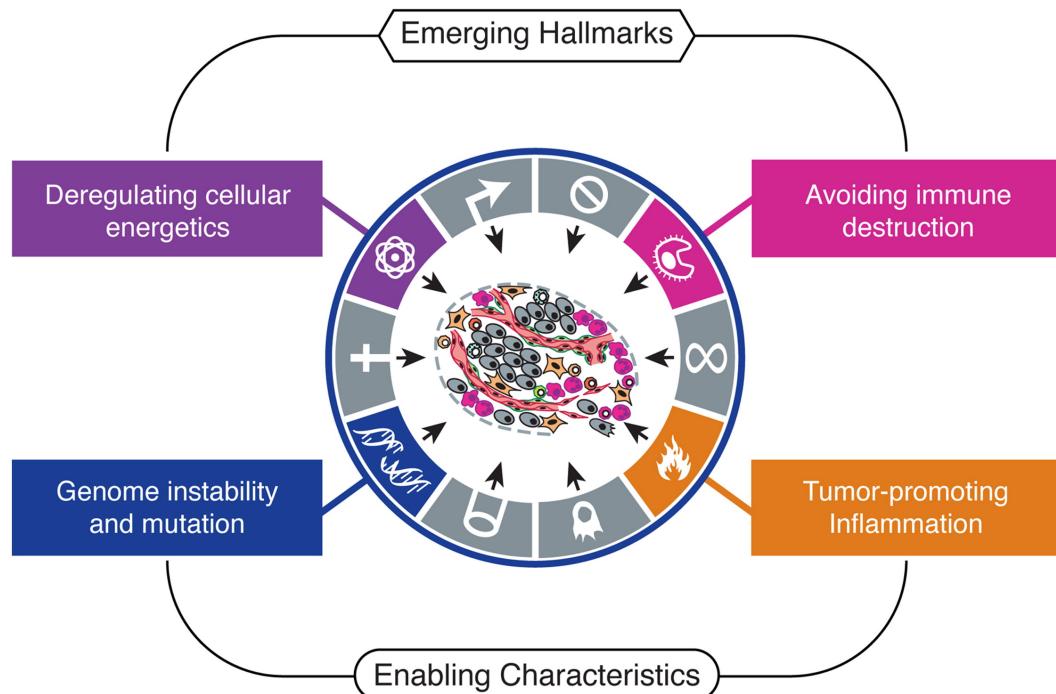


Figure 1.7: Emerging hallmarks and enabling characteristics of cancer; updated version of the hallmarks figure (Figure 1.6) with ; Figure adapted from Hanahan et al.[114]

1.5.0.1 Sustaining proliferative signaling - there are no breaks on the train

1.5.0.2 Evading growth suppressors

1.5.0.3 Activating invasion and metastasis - look at me... I am the organism now

1.5.0.4 Enabling replicative immortality

1.5.0.5 Inducing angiogenesis

1.5.0.6 Resisting cell death

1.5.1 Lungcancer

With around 1.6 million deaths world-wide each year, lung cancer is the number one cause of cancer death [116]. Every year about twelve thousand Australians get diagnosed with lung cancer. These cases can be generally split into two groups: small cell lung cancers (SCLC) and non-small cell lung cancers (NSCLC), which account for around 15% and 85% of cases, respectively. The majority of NSCLC are either lung adenocarcinoma or lung squamous cell carcinoma [117]. Even though smoking is highly associated with lung cancers, there is a big group of never smokers, with a high risk of

lung cancers in East Asia, especially women, which is correlated with outside influences like pollution and occupational carcinogens and paired with genetic susceptibility [118]. This group usually shows *EGFR* (epidermal growth factor receptor) driven tumours. *EGFR* is a transmembrane receptor tyrosine kinase, which is usually only expressed in epithelial, mesenchymal, and neurogenic tissue, but its overexpression in other tissues is a hallmark of many human malignancies, not just NSCLC.

Possibly change this to cancer in general

1.6 Overview

add short description of each chapter

“It is the main source of our mistakes, when making making decision, that we only look at life piece by piece and not as a whole.“

— Lucius Annaeus Seneca, *Epistulae morales ad Lucilium*

2

Joint somatic variant calling - if germline can do it,
so can we

2.1 Introduction

When I started exploring the somatic variant calling methods in the beginning of my PhD in 2018, I was surprised about the stark difference between germline and somatic variant calling methods. Where all "modern" germline variant callers, like Strelka2 [106], HaplotypeCaller [119], DRAGEN [120] and DeepVariant [121], have the built-in capability to joint call multiple related samples, for example from family trios, virtually no somatic variant caller had this function.

The joint analysis of smaller cohorts improves the performance of germline variant calling methods significantly, by allowing to assess technical artifacts, which might be unique for the individual sequencing machine or the researcher handling the DNA [122, 123]. Additionally, as certain parts of the genome are more problematic to sequence (Section 1.3) and map (Section 1.4.1), a “control“ sample can help to distinguish if a certain observed change occurs commonly is a technical issue or in fact a real change.

For somatic variant calling, this concept has been adopted on in the genome analysis toolkit (GATK) [124] to allow the use of panel of normals (PON), which contains frequently seen changes in healthy (“normal“) individuals analysed with the same sequencing technology [125]. However in contrast to the more intricate model for the germline equivalent, this is a post processing step of the analysis. Mutect2, which is the most recent somatic variant calling algorithm provided by the Broad institute, however also provides a multi-sample mode, for which all tumour samples need to be from the same patient, either longitudinal or spatial different [126]. This mode is hidden quite well and all

tutorials published by the developers state that “there is currently no way to perform joint calling for somatic variant discovery” [109], so while all methods in the GATK are considered a beta feature, the multi sample mode needs to be used with care.

There are only two methods currently, which have documented and published capabilities to jointly analyse tumour samples from the same patient to call somatic variants. The first one is a specialised method built on a joint bayesian model for SNVs to occur in multiple samples called multiSNV [127]. However it has multiple shortcomings, which make it not usable for our data. First, as the name suggests, the method can only jointly evaluate SNVs and completely ignores INDELs and structural variants, which would be acceptable for the superior performance it provides. However, multiSNV was optimised only for WES and not for the very deep WGS that is now available and part of this work. This mismatch of input types means exceptionally high runtimes on our data. Even with custom parallelisation that was attempted in this work, the predicted runtime for just one multi sample patient would have been longer than 3 years. This shows, that while multiSNV was a great step forward at the time, there is a real need for new methods to stem the tide of sequencing data available due to the ever decreasing sequencing cost.

multiSNV has been the only software available for multi sample analysis for almost five years, but during this work, superFreq [128] was published. It combines all standard analysis steps for tumour analysis, like quality assessment, variant calling, copy number analysis and clonal deconvolution, into one program and is even able to jointly analyse samples. However similar to multiSNV, its focus during optimisation and development was on WES and RNAseq data, so when applied to our data, we could not find a server node with enough memory to execute the workflow.

This then prompted us to investigate possible workflows to enable the analysis of high depth WGS, which we estimate to become more and more normal, with the ever dropping prices of sequencing. The following sections will first show the publication in part and then the impact of the joint analysis on downstream methods (Section 2.3) then discuss additional analysis done after the the publication of the manuscript (Section 2.4) and end with a section about the usage of the methods by others (Section 2.5).

2.2 Publication

The full publication about joint somatic variant calling can be found at <https://doi.org/10.1093/bioinformatics/btab606> and non-journal formatted version is also attached as Appendix A with all

supplementary methods.

However in this section, I will parts of the paper to make this work as standalone and easy to read as possible. References to supplementary data will be prefixed with the letter A in the text.

2.2.1 Summary

To enable highly sensitive, fast and accurate variant detection from multiple related tumour samples, we have developed joint variant calling extensions to two widely used single-sample algorithms, FreeBayes [104] and Strelka2 [106]. Using both simulated and clinical sequencing data, we show that these workflows are highly accurate and can detect variants at much lower variant allele frequencies than other commonly used methods.

2.2.2 FreeBayesSomatic workflow

The original FreeBayes algorithm can jointly evaluate multiple samples but routinely it does not perform somatic variant calling on tumour-normal pairs. We introduce FreeBayesSomatic which allows concurrent analysis of multiple tumour samples by adapting concepts from SpeedSeq [129] which differentiates the likelihood of a variant between tumour and normal samples instead of imposing an absolute filter for all variants called in the normal. Hence, for each genotype (GT) at SNV sites, FreeBayesSomatic first calculates the difference in likelihoods (LOD) between the normal ([Equation 2.1](#)) and the tumour ([Equation 2.2](#)) samples genotype likelihoods (GL) with g_0 describing the reference genotype.

$$\text{LOD}_{\text{normal}} = \max_{g_i \in \text{GT}} (\text{GL}(g_0) - \text{GL}(g_i)) \quad (2.1)$$

$$\text{LOD}_{\text{tumour}} = \min_{s \in \text{Samples}} \left(\min_{g_i \in \text{GT}} (\text{GL}_s(g_i) - \text{GL}_s(g_0)) \right) \quad (2.2)$$

$$\text{somaticLOD} := (\text{LOD}_{\text{normal}} \geq 3.5 \wedge \text{LOD}_{\text{tumour}} \geq 3.5) \quad (2.3)$$

Next, the variant allele frequencies (VAF) in both the tumour and the normal samples are compared at each site.

$$\text{VAF}_{\text{tumour}} = \max_{s \in \text{Samples}} (\text{VAF}_s) \quad (2.4)$$

$$\begin{aligned} \text{somaticVAF} := & (\text{VAF}_{\text{normal}} \leq 0.001 \vee \\ & (\text{VAF}_{\text{tumour}} \geq 2.7 \cdot \text{VAF}_{\text{normal}})) \end{aligned} \quad (2.5)$$

A variant is classified as somatic when both somaticLOD as well as somatic VAF pass the criteria somaticLOD (Equation 2.3) and somaticVAF (Equation 2.5).

The thresholds chosen for both LOD and VAF calculations were previously fitted by the blue-collar bioinformatics workflow for the DREAM synthetic 3 dataset using the SpeedSeq likelihood difference approach [130] and were selected to identify high confidence variants.

2.2.3 Strelka2Pass workflow

In contrast to FreeBayes, whilst Strelka2 has a multiple-sample mode for germline analysis and tumour-normal pair somatic variant calling capabilities, it cannot jointly analyse multiple related tumour samples. We enable this feature by adapting a two-pass strategy previously used for RNA-seq data [131]. First, somatic variants are called from each tumour-normal pair. All detected variants across the cohort are then used as input for the second pass of the analysis where we re-iterate through each tumour-normal pair but assess allelic information for all input genomic sites.

The method re-evaluates the likelihood of each variant, by integrating every genotype from each tumour-normal pair. This step can "call" a variant (v) in a sample that initially did not present enough evidence to pass the Strelka2 internal filtering using two conditions: 1) if this variant was called as a proper "PASS" by Strelka2 in any other tumour sample, or 2) if the integrated evidence for this variant across all tumour-normal pairs reached a sufficiently high level. The second condition was based on the somatic evidence score (SomEVS) reported by Strelka2, which is the logarithm of the probability of the variant v being an artefact.

$$p_{\text{error}}(v) = 10^{\left(\frac{-\text{SomEVS}(v)}{10}\right)} \quad (2.6)$$

While the germline sample is shared between all processes, we can approximate these individual probabilities as being independent, since one variant calling process is agnostic of the other. Hence, we derive the following:

$$p_{error}(v_{s_1}, v_{s_2}, \dots, v_{s_n}) = \prod_{s \in \text{Samples}} p_{error}(v_s) \quad (2.7)$$

And therefore:

$$\text{SomEVS}(v_{s_1}, v_{s_2}, \dots, v_{s_n}) = \sum_{s \in \text{Samples}} \text{SomEVS}(v_s) \quad (2.8)$$

This allows the summation (Equation 2.8) of the SomEVS score across all supporting variants to assign a "PASS" filter, if it reached a joint SomEVS score threshold. This threshold can be set by the user and is 20 by default, which corresponds to an estimated error rate of 1%. These "recovered" variants need to pass a set of additional quality metrics related to depth of coverage, mapping quality and read position rank sum score.

As an additional improvement, we also built multiallelic support into Strelka2 which originally only reports the most prevalent variant at a specific site. Within the two-pass analysis, we reconstruct the available evidence for a multiallelic variant at a called site from the allele-specific read counts and report the minor allele at this site, if there is sufficient support from other samples. This method allows recovery of minor alleles only if another sample has this variant called by Strelka2, as SomEVS scores are not available for minor alleles.

2.2.4 Validation

As new methods are always a good idea, because they challenge previous assumptions and can be a step forward for the field by just showing that simpler models can have the same performance, all methods should be validated against the gold standard methods in the field with data which allows objective measures to be used. For germline variant calling there have been multiple challenges and specifically designed test datasets, but these predefined datasets do not exist for the somatic equivalent. This issue is even more pronounced, as we do not only need a tumour-normal pair, but we need the multiple tumour samples in the dataset to be related. To allow a fair comparison, I first generated a fully synthetic dataset, where every variant is known and fully defined (Section 2.2.4.1)

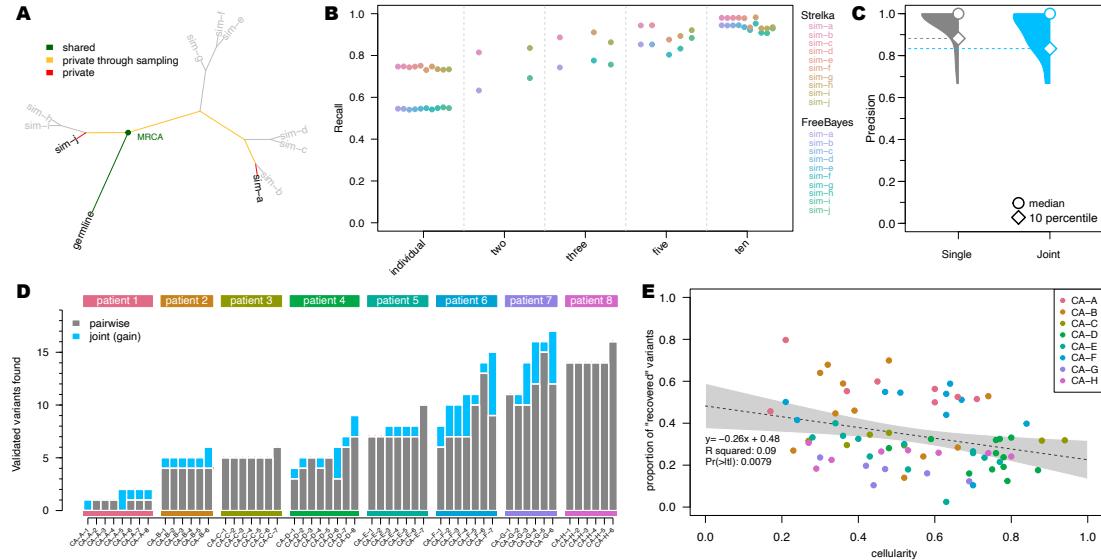


Figure 2.1: Comparison of joint multi-sample variant calling and single tumour-normal paired calling methods; A) Simulated phylogeny highlighting two samples with high evolutionary distance (sim-a and sim-j) where MRCA denotes the most recent common ancestor. B) Recall estimates of FreeBayes and Strelka2, run in individual tumour-normal paired and joint calling configurations using two (sim-a and sim-j), three (sim-a, sim-g and sim-j), five (sim-a, sim-c, sim-f, sim-h and sim-j) and all ten tumour samples. C) Precision of Strelka2 and D) Number of variants called by Strelka2 run in both tumour-normal paired (grey) and added with joint calling configurations (blue), which have been validated by targeted amplicon sequencing (TAS). E) Correlation between cellularity and proportion of variants found only with joint calling using Strelka2Pass for clinical samples; grey area shows the "95%" confidence interval for the linear model fit (dotted line).

to allow a general performance assessment of the methods. Then to ensure that these metrics also hold true in real world data, we then re-analysed previously published datasets which have orthogonal validation in the form of targeted amplicon sequencing (TAS) (Section 2.2.4.2).

2.2.4.1 Simulated data

We first simulated a phylogeny with somatic and germline variants from ten tumour samples and one normal (Figure 2.1A and Figure A.2A, B). Germline variants were simulated at a uniform allele frequency of 0.5. Somatic VAFs were sampled from a custom distribution, modelled to favour low allele frequency variants to closely represent real world data (min VAF: 0.001; max VAF: 1; Fig. S1C, D). Paired-end sequencing reads with realistic error profiles were simulated for WGS data at 160X average coverage using the ART-MountRainier software [132]. The simulated reads were aligned to GRCh38 and both germline and somatic variants from the phylogeny were spiked into the aligned reads using Bamsurgeon [133]. We compared the workflows for FreeBayes and Strelka2 with and

without our extensions for joint variant calling on the simulated datasets. The performance of Mutect2 joint variant calling was also assessed using its proposed best practice workflow. As both Mutect2 and FreeBayes do not return a verdict for each individual sample, we needed to assign each sample in the multi-sample VCF its own FILTER value. We called a somatic variant as present in a sample, if there were at least two reads supporting it for this sample and the overall FILTER showed a "PASS", which was the same cut-off used in the refiltering step in the Strelka2-pass workflow.

While the precision of each method without our extensions was greater than 99.8%, they all missed at least 25% of all variants in the samples (i.e recall \leq 75%). In contrast, the recall of the modified workflows increased to \approx 95% with only a minute decrease in the precision for both FreeBayes and Strelka2 ([Figure A.3](#)). Mutect2 however, had virtually no change in precision, but the recall actually decreased from \approx 75% to \approx 41% when analysing the samples jointly ([Figure A.3B](#)). Additionally, with our modified workflows, true positive variants were called with VAFs as low as 0.008 (median detected VAF \geq 0.14 for joint sample analysis and \geq 0.21 for single tumour-normal pair analysis), enabling improved distinction between true variants and technical errors ([Figure A.4](#)). This improvement in performance for Strelka2 is only achieved after the refiltering step and not just a result of the second pass ([Figure A.5](#), [Section A.5.4](#)).

The performance of joint variant calling in Mutect2 was inferior compared to all other methods ([Figure A.3A, B](#)). This was primarily due to the "clustered_events" filter in Mutect2, which excluded the majority of false negative variants, with negligible contribution to the exclusion of true negative variants ([Figure A.6A, B](#)). This result was unexpected as the simulated variants were evenly distributed along the genome and the corresponding allele frequencies were sampled randomly ([Figure A.2D](#)).

Since the extent of the improvement in our joint calling workflows is bound by the number of shared variants between samples, we sub-sampled the simulated dataset, to show the effect of incomplete sampling on our methods, which is more likely in clinical settings. Furthermore, the evolutionary distance between the related samples in addition to the number of samples, has a major impact on the number of shared variants, as only variants acquired between the germline and the most recent common ancestor (MRCA), will benefit from the joint analysis. Therefore, we selected three sample subsets which included two, three and five samples with high evolutionary distance to show the minimum expected improvement ([Figure 2.1A, B](#)). There was a clear linear improvement for both FreeBayesSomatic and Strelka2Pass when increasing the number of samples even if they had a distant evolutionary relationship. In contrast, when using only two samples with a small evolutionary

distance, the increase in performance was almost as large as when jointly analysing all 10 available samples. This shows that samples with a high number of shared variants will perform better in joint calling workflows ([Figure A.7](#)).

2.2.4.2 Clinical data

To validate the performance of our new workflows, we then analysed WGS and whole-exome sequencing (WES) data of multi-region tumour samples from eight patients, with multiple tumour sites (average 7 samples per patient; total number of samples 55), enrolled in a rapid autopsy program conducted at the Peter MacCallum Cancer Centre ([Table A.1](#) and [Section A.5.2](#)) [[134](#), [135](#)]. The published studies had multiple somatic variants from the clinical samples orthogonally validated through targeted amplicon sequencing (TAS). We used these TAS-validated variants as the gold standard to evaluate the performance of different workflows, acknowledging that the technical biases inherent to TAS data are different to those present in WGS and WES ([Figure A.8](#)) and that there would be sampling biases depending on different tumour cells analysed in each data type.

In concordance with the results of the simulated data, our improved workflows found additional variants in all but one patient ([Figure 2.1D](#), [Figure A.9](#)) (total additional variants Strelka2Pass: 64; FreeBayesSomatic: 85) with only a slight drop in precision for FreeBayesSomatic (mean: 0.94 vs. 0.88) and Strelka2Pass (mean: 0.97 vs. 0.92). Since the panel of variants validated by TAS was limited (7108 bp for patients CA-B through -H), this increase in detected variants suggests that a high number of shared variants in samples are missed with current approaches, which in turn leads to an overestimation of tumour heterogeneity between samples, as these variants are thought to not be present rather than undetected.

Even though the number of shared variants is a major influencing factor when jointly calling variants, low cellularity samples benefit more from the joint calling, as conventional methods cannot reliably distinguish low allele frequency variants from noise. Through a joint analysis approach, the number of recovered variants is higher in low cellularity samples, which indicates, that especially for clinical samples with variable tumour purity, joint analysis can have a major impact on improving performance ([Figure 2.1E](#), [Figure A.10](#)).

Mutect2 in contrast, did not show significant improvement in any sample in its joint calling configuration, but showed inferior performance compared to the tumour-normal pairwise approach in two samples ([Figure A.9E](#)), similar to its decreased performance in the simulated data ([Figure A.3](#)).

This was due to true variants being removed by the internal filters of the tool ([Figure A.6](#)C, D). This is in stark contrast to our novel workflows, where the joint analysis preserves all called sites from the pairwise method and finds additional variants. Overall, Mutect2 found less validated variants in all patients than both Strelka2Pass (mean: 2.2) and FreeBayesSomatic (mean: 2.5) with comparable levels of precision ([Figure A.9](#), [Figure A.11](#)) but longer run times ([Table A.2](#)).

Our improved workflow also enabled the discovery of multiallelic variants with Strelka2, which led to the discovery of on average 42 additional variants (min: 1; max: 535) in the analysed WES and 987 additional variants in the WGS (min: 81; max 2329). These variants are strong indicators of sub clonal structure and are invaluable for the study of evolutionary trajectories in cancer as shown in the following sections.

2.3 Effects on downstream analysis - not quite the missing link, but close

The ability to find additional shared variants has significant impact on our understanding of cancer evolution and the timing of initiation and metastatic seeding. Recent work has shown, that similar to the well known genetic heterogeneity, there is heterogeneity when it comes to the timing of metastatic seeding. While traditionally it was thought that tumours only metastasise after they reach a certain size, to escape the restrictions of the niche, like reduced nutrition, recent publications showed, there is also very early metastatic seeding [[136](#)]. But all methods analysing heterogeneity and evolutionary timing and history are ultimately based on the somatic variants found in the data, so if we improve on the input of the analysis methods, we can expect a clearer and possibly more granular result.

In the following sections I will highlight how big the effect can be for methods, like phylogenetic reconstruction and clonal decomposition, which use somatic variants as input.

2.3.1 Phylogenetic reconstruction

As this work is not about the advantages and shortcomings of different phylogenetic reconstruction tools, I will not show a comprehensive amount of these tools, but rather focus on the results. For this reason, I chose to use neighbour joining (NJ) [[137](#)], because it is fast, readily available in most phylogenetic reconstruction tool kits and if the input distance is correct, the output will be correct. And

even, if the distance is not 100% correct, if the distance is “nearly additive” and the input distances are not far off from the real distance, the tree topology will still be reconstructed correctly [138]. Lastly, in contrast to many other methods like UPGMA and WPGMA [139], NJ does not assume an equal mutation rate of each sample, because we know, that the molecular clock hypothesis [140] is not valid for different lineages of cancers [141].

The only thing that NJ requires as an input is a distance matrix of all samples, so the next step was the selection of the right distance metric. While there are lots of distance measures for DNA sequences, which allow accounting for different probabilities of transitions and transversions as well as uneven base composition, models like F81 [142] or HKY85 [143] are only really designed for germline mutations and are not easily applicable for subclonal somatic mutations, which is why I decided to first transform the variants present in all samples into a binary occurrence vector and then calculating the Hamming distance [144] between all samples. This generates a maximum parsimony approach and the branch length of the trees will be directly translatable to the amount of variants which are different between samples.

[Figure 2.2](#) shows both the reconstructed phylogenies of the autopsy samples of the late stage melanoma patient “CA-F“ from the manuscript ([Appendix A, Table A.1](#)), using the variants found with the default tumour-normal method on the left and our improved joint method on the right. The exact same reconstruction methodology was used otherwise, such that only the difference inputs lead to the final difference.

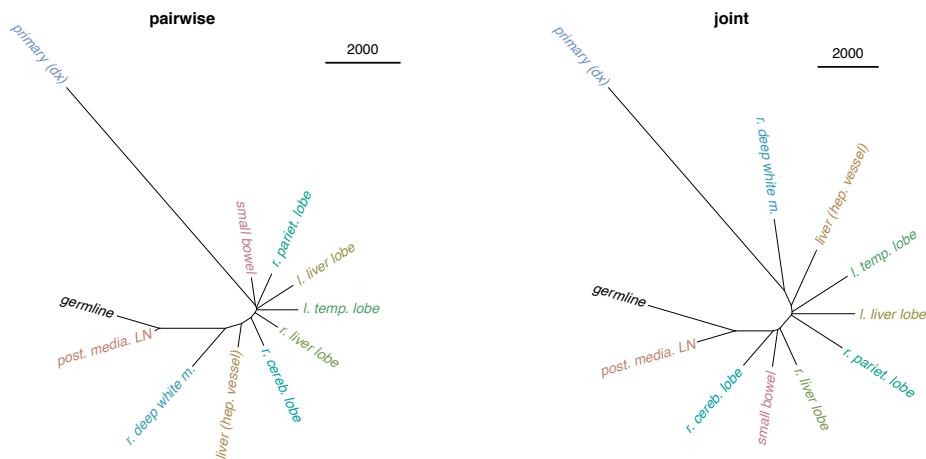


Figure 2.2: Reconstructed phylogenies of a patient with multiple spatially distinct samples; Neighbour joining on Hamming distance on variant occurrence vector. Tip labels describe the location of the sample in the patient. Trees are shown as unrooted with germline as fixated origin point; black line ruler shows the length of an edge with 2000 mutations

Maybe adjust the font size in the trees to make it more readable

There are several obvious changes, first, the longer edge connecting the germline, which we consider as the state of no somatic variants, to all other samples. This shows that there are many more shared mutations in all samples, than what would have been anticipated with the default method, which corresponds to an overestimation of the heterogeneity of the samples. As the accumulation of somatic variants is still used as a proxy for timing and cell divisions, when assuming a high mutation rate for lung cancer ($5.3 \cdot 10^{-8}$ from Werner et al. [145]) this difference of ≈ 36000 variants is equivalent to ≈ 2000 cell divisions. While the cell doubling rate of lung cancers is highly dependent on the type [146], this difference makes a huge difference when assessing the timing of the tumour initiation and further evolution.

Secondly, there have been topological changes, which generate a longer bifurcating edge between the olive coloured “r. liver lobe” and the “r. pariet. lobe” showing a bottle neck in cancer evolution, which fits very well with the clinical history, where the patient lived with stable disease for almost ten years, before progressing and dying. The extreme distance of the primary/diagnostic sample from the rest of the samples could be either a difference in sequencing quality, or due to the exposure to FFPE for the ten years between tumour diagnosis and death. However, as this feature is preserved between both the joint and the pairwise analysis it is no result of our new method.

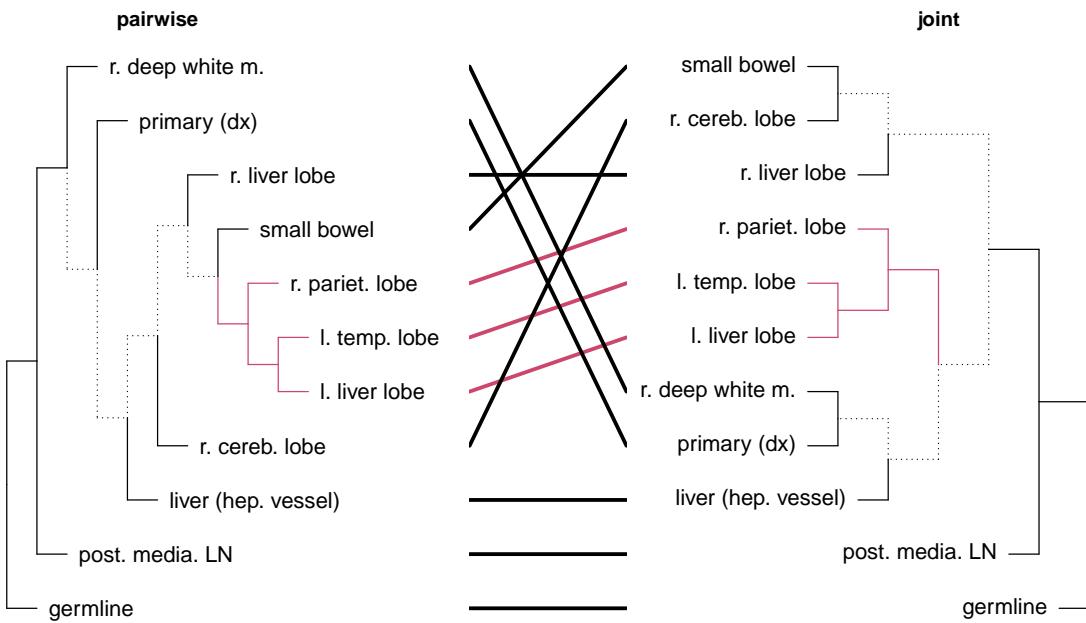


Figure 2.3: Side by side view of the reconstructed trees from Figure 2.2; internal edges, which are distinct between trees are shown as dotted lines; common subtree is shown in red Tree labels have been sorted to minimise distance between labels; Visualisation generated with dendextend [25]

maybe increase the line width of the edges

[Figure 2.3](#) shows a topology focused view of the two trees, which highlights the breaks which are needed to morph one tree into the other with dotted edges [147]. The common subtrees are coloured the same on both sides and connecting lines show identical labels. This format shows that while the trees look quite similar at first glance, they show vastly different topologies.

One example of this is “small bowel” which was connected to the red common subtree, but is now much closer to the “r. cereb. lobe” and forms a parallel clade with the “r.liver lobe”. In general, where the pairwise tree shows a very linear topology, which leaves only branching out of the main with no disjunct subclades, which are clearly present in the joint reconstruction. ([Figure 2.3](#)).

2.4 Longitudinal analysis - something for the ages

While the initial motivation for the development of these workflows was the analysis of multi-region, so spatial, samples from the same patient coming from the CASCADE rapid autopsy program, a longitudinal application of these methods for the joint analysis of diagnostic and relapse sample, or even the repeated testing of ctDNA are quite worth thinking about. In this part, I will present work using the published workflows on a longitudinal dataset, which highlights the flexibility and wide spread usability of the new methods.

In addition to their autopsy, Patient “CA-F“ also had three longitudinal blood samples taken, from which ctDNA was extracted and WES performed. In a study of late stage melanoma patients, Tan et al. identified ctDNA sequencing as a way to stratify patients into high and low risk of relapse and therefore inform adjuvant therapy [148], which makes this patient a prime example to show the improvement with joint variant calling. Similar to the spatially related samples, the joint analysis can improve the performance, which then in term enables the detection of lower allele frequency variants, either through lower tumour burden or through the limited availability of DNA fragments from brain lesions due to the blood brain barrier [149].

To show that even in longitudinal data, the joint analysis can boost the signal, we jointly variant called the diagnostic biopsy with only the three ctDNA samples and compared them with the results from the pairwise analysis. On average we found 2905 additional variants in each of the ctDNA samples which is more than double than the average number of variants found with the pairwise analysis (2414). Out of those, we found 534 variants in the ctDNA samples, which were found as

a high confidence variant in the diagnostic sample, indicating that these findings are high quality calls.

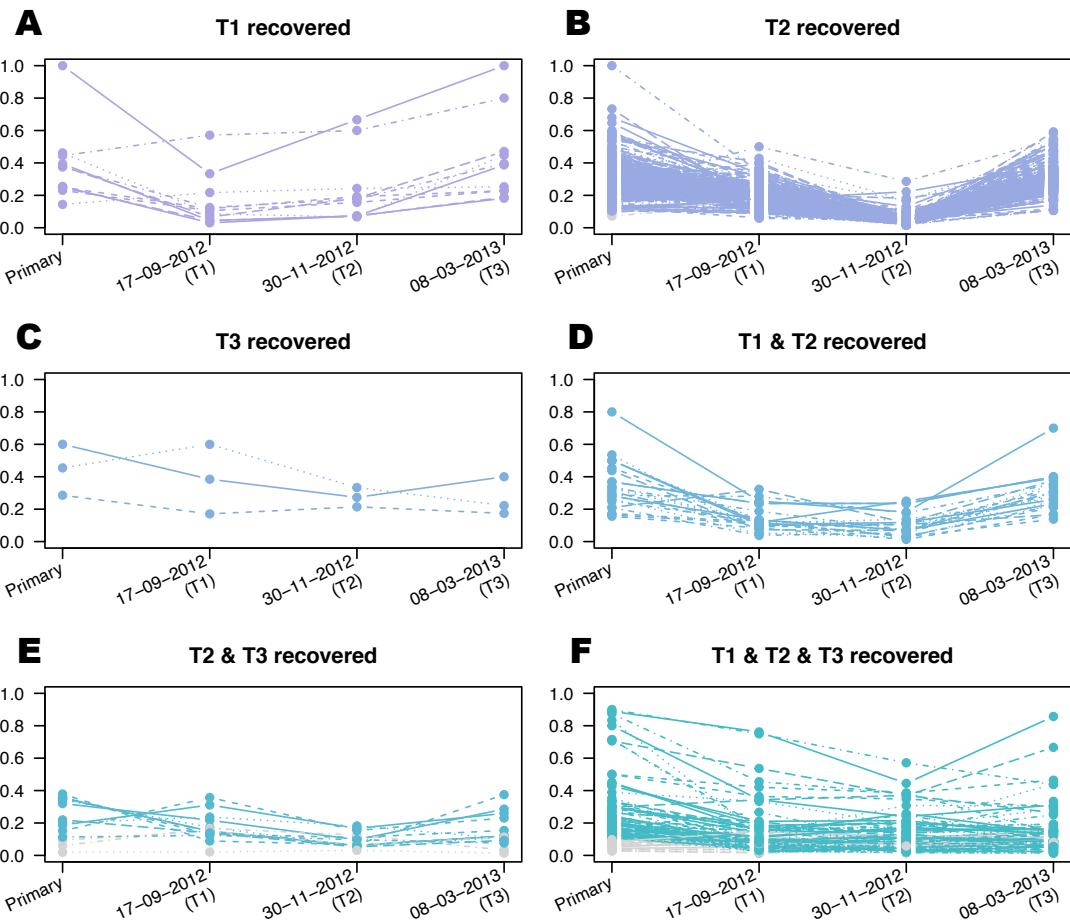


Figure 2.4: Improved somatic variant calling in longitudinal data: Variant allele frequency (VAF) of variants found additionally through joint variant calling which were found as high confidence variants in the primary sample; Variants with less than 0.1 VAF in the primary are coloured grey; “T1 recovered“ shows variants, which were high confidence in all ctDNA samples but T1 and were only found through joint calling there; Axis label show the date of blood collection

Exactly like in the spatially different samples, in longitudinal data lower tumour purity samples benefit more from the joint analysis. We see that time point 2 (T2) has the highest amount of recovered variants (377) which are found as high confidence variants in both other time points (Figure 2.4 A vs. B vs. C) and T2 also has the lowest tumour purity in the cfDNA recorded (T1: 60%; T2: 20%; T3: 60%) however, there are still 106 variants, which were not found in the ctDNA samples at all with the pairwise analysis at all, even though they were high confidence variants in the primary sample (Figure 2.4 F). These variants usually show a lower depth of coverage (dp) in the ctDNA samples, which might indicate a problematic region in the genome, but rather than it not being called a variant, it is just a sign of incomplete data, which can be used with our joint approach.

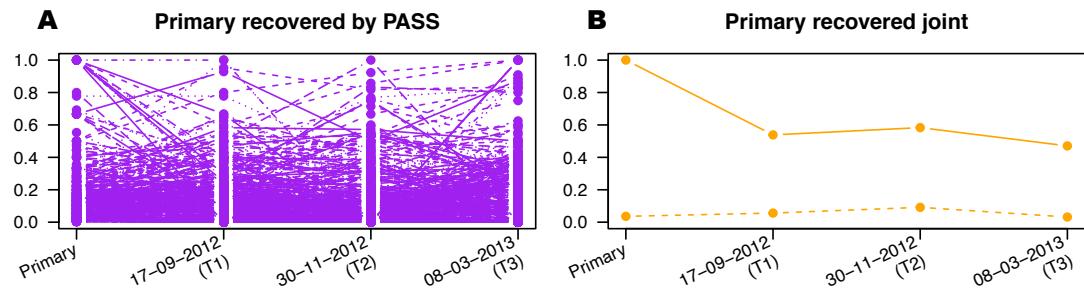


Figure 2.5: Longitudinal data informs diagnostic variant calling: Vafs of variants additionally found through joint calling in the primary samples; Primary recovered by PASS shows variants which were high confidence in at least one ctDNA sample; Primary recovered joint shows variant which were low confidence in all samples in the pairwise analysis; Axis label show the date of blood collection

Finally, we can also find 398 additional variants in the primary sample. 398 were discarded due to missing data in the tissue sequencing, but could be found with a high confidence in the longitudinal data and two of the variants were included, as all 4 samples had this variant below the detection threshold (Figure 2.5). The missing depth in the primary also leads to the occasional very high allele frequency of the variant, as all available reads show the variant, but their numbers are below the threshold normal variant callers will report variants.

This shows that both spatially and longitudinal related samples should be analysed jointly, as it substantially increases the amount of true variants found, which as shown before have a big impact on downstream analysis of the samples.

2.4.1 Clonal deconvolution

Finally, the holy grail of analysis of multiple related samples from the same patient is the clonal deconvolution, where subclonal reoccurring patterns of mutations (clones) are resolved both spatially and longitudinally. These reoccurring clones can be linked to either parallel evolution through positive selection pressure, like a targeted drug, or to the process of developing Metastases where a piece of the cancer “breaks” off and grows at a different site. Surprisingly, as it shares the same issue as the joint somatic variant calling of needing deeply sequenced data from multiple samples of the same organism, there is a plethora of algorithms and methods available for clonal deconvolution. Since 2015 PhyloWGS [150], Canopy [151], CLOE [152], CloneFinder [153], MACHINA [154] and MOBSTER [155] were published, to name a few. Underlying all these models is a form of clustering variants with similar variant allele frequency together, to reduce the combinatorial space and enhance the confidence in the signal [156]. However due to the high number of tools, the challenge to

select the right tool is substantial, especially since all of them have up and downsides [157]. In this work I decided to use PhylogicNDT [158] as it has been shown to work well on clinical samples [159] and does not have the restriction for the input to be from copy number neutral areas which many of the other tools have.

The analysis was conducted by transforming the variants called with the joint workflows as well as the default pairwise analysis into the required file format without the cancer cell fraction part, in order to let PhylogicNDT calculate those on the fly. The local copy numbers for each variants are generated by intersecting the called segments from sequenza with the position of the variants. Variants which could not be assigned a local copy number change were discarded from the analysis.

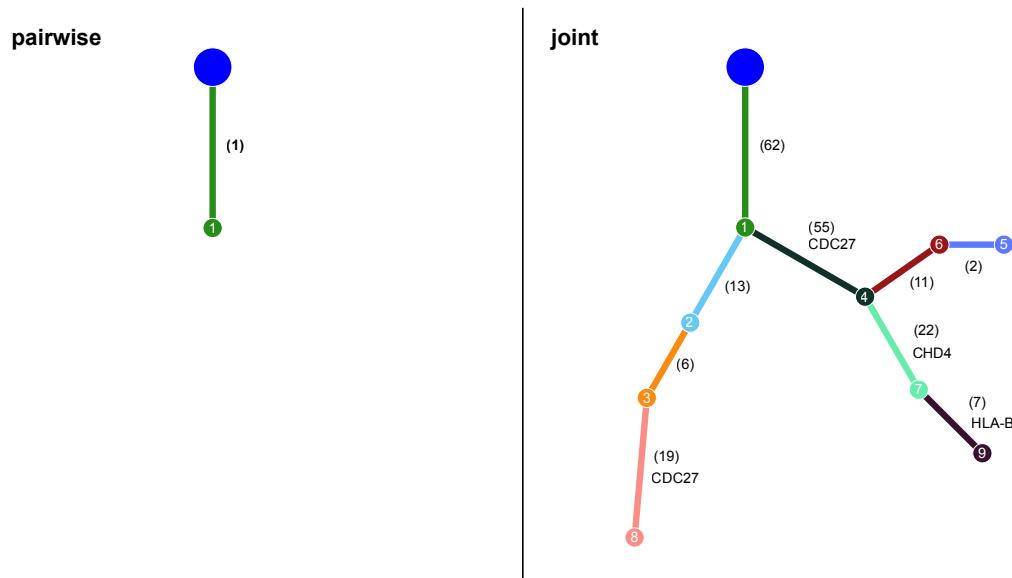


Figure 2.6: Reconstructed clonal trees from PhylogicNDT; Blue circle at top depicts the germline/normal state. The coloured edges with the same coloured circle represents a distinct subclone of the parent from which the edge emerges; The number in braces next to the edge is the number of mutations which define this subclone with an added gene symbol added, if there is a cancer driver gene mutation. The left part shows the result when using the default pairwise method of Strelka2 and the right side uses the results from the Strelka2Pass workflow

[Figure 2.6](#) shows the highest parsimony clonal tree reconstructed by PhylogicNDT for the pairwise as well as the joint variant calling. As the copynumber calling information is the same for both inputs, the only difference is in the called variants. While there was no subclonal structure detected at all for the pairwise analysis, there is a highly variable structure detected using the jointly called variants. As this is a clinical sample, we cannot 100% be sure that the more branched model is the actual truth, but its biologically more logical that a late stage cancer has developed several subclones, rather than it being a very homogenous disease at all of the 10 sites at autopsy with no evolution

over ten years of disease [159]. It is of particular interest, that the *CDC27* gene got mutated at different timepoints in different clones (clone 8 vs. clone 4), which a clear sign of parallel evolution, which would definitely be missed without the joint analysis.

2.4.2 Longitudinal enriched phylogeny

Of course it is finally also possible to build a phylogeny with both the spatial tissue samples as well as the longitudinal ctDNA samples. However, as the ctDNA give a holistic view of all cancer metastases (Section 1.2) the interpretation needs to accommodate for that.

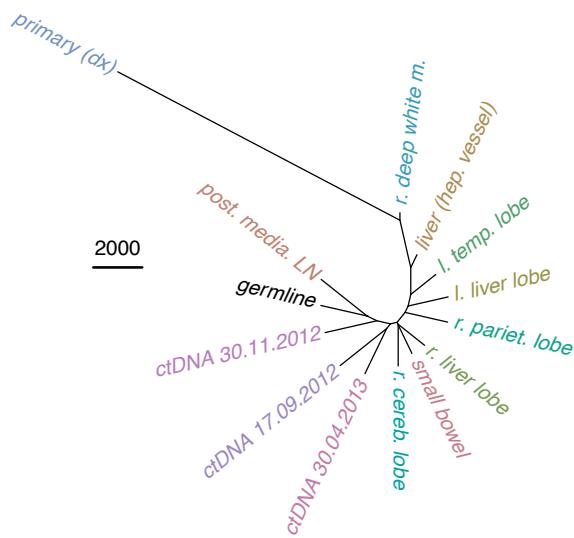


Figure 2.7: Reconstructed phylogeny with longitudinal ctDNA samples: Tree from Figure 2.2 with three additional ctDNA samples from different time points about one year prior to death. Ruler shows the equivalent of 2000 mutations

The maybe most surprising thing is that the more temporally distant ctDNA samples from 17.09.2012 and 30.04.2013 are in a subclade together, away from the “ctDNA 30.11.2012“ sample. Secondly, the addition of the ctDNA samples also lead to a further bipartition edge, which separates “r. liver lobe“, “small bowel“ and “r. cereb. lobe“ from the rest of the tree (Figure 2.7). This was already inferable from the topology of the previous tree in Figure 2.3 “joint“, but is even more pronounced with the inclusion of the ctDNA samples.

This shows again, that the addition of more samples helps to refine and improve the trajectory and history of cancer samples and it is vital to do this analysis jointly to generate the optimal result.

2.5 Usage - its not just me that thinks it is good

As amazing as published open source software is, its real value is in the re-usability and portability. Many published software packages are not maintained or not even functional even though they are published. While I developed these joint somatic variant calling workflows to deal with a challenge I faced, many other people even just at the same institute have already shown interest and some research groups are already using the software to analyse their multi-sample data.

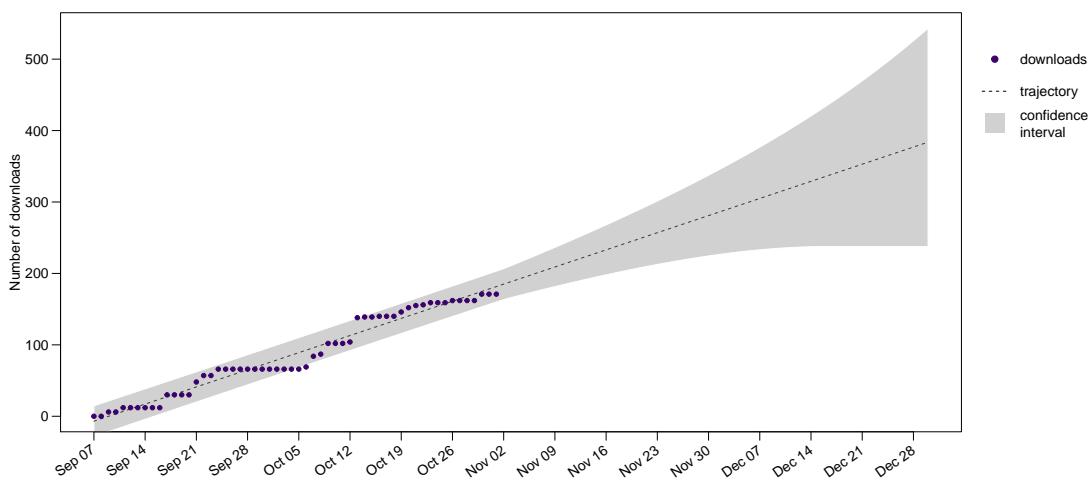


Figure 2.8: Cumulative download numbers of the “dawsontoolkit” docker container since publication of the manuscript; Actual counts are shown as dots, with smoothed trajectory depicted as dotted line with the 95% confidence interval shown as a grey background; confidence interval has been adjusted with exponential decay of prediction accuracy with distance from the last data point; Start date 7th September 2021

To have some proxy of the usage statistics of the workflows, I recorded the download numbers of the “dawsontoolkit” docker container after the publication of the manuscript. The container only consists of software for refiltering and joint analysis of the workflows . Obviously, this is an imperfect measurement, as people can reuse a downloaded container as often as they want, which would not appear in the count and similarly, just because the container was downloaded, the analysis might not have been used. However it still shows an interaction and an interest in the methods. The download numbers show a sustained stable increase in downloads (Figure 2.8). This suggests, that there is a need in the methods, rather than a simple curiosity after publication, which hopefully will facilitate a higher quality analysis of future projects and therefore lead to a better understanding of cancer evolution and heterogeneity.

Update the plot at a later timepoint)

“Death is a release from and an end of all pains: beyond it our sufferings cannot extend: it restores us to the peaceful rest in which we lay before we were born”

— Lucius Annaeus Seneca, *De Consolatione ad Marciam*

3

CASCADE - Late stage lung cancer in the spotlight

3.1 Introduction

3.2 Publication

This chapter includes the data analysis for two publications. The first publication features the resistance mechanism of small cell transformation ([https://doi.org/10.1016/j.ccell.2019.08.008\[160\]](https://doi.org/10.1016/j.ccell.2019.08.008[160])) and the second shows the discovery of resistance to a targeted RET-fusion driven cancer ([https://doi.org/10.1016/j.jtho.2020.01.006\[134\]](https://doi.org/10.1016/j.jtho.2020.01.006[134]))

Cant include papers like this, will have to write the chapter as a whole

3.3 Cohort analysis

3.4 Mitochondrial phylogenetic reconstruction - the power house of the phylogenies

3.5 Outlook

“Many a mickle makes a muckle.”

— proverb

4

MisMatchFinder - hope springs eternal

4.1 Introduction

“As you think, so you become. Our busy minds are forever jumping to conclusions, manufacturing and interpreting signs that aren’t there.“

— Epictetus, *The Enchiridion*

5

Conclusion

A

Custom workflows to improve joint variant calling from multiple related tumour samples: FreeBayesSomatic and Strelka2Pass

This appendix contains the manuscript published at *Bioinformatics* in a non journal style format with the supplementary methods and figures. It can also be found at [10.1093/bioinformatics/btab606/6361543](https://doi.org/10.1093/bioinformatics/btab606/6361543) for a paper style version.

Hollizeck S.^{1,2}, Wong S.Q.^{1,2}, Solomon B.^{1,2}, Chandrananda D.^{1,2}, and Dawson S-J.^{1,2,3,*}

¹ Peter MacCallum Cancer Centre, Melbourne 3000, Victoria, Australia

² Sir Peter MacCallum Department of Oncology, University of Melbourne, Melbourne 3000, Victoria, Australia

³ Centre for Cancer Research, University of Melbourne, Melbourne 3000, Victoria, Australia

* D.C and S.J.D are co-senior authors and contributed equally to this article

Received on 27-Jan-2021; revised on 13-Jul-2021; accepted on 12-Aug-2021

Abstract

Summary: This work describes two novel workflows for variant calling that extend the widely used algorithms of Strelka2 and FreeBayes to call somatic mutations from multiple related tumour samples and one matched normal sample. We show that these workflows offer higher precision and

recall than their single tumour-normal pair equivalents in both simulated and clinical sequencing data.

Availability and Implementation: Source code freely available at the following link: <https://atlassian.petermac.org.au/bitbucket/projects/DAW/repos/multisamplevariantcalling> and executable through Janis (<https://github.com/PMCC-BioinformaticsCore/janis>) under the GPLv3 licence.

Contact: Dineika.Chandrananda@petermac.org, Sarah-Jane.Dawson@petermac.org

Supplementary information: Supplementary data are available at *Bioinformatics* online.

A.1 Introduction

Joint variant calling methods are routinely used to call germline variants by leveraging population-wide information across multiple related samples [161, 162]. This concept is also advantageous for somatic variant calling to potentially overcome the challenges of spatial heterogeneity and low tumour purity. However, there is a critical lack of robust algorithms that allow multi-sample somatic calling. Most studies still rely on variant calling of separate tumour-normal pairs, subsequently combining the results across a sample cohort [136, 3, 163].

There are two major pitfalls for combining variants called from individual tumour samples. First, it is very difficult to differentiate between a false negative result due to "missing data" versus the true absence of a variant. Second, there is limited sensitivity for low allele frequency variants thus, decreasing the ability to detect minor clones, particularly in samples with low tumour purity.

Currently, only three algorithms claim to have the functionality to jointly analyse multiple samples: multiSNV [127], SuperFreq [128], and Mutect2 [107], each presenting different limitations. For instance, multiSNV cannot call indels and along with SuperFreq, is not optimised for analysis of deep coverage whole-genome sequencing (WGS) data. Mutect2 has previously been shown to be disadvantageously conservative as well as computationally inefficient [164].

To enable highly sensitive, fast and accurate variant detection from multiple related tumour samples, we have developed joint variant calling extensions to two widely used single-sample algorithms, FreeBayes [104] and Strelka2 [106]. Using both simulated and clinical sequencing data, we show that these workflows are highly accurate and can detect variants at much lower variant allele frequencies than commonly used methods.

A.2 Materials and methods

A.2.1 FreeBayesSomatic workflow

The original FreeBayes algorithm can jointly evaluate multiple samples but routinely it does not perform somatic variant calling on tumour-normal pairs. We introduce FreeBayesSomatic which allows concurrent analysis of multiple tumour samples by adapting concepts from SpeedSeq [129] which differentiates the likelihood of a variant between tumour and normal samples instead of imposing an absolute filter for all variants called in the normal. Hence, for each genotype (GT) at SNV sites, FreeBayesSomatic first calculates the difference in likelihoods (LOD) between the normal (Equation A.1) and the tumour (Equation A.2) samples genotype likelihoods (GL) with g_0 describing the reference genotype.

$$\text{LOD}_{\text{normal}} = \max_{g_i \in \text{GT}} (\text{GL}(g_0) - \text{GL}(g_i)) \quad (\text{A.1})$$

$$\text{LOD}_{\text{tumour}} = \min_{s \in \text{Samples}} \left(\min_{g_i \in \text{GT}} (\text{GL}_s(g_i) - \text{GL}_s(g_0)) \right) \quad (\text{A.2})$$

$$\text{somaticLOD} := (\text{LOD}_{\text{normal}} \geq 3.5 \wedge \text{LOD}_{\text{tumour}} \geq 3.5) \quad (\text{A.3})$$

Next, the variant allele frequencies (VAF) in both the tumour and the normal samples are compared at each site.

$$\text{VAF}_{\text{tumour}} = \max_{s \in \text{Samples}} (\text{VAF}_s) \quad (\text{A.4})$$

$$\begin{aligned} \text{somaticVAF} := & (\text{VAF}_{\text{normal}} \leq 0.001 \vee \\ & (\text{VAF}_{\text{tumour}} \geq 2.7 \cdot \text{VAF}_{\text{normal}})) \end{aligned} \quad (\text{A.5})$$

A variant is classified as somatic when both somaticLOD as well as somatic VAF pass the criteria somaticLOD (Equation A.3) and somaticVAF (Equation A.5).

The thresholds chosen for both LOD and VAF calculations were previously fitted by the blue-collar bioinformatics workflow for the DREAM synthetic 3 dataset using the SpeedSeq likelihood difference approach [130] and were selected to identify high confidence variants.

A.2.2 Strelka2Pass workflow

In contrast to FreeBayes, whilst Strelka2 has a multiple-sample mode for germline analysis and tumour-normal pair somatic variant calling capabilities, it cannot jointly analyse multiple related tumour samples. We enable this feature by adapting a two-pass strategy previously used for RNA-seq data [131]. First, somatic variants are called from each tumour-normal pair. All detected variants across the cohort are then used as input for the second pass of the analysis where we re-iterate through each tumour-normal pair but assess allelic information for all input genomic sites.

The method re-evaluates the likelihood of each variant, by integrating every genotype from each tumour-normal pair. This step can "call" a variant (v) in a sample that initially did not present enough evidence to pass the Strelka2 internal filtering using two conditions: 1) if this variant was called as a proper "PASS" by Strelka2 in any other tumour sample, or 2) if the integrated evidence for this variant across all tumour-normal pairs reached a sufficiently high level. The second condition was based on the somatic evidence score (SomEVS) reported by Strelka2, which is the logarithm of the probability of the variant v being an artefact.

$$p_{error}(v) = 10^{\left(\frac{-\text{SomEVS}(v)}{10}\right)} \quad (\text{A.6})$$

While the germline sample is shared between all processes, we can approximate these individual probabilities as being independent, since one variant calling process is agnostic of the other. Hence, we derive the following:

$$p_{error}(v_{s_1}, v_{s_2}, \dots, v_{s_n}) = \prod_{s \in \text{Samples}} p_{error}(v_s) \quad (\text{A.7})$$

And therefore:

$$\text{SomEVS}(v_{s_1}, v_{s_2}, \dots, v_{s_n}) = \sum_{s \in \text{Samples}} \text{SomEVS}(v_s) \quad (\text{A.8})$$

This allows the summation (Equation A.8) of the SomEVS score across all supporting variants to assign a "PASS" filter, if it reached a joint SomEVS score threshold. This threshold can be set by the user and is 20 by default, which corresponds to an estimated error rate of 1%. These "recovered" variants

need to pass a set of additional quality metrics related to depth of coverage, mapping quality and read position rank sum score.

As an additional improvement, we also built multiallelic support into Strelka2 which originally only reports the most prevalent variant at a specific site. Within the two-pass analysis, we reconstruct the available evidence for a multiallelic variant at a called site from the allele-specific read counts and report the minor allele at this site, if there is sufficient support from other samples. This method allows recovery of minor alleles only if another sample has this variant called by Strelka2, as SomEVS scores are not available for minor alleles.

A.3 Validation

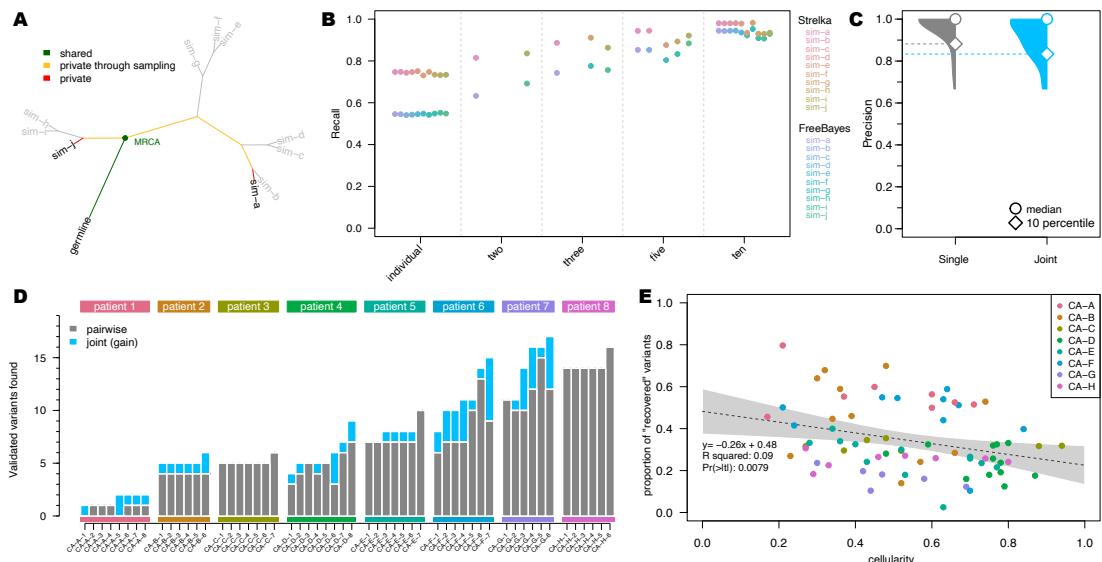


Figure A.1: Comparison of joint multi-sample variant calling and single tumour-normal paired calling methods; A) Simulated phylogeny highlighting two samples with high evolutionary distance (sim-a and sim-j) where MRCA denotes the most recent common ancestor. B) Recall estimates of FreeBayes and Strelka2, run in individual tumour-normal paired and joint calling configurations using two (sim-a and sim-j), three (sim-a, sim-g and sim-j), five (sim-a, sim-c, sim-f, sim-h and sim-j) and all ten tumour samples. C) Precision of Strelka2 and D) Number of variants called by Strelka2 run in both tumour-normal paired (grey) and added with joint calling configurations (blue), which have been validated by targeted amplicon sequencing (TAS). E) Correlation between cellularity and proportion of variants found only with joint calling using Strelka2Pass for clinical samples; grey area shows the "95%" confidence interval for the linear model fit (dotted line).

A.3.1 Simulated data

We first simulated a phylogeny with somatic and germline variants from ten tumour samples and one normal (Fig. 1A, S1A, B) (Supplementary methods). Germline variants were simulated at a uniform allele frequency of 0.5. Somatic VAFs were sampled from a custom distribution, modelled to favour low allele frequency variants to closely represent real world data (min VAF: 0.001; max VAF: 1; Fig. S1C, D). Paired-end sequencing reads with realistic error profiles were simulated for WGS data at 160X average coverage using the ART-MountRainier software [132]. The simulated reads were aligned to GRCh38 and both germline and somatic variants from the phylogeny were spiked into the aligned reads using Bamsurgeon [133]. We compared the workflows for FreeBayes and Strelka2 with and without our extensions for joint variant calling on the simulated datasets. The performance of Mutect2 joint variant calling was also assessed using its proposed best practice workflow. As both Mutect2 and FreeBayes do not return a verdict for each individual sample, we needed to assign each sample in the multi-sample VCF its own FILTER value. We called a somatic variant as present in a sample, if there were at least two reads supporting it for this sample and the overall FILTER showed a “PASS”, which was the same cut-off used in the refiltering step in the Strelka2-pass workflow.

While the precision of each method without our extensions was greater than 99.8%, they all missed at least 25% of all variants in the samples (i.e recall $\leq 75\%$). In contrast, the recall of the modified workflows increased to $\approx 95\%$ with only a minute decrease in the precision for both FreeBayes and Strelka2 (Fig. S2). Mutect2 however, had virtually no change in precision, but the recall actually decreased from $\approx 75\%$ to $\approx 41\%$ when analysing the samples jointly (Fig. S2B). Additionally, with our modified workflows, true positive variants were called with VAFs as low as 0.008 (median detected VAF ≥ 0.14 for joint sample analysis and ≥ 0.21 for single tumour-normal pair analysis), enabling improved distinction between true variants and technical errors (Fig. S3). This improvement in performance for Strelka2 is only achieved after the refiltering step and not just a result of the second pass (Fig. S4) (Supplementary Methods).

The performance of joint variant calling in Mutect2 was inferior compared to all other methods (Fig. S2A, B). This was primarily due to the "clustered_events" filter in Mutect2, which excluded the majority of false negative variants, with negligible contribution to the exclusion of true negative variants (Fig. S5A, B). This result was unexpected as the simulated variants were evenly distributed along the genome and the corresponding allele frequencies were sampled randomly (Fig. S1D).

Since the extent of the improvement in our joint calling workflows is bound by the number of shared variants between samples, we sub-sampled the simulated dataset, to show the effect of incomplete sampling on our methods, which is more likely in clinical settings. Furthermore, the evolutionary distance between the related samples in addition to the number of samples, has a major impact on the number of shared variants, as only variants acquired between the germline and the most recent common ancestor (MRCA), will benefit from the joint analysis. Therefore, we selected three sample subsets which included two, three and five samples with high evolutionary distance to show the minimum expected improvement (Fig. 1A, B). There was a clear linear improvement for both FreeBayesSomatic and Strelka2Pass when increasing the number of samples even if they had a distant evolutionary relationship. In contrast, when using only two samples with a small evolutionary distance, the increase in performance was almost as large as when jointly analysing all 10 available samples. This shows that samples with a high number of shared variants will perform better in joint calling workflows (Fig. S6).

A.3.2 Clinical data

To validate the performance of our new workflows, we then analysed WGS and whole-exome sequencing (WES) data of multi-region tumour samples from eight patients, with multiple tumour sites (average 7 samples per patient; total number of samples 55), enrolled in a rapid autopsy program conducted at the Peter MacCallum Cancer Centre (Table S1 and Supplementary methods) [134, 135]. The published studies had multiple somatic variants from the clinical samples orthogonally validated through targeted amplicon sequencing (TAS). We used these TAS-validated variants as the gold standard to evaluate the performance of different workflows, acknowledging that the technical biases inherent to TAS data are different to those present in WGS and WES (Fig. S7) and that there would be sampling biases depending on different tumour cells analysed in each data type.

In concordance with the results of the simulated data, our improved workflows found additional variants in all but one patient (Fig. 1D, S8) (total additional variants Strelka2Pass: 64; FreeBayesSomatic: 85) with only a slight drop in precision for FreeBayesSomatic (mean: 0.94 vs. 0.88) and Strelka2Pass (mean: 0.97 vs. 0.92). Since the panel of variants validated by TAS was limited (7108 bp for patients CA-B through -H), this increase in detected variants suggests that a high number of shared variants in samples are missed with current approaches, which in turn leads to an overestimation of tumour heterogeneity between samples, as these variants are thought to not be present rather than undetected.

Even though the number of shared variants is a major influencing factor when jointly calling variants, low cellularity samples benefit more from the joint calling, as conventional methods cannot reliably distinguish low allele frequency variants from noise. Through a joint analysis approach, the number of recovered variants is higher in low cellularity samples, which indicates, that especially for clinical samples with variable tumour purity, joint analysis can have a major impact on improving performance (Fig. 1E, S9).

Mutect2 in contrast, did not show significant improvement in any sample in its joint calling configuration, but showed inferior performance compared to the tumour-normal pairwise approach in two samples (Fig. S8E), similar to its decreased performance in the simulated data (Fig. S2). This was due to true variants being removed by the internal filters of the tool (Fig. S5C, D). This is in stark contrast to our novel workflows, where the joint analysis preserves all called sites from the pairwise method and finds additional variants. Overall, Mutect2 found less validated variants in all patients than both Strelka2Pass (mean: 2.2) and FreeBayesSomatic (mean: 2.5) with comparable levels of precision (Fig. S8, S10) but longer run times (Table S2).

Our improved workflow also enabled the discovery of multiallelic variants with Strelka2, which led to the discovery of on average 42 additional variants (min: 1; max: 535) in the analysed WES and 987 additional variants in the WGS (min: 81; max 2329). These variants are strong indicators of sub clonal structure and could be invaluable for the study of evolutionary trajectories in cancer.

A.4 Discussion

Here we present an extension to two widely used variant callers, enabling them to analyse multiple related tumour samples and improve the sensitivity of detecting low allele frequency variants. This is highly relevant in clinical settings where low tumour purities in samples is a common occurrence. These workflows are an important step to satisfy the current unmet need for multi-sample tumour variant calling. While we have showcased their improvements in patient sequencing data, additional validation on larger clinical datasets is warranted to ensure the methodology performs robustly in real world settings. Importantly, these workflows are fully containerised and can be run through Janis [165] on almost any high-performance computing environment, as well as cloud services. Each workflow is highly optimised and parallelised to facilitate the analysis of the large

amount of data joint variant calling requires. The workflow specification also allows the easy adjustment of parameters to enable customisation for the user's needs and priorities, whereas building an ensemble workflow using multiple callers is up to the discretion of the user (Fig. S11).

Acknowledgements

The authors would like to thank all patients who provided tissue samples utilised in this study. The authors acknowledge Dr Lavinia Tan for assistance provided with the collection of patient clinical samples.

Funding

This work was supported by the National Health and Medical Research Council [grant numbers 1196755 to S.J.D, 1158345 to S.J.D and B.J.S, 1194783 to S.Q.W, 1173450 to B.J.S]; and CSL Centenary Fellowship to S.J.D; Victorian Cancer Agency [grant numbers 19008 to D.C, 19002 to S.Q.W]

Conflicts of Interest

S.J.D has been a member of advisory boards for AstraZeneca and Inivata. The S.J.D. lab has received funding from Cancer Therapeutics CRC and Roche-Genentech. B.J.S. has been a member of advisory boards for AstraZeneca, Roche-Genentech, Pfizer, Novartis, Amgen, Bristol Myers Squibb and Merk

Data availability

The simulated data and the respective final variant calling files underlying this article are available from Figshare at <https://melbourne.figshare.com>, and can be accessed with <https://doi.org/10.26188/13635186> for the dataset and <https://doi.org/10.26188/13635187> for the called variants.

The biological data underlying this article are available at the European Genome-Phenome Archive (EGA) at <https://ega-archive.org>, and can be accessed with study id [EGAS00001004023](https://ega-archive.org/study/EGAS00001004023) and [EGAS00001004950](https://ega-archive.org/study/EGAS00001004950).

Supplementary data

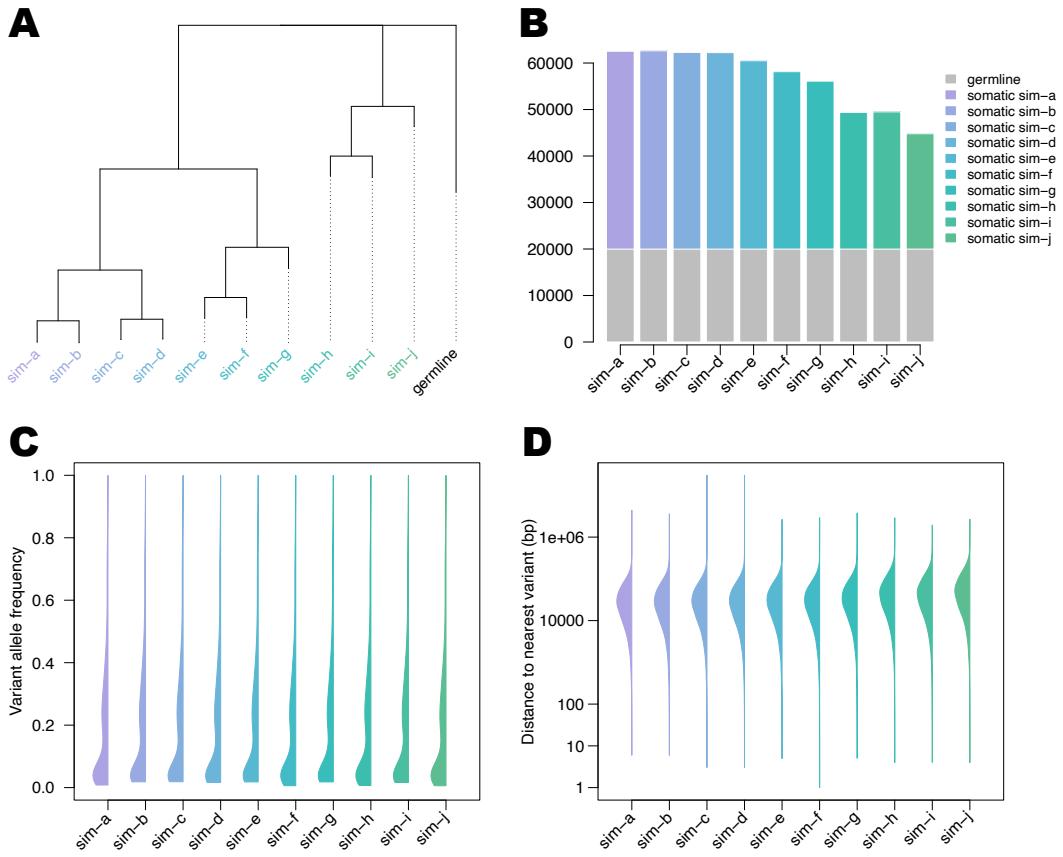


Figure A.2: Characteristics of simulated data: A) Simulated phylogeny of samples B) Number of simulated germline and somatic variants per sample C) Variant allele frequency distribution of simulated variants per sample D) Distance to nearest variant in each sample.

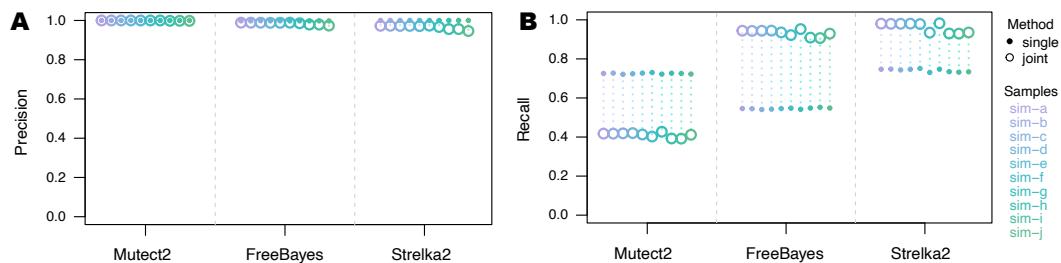


Figure A.3: Performance of workflows using simulated data: A) Precision and B) Recall of Mutect2, FreeBayes and Strelka2, run in single tumour-normal paired and joint calling configurations.

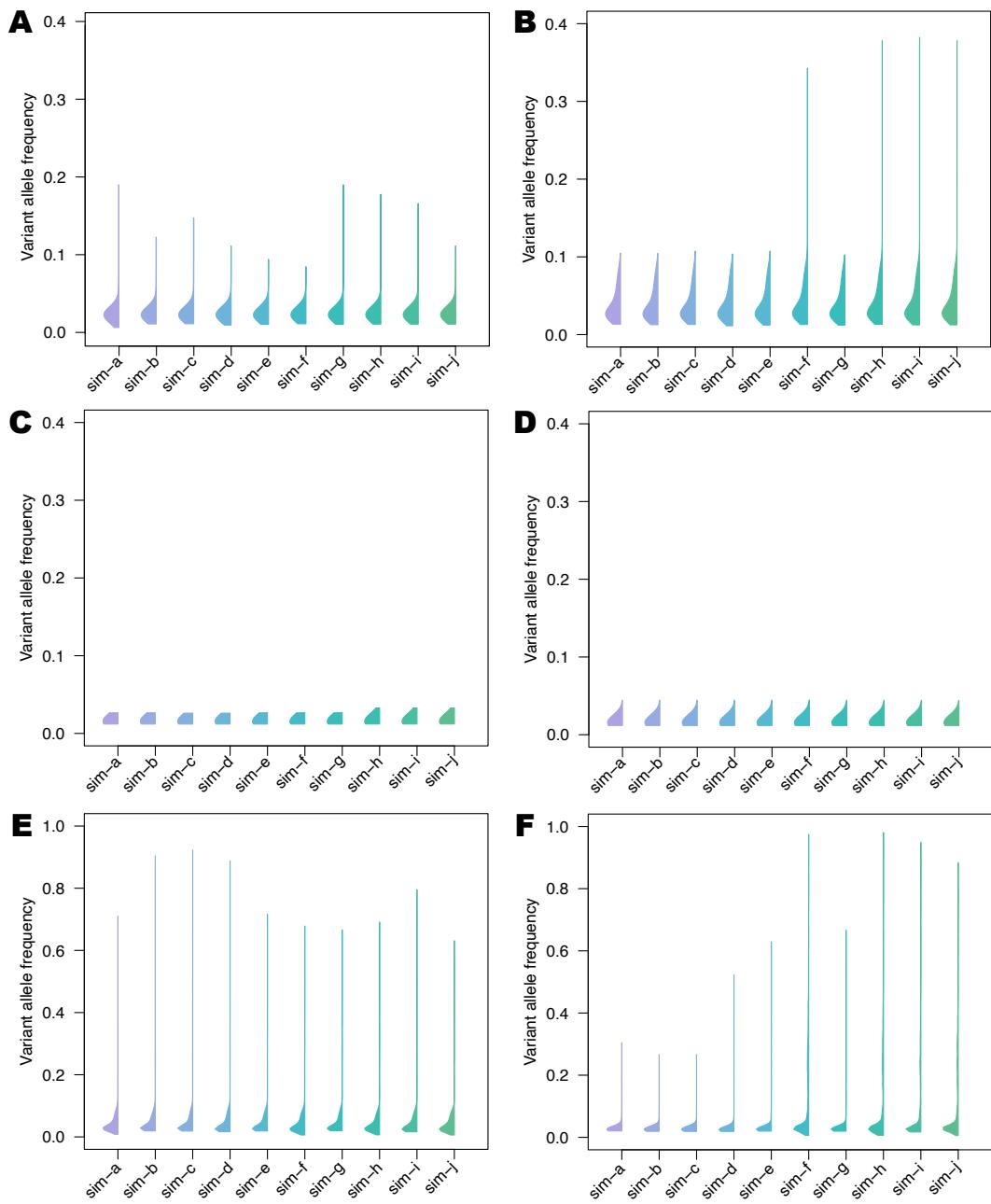


Figure A.4: Variant allele frequencies (VAF) of variants detected by joint sample analysis; A) VAF distribution of true positive variants additionally detected by Strelka2pass B) and FreeBayesSomatic C) VAF distribution of false positive variants additionally detected by FreeBayesSomatic D) and Strelka2pass E) VAF distribution of false negatives not called by FreeBayesSomatic F) and Strelka2pass.

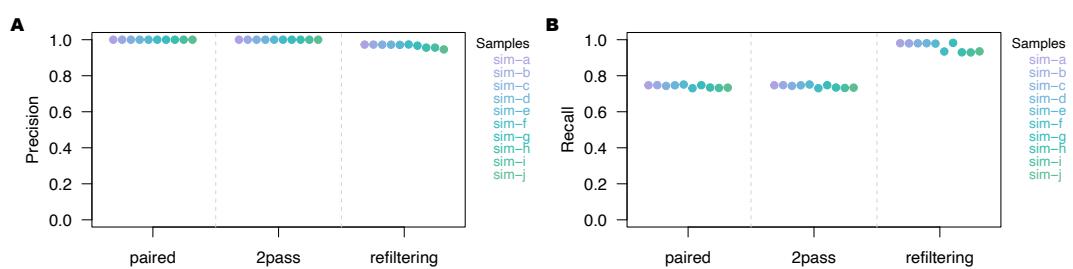


Figure A.5: Performance of individual steps in the Strelka2pass workflow using the simulated data:
A) Precision and B) Recall of tumour-normal paired analysis, two-pass step without refiltering (supplying variants from all tumour-normal pairs for evaluation) and two-pass step with refiltering (the final workflow)

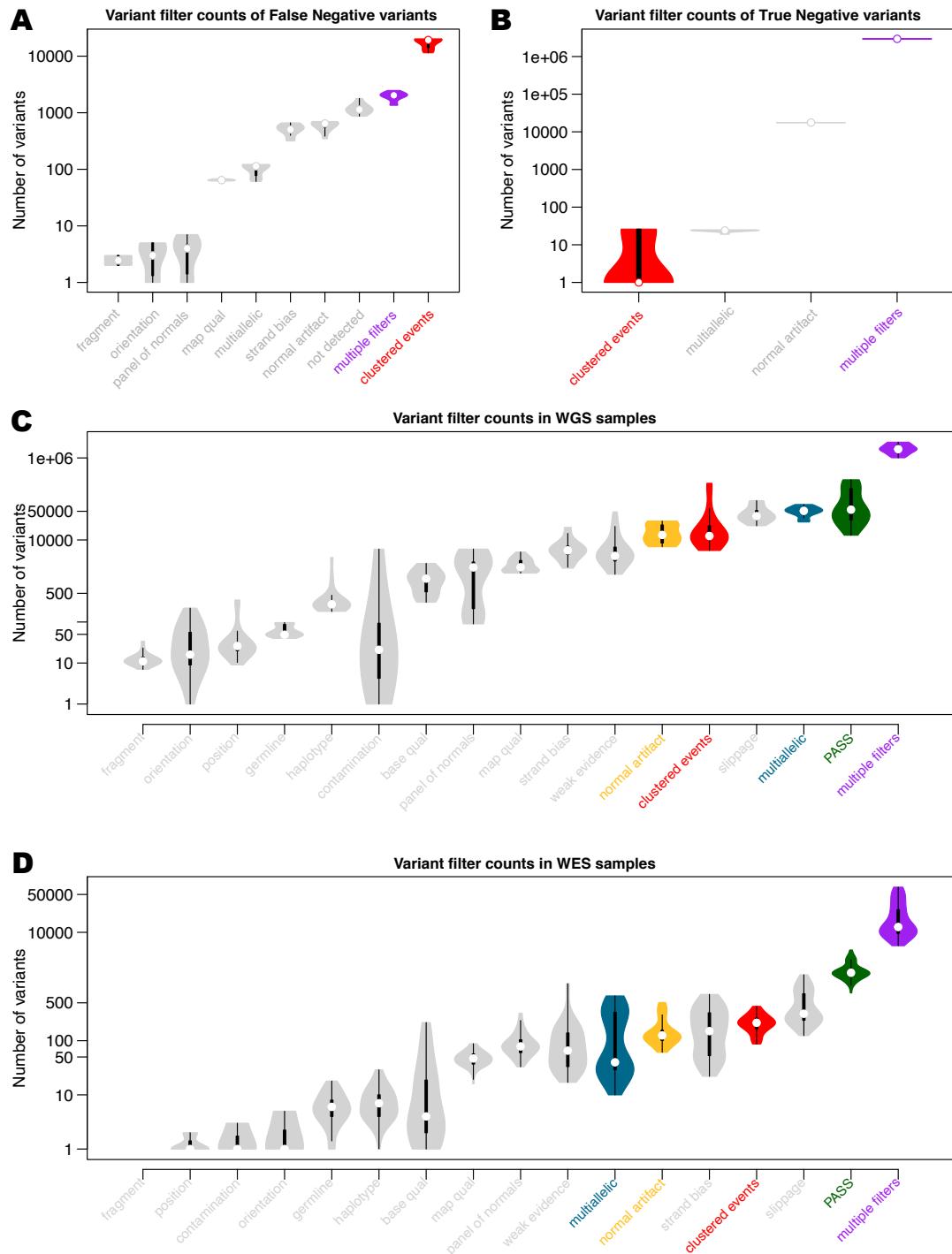


Figure A.6: Summary of variant filters assigned by Mutect2; The counts for each filter type are denoted by black boxplots with white circles depicting the median values. The fitted distribution of variant counts outlines each boxplot; A) Counts of filter assignments for false negative variants and B) true negative variants called by Mutect2 C) Filter assignment for all variants reported for sequenced patient data sequenced with WGS or D) WES.

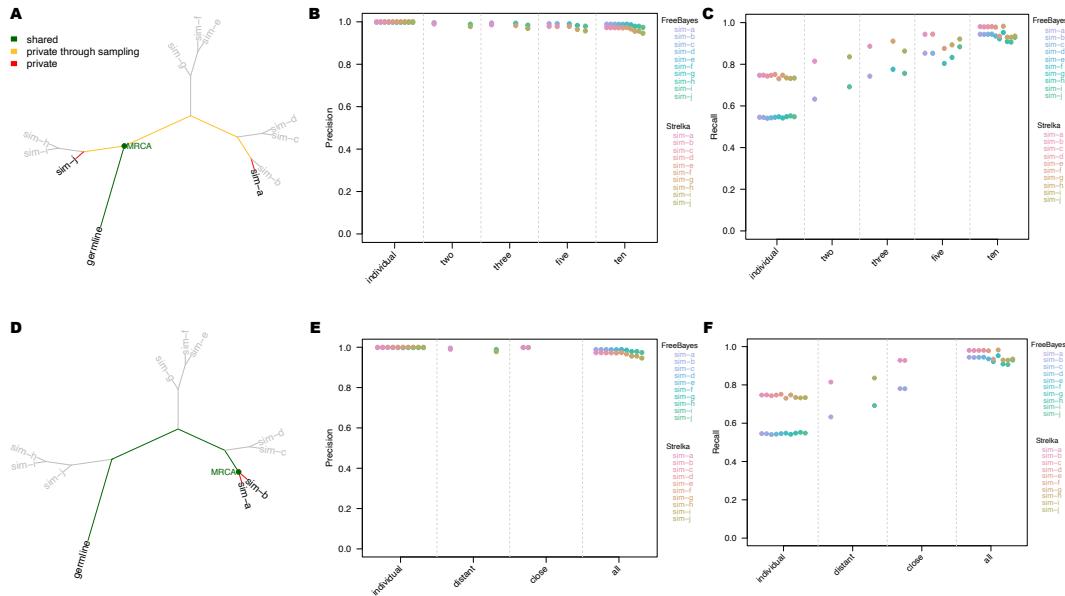


Figure A.7: Assessing the performance of different workflows using tumour samples with different evolutionary relationships in the simulated data; A) Simulated phylogeny highlighting two samples with high evolutionary distance (sim-a and sim-j) where MRCA denotes the most recent common ancestor. B) Precision and C) Recall estimates of FreeBayes and Strelka, run in individual tumour-normal paired and joint calling configurations using two (sim-a and sim-j), three (sim-a, sim-g and sim-j), five (sim-a, sim-c, sim-f, sim-h and sim-j) and all ten tumour samples D) Simulated phylogeny highlighting two samples with low evolutionary distance (sim-a and sim-b). E) Precision and F) Recall estimates for FreeBayes and Strelka run in individual tumour-normal paired and joint calling configurations. The plots compare the performance of these workflows when using two evolutionary distant samples (sim-a and sim-j), two evolutionary close samples (sim-a and sim-b) and all ten tumour samples.

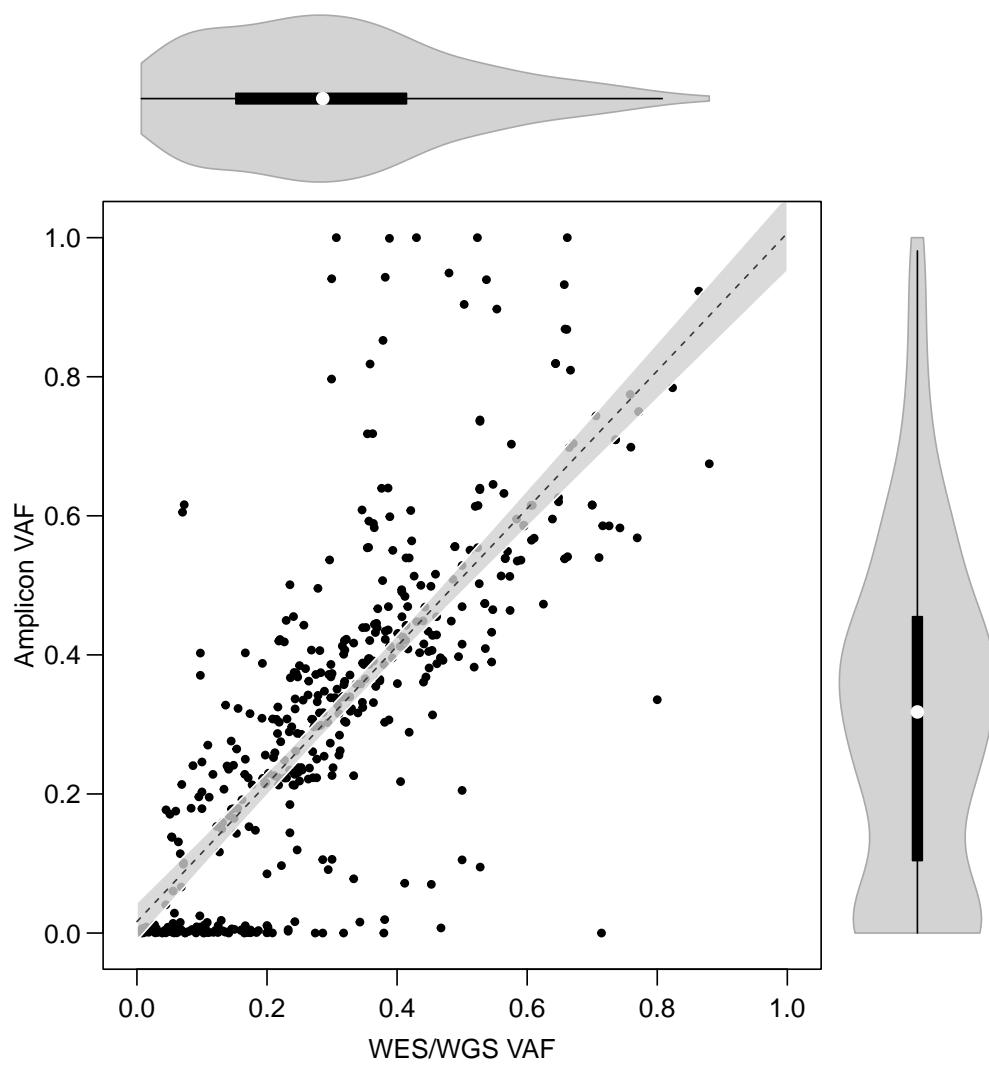


Figure A.8: Correlation of variant allele frequencies (VAF) from WES and WGS data against targeted amplicon sequencing VAF values with fitted violin plots of each individual distribution. Grey background shows 95% confidence interval for the fit of the linear model (dotted line).

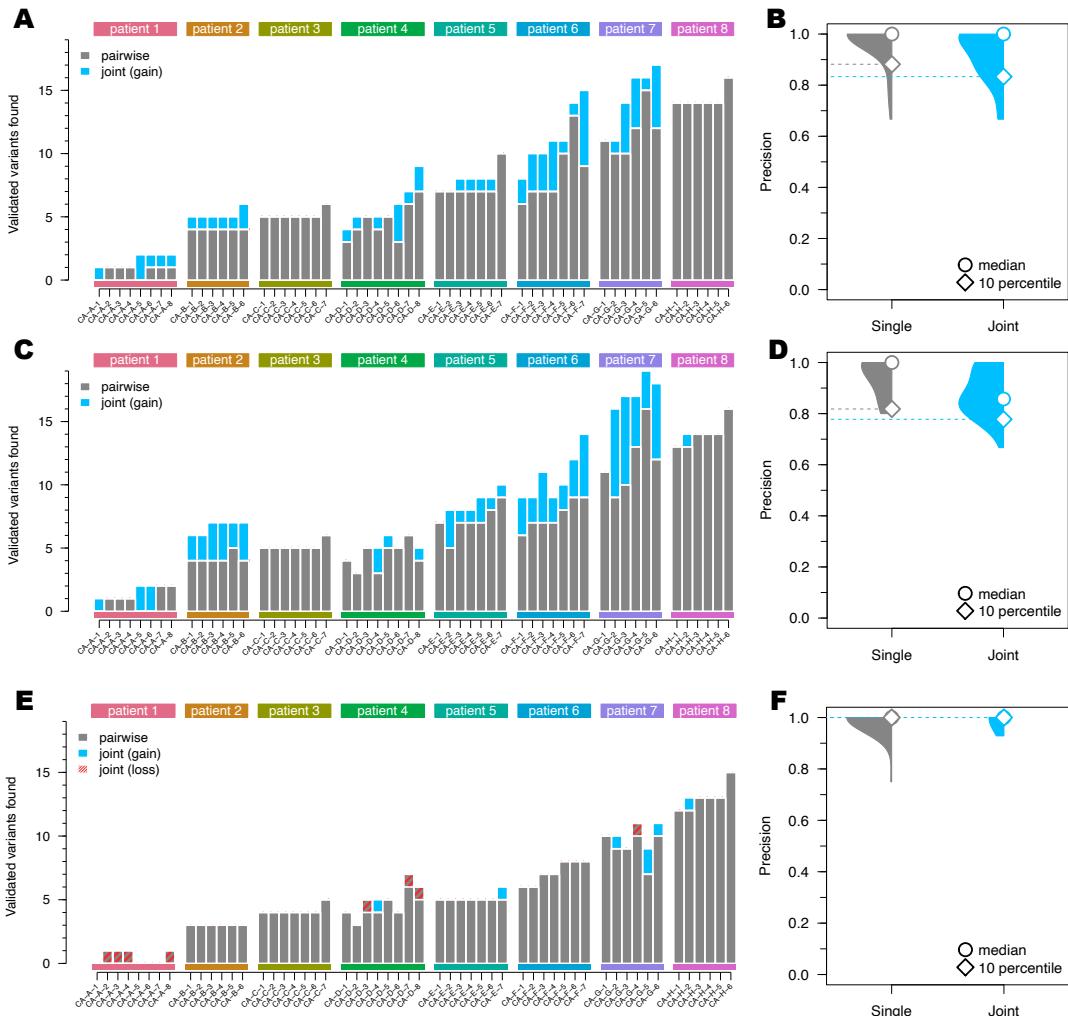


Figure A.9: Performance of the different workflows using clinical samples from eight cancer patients: A) Number of variants called by Strelka2 run in the tumour-normal paired (grey) and joint calling configurations, which have been validated by targeted amplicon sequencing (TAS). The same for C) FreeBayes and E) Mutect2 workflows. Precision of tumour-normal paired and joint analysis of TAS validated clinical data for B) Strelka2, D) FreeBayes and F) Mutect2; Sup. Table 1 provides the sample naming map to the original publications.

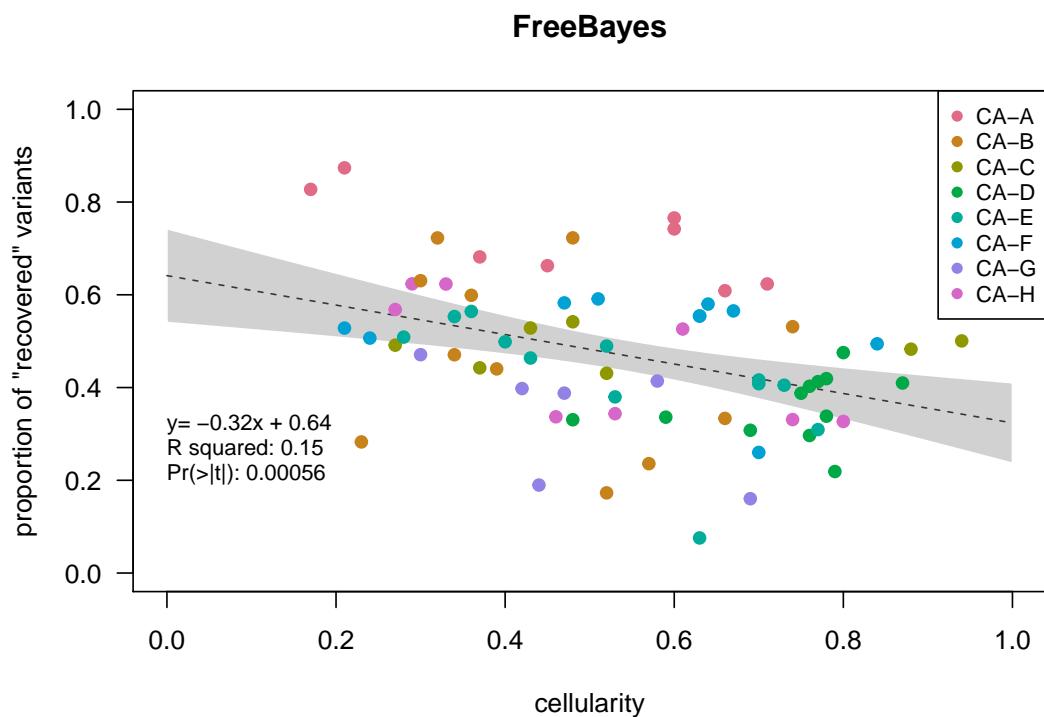


Figure A.10: Correlation between cellularity and proportion of variants found only with joint calling using FreeBayesSomatic. Grey background shows 95% confidence interval for fit of linear model (dotted line)

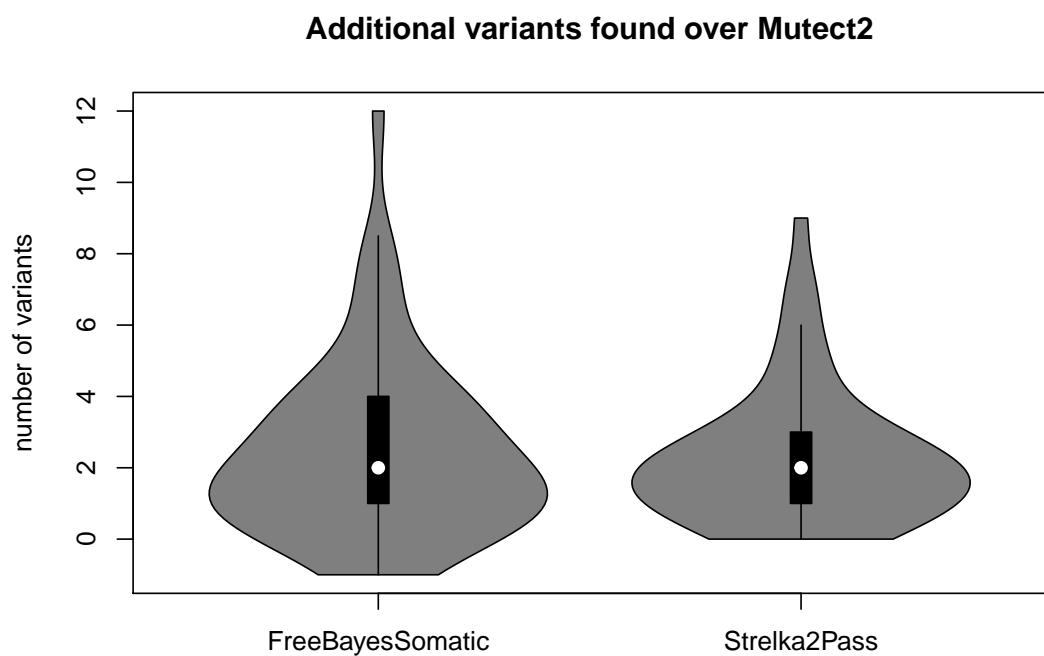


Figure A.11: Improvement in recall using FreeBayesSomatic and Strelka2pass over Mutect2 in the clinical samples.

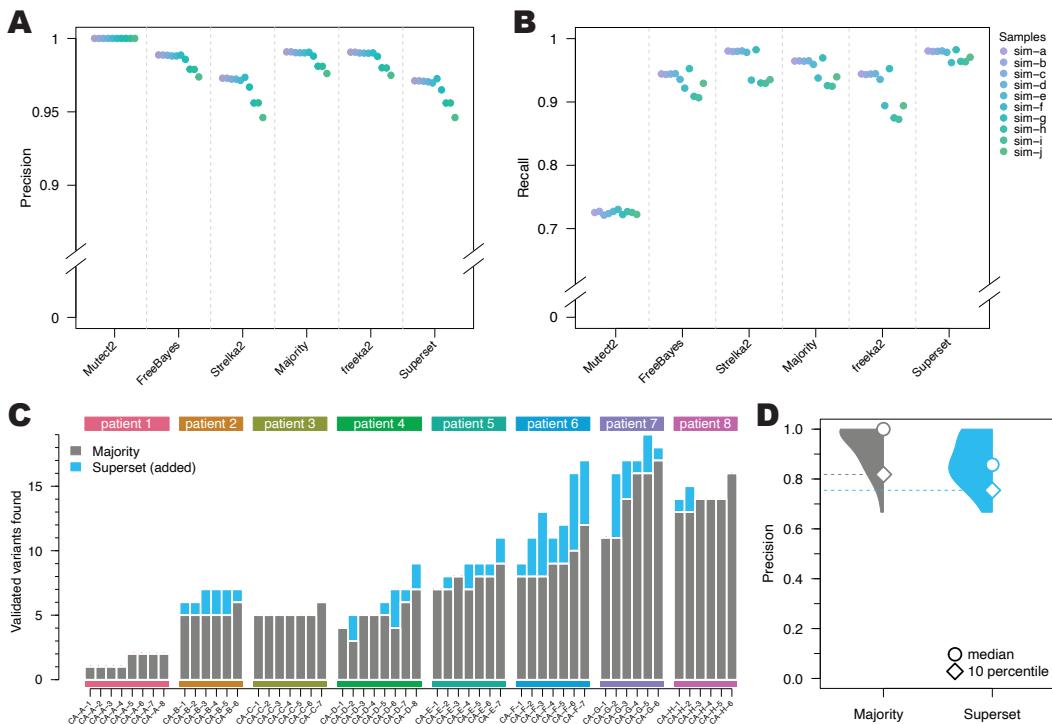


Figure A.12: Performance of ensemble variant calling strategies. A) Precision and B) Recall of variant detection using the joint multi-sample calling of each tool separately and compared to using Majority-vote ensemble calling (variant is called by at least two callers), Freek2 (variant is called by both FreeBayesSomatic and Strelka2pass) and Superset (variant is called by either FreeBayesSomatic or Strelka2pass) for the simulated dataset D) Number of TAS validated variants found in the clinical samples with Majority-vote and Superset methods and the corresponding D) Precision estimates.

Table A.1: Sample naming map relating to previously published datasets. The first column contains sample names as they appear in this work, and the third column denotes how the samples are referred to in the original studies. Forth column shows the type of sequencing WES: whole-exome sequencing; WGS: whole genome sequencing.

SAMPLE NAME	PUBLISHED STUDY	ORIGINAL NAME	SEQUENCING TYPE
CA-A-1	Solomon et al. [134]	Case 1 Left liver 1	WGS
CA-A-2		Case 1 Right occipital	
CA-A-3		Case 1 Right liver 2	
CA-A-4		Case 1 Right pleura	
CA-A-5		Case 1 Left lower lung lobe	
CA-A-6		Case 1 Left liver 5	
CA-A-7		Case 1 Right liver 3	
CA-A-8		Case 1 Left liver 2	
CA-B-1	Vergara et al. [135]	CAS-B-21-L-LUNG	WES
CA-B-2		CAS-B-22-R-LUNG	
CA-B-3		CAS-B-14B37035-1B	
CA-B-4		CAS-B-Primary-1	
CA-B-5		CAS-B-15Bo8317-3A	
CA-B-6		CAS-B-14B37035-1C	
CA-C-1		CAS-A-FR07935894	WGS
CA-C-2		CAS-A-FR07935905	
CA-C-3		CAS-A-FR07935906	
CA-C-4		CAS-A-FR07935907	
CA-C-5		CAS-A-FR07935908	
CA-C-6		CAS-A-FR07935916	
CA-C-7		CAS-A-FR07935918	
CA-D-1		CAS-G-91-2	WES
CA-D-2		CAS-G-75	
CA-D-3		CAS-G-74	
CA-D-4		CAS-G-71	
CA-D-5		CAS-G-91	
CA-D-6		CAS-G-76	WES
CA-D-7		CAS-G-94	
CA-D-8		CAS-G-72	
CA-E-1		CAS-D-70	
CA-E-2		CAS-D-61-3	
CA-E-3		CAS-D-66	
CA-E-4		CAS-D-68	WES
CA-E-5		CAS-D-64	
CA-E-6		CAS-D-61-2	
CA-E-7		CAS-D-62	
CA-F-1		CAS-C-41	WES
CA-F-2		CAS-C-40-Fresh	
CA-F-3		CAS-C-37	
CA-F-4		CAS-C-44	
CA-F-5		CAS-C-42-Fresh	
CA-F-6		CAS-C-43-Fresh	
CA-F-7		CAS-C-46-Primary	
CA-G-1		CAS-F-FR07935922	WGS
CA-G-2		CAS-F-FR07935915	
CA-G-3		CAS-F-FR07935913	
CA-G-4		CAS-F-FR07935909	
CA-G-5		CAS-F-FR07935904	
CA-G-6		CAS-F-FR07935903	
CA-H-1		CAS-E-1	WES
CA-H-2		CAS-E-3	
CA-H-3		CAS-E-4	
CA-H-4		CAS-E-10	
CA-H-5		CAS-E-6	
CA-H-6		CAS-E-8	

Table A.2: Runtime of different workflows on simulated data; The runtimes were generated on the Peter MacCallum Cancer Centre HPC cluster with Intel(R) Xeon(R) CPU E5-2660 v3 @ 2.60GHz. The times are displayed in single CPU runtime, but each workflow is highly parallelised, such that the user runtime is far lower.

Method	Number of tumour samples used for joint calling			
	2	3	5	10
FreeBayesSomatic	562h	811h	1185h	2292h
Strelka2Pass	310h	465h	776h	1552h
Mutect2	-	-	-	28418h

A.5 Supplementary methods

A.5.1 Alignment of clinical data

Detailed information on processing of the clinical sequencing datasets was published previously [134, 135]. Briefly, reads were aligned to GRCh38 for patient CAS-A and GRCh37 for patients CAS-B through CAS-H using BWA version 0.7.17 [166] allowing the use of alternative contigs. Reads were then marked as duplicates with Picard software (v2.17.3).

A.5.2 Validation of clinical data

Detailed information on targeted amplicon sequencing of patient samples can be found in the original publications [134, 135]. A SNV called in WES with any workflow was considered a true positive when the adjusted p-value calculated through an exact binomial test was lower than 0.05 on the TAS data. The probability of success for this test was estimated as the number of bases different from the reference divided by the total number of sequenced bases (0.001) and the number of trials was the read depth covering the variant. For indels, a variant was considered to be validated if either of the panel variant callers primal (in house) or canary [167] called the same variant.

Only amplicons with an average mapping rate of at least 80% over all samples, as well as an average coverage of more than 300 were considered for further analysis. WES variants were first subsetted to be within the area of the respective amplicons.

A.5.3 Purity estimation with sequenza

For CA-A the sequenza-utils python program was used to generate input files for the sequenza R program on the aligned BAM files [13]. Kmin and gamma were set to 100 and 500 respectively to discourage a highly fragmented result. For CA-B through -H the reported tumour purities were used from the publication [135].

A.5.4 Performance of individual steps in Strelka2Pass

As each of the three steps potentially has implications for the performance, we assessed the improvement provided by each step in the Strelka2pass workflow. [Figure A.5](#) shows, that there is no change in either precision or recall just by supplying variants from all tumour-normal pairs for a second round of evaluation. However, there is a >20% improvement in recall when coupling this to the refiltering step that we have built into the workflow.

A.5.5 Ensemble workflows – user suggestions

An overall workflow can contain any number of additional variant callers, when not restricted to callers with joint analysis capability. Importantly, there is no benefit of jointly analysing samples with Mutect2, and it may decrease the performance in some cases. Each of our presented workflows outperformed Mutect2 on the data shown here, so when assembling an ensemble method, these methods, should have a higher confidence assigned to them in joint analysis cases, than tumour-normal pair approaches.

Depending on the end needs of the user, an ensemble workflow can be optimised towards precision or recall. In [Figure A.12](#) we show the performance changes improvement that can be achieved by combining Mutect2 in tumour-normal paired analysis with the two new workflows FreeBayesSomatic and Strelka2Pass. First, in a “best of three” majority vote, where the variant needs to be called by two out of three variant callers, we enhance the precision of each of the individual tools, with slightly lower recall. On the other hand, with the super set approach, where any variant called in either FreeBayesSomatic or Strelka2Pass is included in the end result, this improves the recall even further, but slightly reduces the precision. This approach has the additional benefit of not needing to run Mutect2 which is an order of magnitude slower in our tests, than Strelka2Pass and FreeBayesSomatic ([Table A.2](#)). The usage of these workflows can be easily integrated into existing workflows and can be customised to the needs of the user.

Bibliography

- [1] Ibiayi Dagogo-Jack and Alice T. Shaw. “Tumour heterogeneity and resistance to cancer therapies”. In: *Nature Reviews Clinical Oncology* 15.2 (Nov. 2017), pp. 81–94. doi: [10.1038/nrclinonc.2017.166](https://doi.org/10.1038/nrclinonc.2017.166).
- [2] R. Fisher, L. Pusztai, and C. Swanton. “Cancer heterogeneity: implications for targeted therapeutics.” In: *British journal of cancer* 108 (3 Feb. 2013), pp. 479–485. issn: 1532-1827. doi: [10.1038/bjc.2012.581](https://doi.org/10.1038/bjc.2012.581). ppublish.
- [3] Tracy L. Leong et al. “Deep multi-region whole-genome sequencing reveals heterogeneity and gene-by-environment interactions in treatment-naive, metastatic lung cancer”. In: *Oncogene* 38.10 (Oct. 2018), pp. 1661–1675. doi: [10.1038/s41388-018-0536-1](https://doi.org/10.1038/s41388-018-0536-1).
- [4] Ting Yan et al. “Multi-region sequencing unveils novel actionable targets and spatial heterogeneity in esophageal squamous cell carcinoma”. In: *Nature Communications* 10.1 (Apr. 2019). doi: [10.1038/s41467-019-109255-1](https://doi.org/10.1038/s41467-019-109255-1).
- [5] The Centos Project. *Centos* 7. July 2014. url: <https://www.centos.org/> (visited on 10/26/2021).
- [6] Free Software Foundation. *Bash (3.2.48)*. [Unix shell program]. Version 5.1.8(1)-release. 2007. url: <http://ftp.gnu.org/gnu/bash/bash-3.2.48.tar.gz>.
- [7] Ole Tange et al. “GNU Parallel - The Command-Line Power Tool”. In: *login: The USENIX Magazine* 36.1 (Feb. 2011), pp. 42–47.
- [8] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria, 2021. url: <https://www.R-project.org/>.
- [9] Martin Morgan et al. *BiocParallel: Bioconductor facilities for parallel evaluation*. R package version 1.24.1. 2020. url: <https://github.com/Bioconductor/BiocParallel>.
- [10] Martin Morgan. *BiocManager: Access the Bioconductor Project Package Repository*. R package version 1.30.10. 2019. url: <https://CRAN.R-project.org/package=BiocManager>.
- [11] Achim Zeileis, Kurt Hornik, and Paul Murrell. “Escaping RGBland: Selecting colors for statistical graphics”. In: *Computational Statistics & Data Analysis* 53.9 (July 2009), pp. 3259–3270. doi: [10.1016/j.csda.2008.11.033](https://doi.org/10.1016/j.csda.2008.11.033).
- [12] Achim Zeileis et al. “colorspace: A Toolbox for Manipulating and Assessing Colors and Palettes”. In: *Journal of Statistical Software* 96.1 (2020). doi: [10.18637/jss.v096.i01](https://doi.org/10.18637/jss.v096.i01).
- [13] F. Favero et al. “Sequenza: allele-specific copy number and mutation profiles from tumor sequencing data”. In: *Annals of Oncology* 26.1 (Jan. 2015), pp. 64–70. doi: [10.1093/annonc/mdu479](https://doi.org/10.1093/annonc/mdu479).

- [14] Ronglai Shen and Venkatraman E. Seshan. “FACETS: allele-specific copy number and clonal heterogeneity analysis tool for high-throughput DNA sequencing”. In: *Nucleic Acids Research* 44.16 (June 2016), e131–e131. doi: [10.1093/nar/gkw520](https://doi.org/10.1093/nar/gkw520).
- [15] Venkatraman E. Seshan and Ronglai Shen. *facets: Cellular Fraction and Copy Numbers from Tumor Sequencing*. R package version 0.6.0. 2018. url: <https://github.com/mskcc/facets> (visited on 09/29/2021).
- [16] Daniel L. Cameron et al. “GRIDSS, PURPLE, LINX: Unscrambling the tumor genome via integrated analysis of structural variation and copy number”. In: *bioRxiv* (Sept. 2019). doi: [10.1101/781013](https://doi.org/10.1101/781013). url: <https://doi.org/10.1101/781013>.
- [17] Gro Nilsen et al. “Copynumber: Efficient algorithms for single- and multi-track copy number segmentation”. In: *BMC Genomics* 13.1 (Nov. 2012). doi: [10.1186/1471-2164-13-591](https://doi.org/10.1186/1471-2164-13-591).
- [18] Gro Nilsen, Knut Liestol, and Ole Christian Lingjaerde. *copynumber: Segmentation of single- and multi-track copy number data by penalized least squares regression*. R package version 1.29.0.9000. 2021.
- [19] William McLaren et al. “The Ensembl Variant Effect Predictor”. In: *Genome Biology* 17.1 (June 2016). doi: [10.1186/s13059-016-0974-4](https://doi.org/10.1186/s13059-016-0974-4).
- [20] Matt Dowle and Arun Srinivasan. *data.table: Extension of ‘data.frame’*. R package version 1.14.0. 2021. url: <https://CRAN.R-project.org/package=data.table>.
- [21] Daniel Adler and S. Thomas Kelly. *vioplot: violin plot*. R package version 0.3.5. 2020. url: <https://github.com/TomKellyGenetics/vioplot>.
- [22] Zuguang Gu, Roland Eils, and Matthias Schlesner. “Complex heatmaps reveal patterns and correlations in multidimensional genomic data”. In: *Bioinformatics* 32.18 (May 2016), pp. 2847–2849. doi: [10.1093/bioinformatics/btw313](https://doi.org/10.1093/bioinformatics/btw313).
- [23] Emmanuel Paradis and Klaus Schliep. “ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R”. In: *Bioinformatics* 35.3 (July 2018). Ed. by Russell Schwartz, pp. 526–528. doi: [10.1093/bioinformatics/bty633](https://doi.org/10.1093/bioinformatics/bty633).
- [24] Klaus Schliep et al. “Intertwining phylogenetic trees and networks”. In: *Methods in Ecology and Evolution* 8.10 (Apr. 2017). Ed. by Richard Fitzjohn, pp. 1212–1220. doi: [10.1111/2041-210X.12760](https://doi.org/10.1111/2041-210X.12760).
- [25] Tal Galili. “dendextend: an R package for visualizing, adjusting and comparing trees of hierarchical clustering”. In: *bioRxiv* 31.22 (July 2015), pp. 3718–3720. doi: [10.1101/2015.07.09.021428](https://doi.org/10.1101/2015.07.09.021428).

- [26] Jennifer Bryan. *googlesheets4: Access Google Sheets using the Sheets API V4*. R package version 1.0.0. 2021. url: <https://CRAN.R-project.org/package=googlesheets4>.
- [27] Martin Morgan et al. *Rsamtools: Binary alignment (BAM), FASTA, variant call (BCF), and tabix file import*. R package version 2.8.0. 2021. url: <https://bioconductor.org/packages/Rsamtools>.
- [28] Michael Lawrence et al. “Software for Computing and Annotating Genomic Ranges”. In: *PLoS Computational Biology* 9.8 (8 Aug. 2013). Ed. by Andreas Prlic, e1003118. doi: [10.1371/journal.pcbi.1003118](https://doi.org/10.1371/journal.pcbi.1003118). url: <http://www.ploscompbiol.org/article/info%3Adoi%2F10.1371%2Fjournal.pcbi.1003118>.
- [29] Trevor L Davis. *optparse: Command Line Option Parser*. R package version 1.6.6. 2020. url: <https://CRAN.R-project.org/package=optparse>.
- [30] Valerie Obenchain et al. “VariantAnnotation: a Bioconductor package for exploration and annotation of genetic variants”. In: *Bioinformatics* 30.14 (2014), pp. 2076–2078. doi: [10.1093/bioinformatics/btu168](https://doi.org/10.1093/bioinformatics/btu168).
- [31] Marcel Ramos et al. “Software for the integration of multi-omics experiments in Bioconductor”. In: *Cancer Research* 77(21); e39-42 (June 2017). doi: [10.1101/144774](https://doi.org/10.1101/144774).
- [32] Zuguang Gu et al. “circlize implements and enhances circular visualization in R”. In: *Bioinformatics* 30.19 (19 June 2014), pp. 2811–2812. doi: [10.1093/bioinformatics/btu393](https://doi.org/10.1093/bioinformatics/btu393).
- [33] Jitao David Zhang et al. “Detect tissue heterogeneity in gene expression data with BioQC”. In: *BMC Genomics* 18.1 (Apr. 2017), p. 277. doi: [10.1186/s12864-017-3661-2](https://doi.org/10.1186/s12864-017-3661-2). url: <http://accio.github.io/BioQC/>.
- [34] H. Pagès et al. *Biostrings: Efficient manipulation of biological strings*. R package version 2.58.0. 2020. url: <https://bioconductor.org/packages/Biostrings>.
- [35] Rachel Rosenthal. *deconstructSigs: Identifies Signatures Present in a Tumor Sample*. R package version 1.8.0. 2016. url: <https://CRAN.R-project.org/package=deconstructSigs>.
- [36] Hervé Pagès. *BSgenome: Software infrastructure for efficient representation of full genomes and their SNPs*. R package version 1.58.0. 2020. url: <https://bioconductor.org/packages/BSgenome>.
- [37] Ilari Scheinin et al. “DNA copy number analysis of fresh and formalin-fixed specimens by shallow whole-genome sequencing with identification and exclusion of problematic regions in the genome assembly”. In: *Genome Research* 24.12 (Sept. 2014), pp. 2022–2032. doi: [10.1101/gr.175141.114](https://doi.org/10.1101/gr.175141.114).

- [38] Erich Neuwirth. *RColorBrewer: ColorBrewer Palettes*. R package version 1.1-2. 2014. url: <https://CRAN.R-project.org/package=RColorBrewer>.
- [39] Raivo Kolde. *pheatmap: Pretty Heatmaps*. R package version 1.0.12. 2019. url: <https://CRAN.R-project.org/package=pheatmap>.
- [40] Valerie Obenchain and Lori Shepherd. *ensemblVEP: R Interface to Ensembl Variant Effect Predictor*. R package version 1.32.0. 2020.
- [41] M.P.J. van der Loo. “The stringdist Package for Approximate String Matching”. In: *The R Journal* 6.1 (1 2014), pp. 111–122. doi: [10.32614/rj-2014-011](https://doi.org/10.32614/rj-2014-011). url: <https://CRAN.R-project.org/package=stringdist>.
- [42] Yang Liao, Gordon K. Smyth, and Wei Shi. “The R package Rsubread is easier, faster, cheaper and better for alignment and quantification of RNA sequencing reads”. In: *Nucleic Acids Research* 47.8 (8 Feb. 2019), e47–e47. doi: [10.1093/nar/gkz114](https://doi.org/10.1093/nar/gkz114).
- [43] Hadley Wickham et al. *svglite: An 'SVG' Graphics Device*. R package version 2.0.0. 2021. url: <https://CRAN.R-project.org/package=svglite>.
- [44] Paul Murrell. “Importing Vector Graphics: The grImport Package for R”. In: *Journal of Statistical Software* 30.4 (2009), pp. 1–37. doi: [10.18637/jss.v030.i04](https://doi.org/10.18637/jss.v030.i04). url: <http://www.jstatsoft.org/v30/i04/>.
- [45] Duncan Temple Lang. *XML: Tools for Parsing and Generating XML Within R and S-Plus*. R package version 3.99-0.5. 2020. url: <https://CRAN.R-project.org/package=XML>.
- [46] Hao Zhu. *kableExtra: Construct Complex Table with 'kable' and Pipe Syntax*. R package version 1.3.4. 2021. url: <https://CRAN.R-project.org/package=kableExtra>.
- [47] Fridolin Wild. *lsa: Latent Semantic Analysis*. R package version 0.73.2. 2020. url: <https://CRAN.R-project.org/package=lsa>.
- [48] Jim Baglama, Lothar Reichel, and B. W. Lewis. *irlba: Fast Truncated Singular Value Decomposition and Principal Components Analysis for Large Dense and Sparse Matrices*. R package version 2.3.3. 2019. url: <https://CRAN.R-project.org/package=irlba>.
- [49] Hadley Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, June 8, 2016. 260 pp. isbn: 978-3-319-24277-4. url: <https://ggplot2.tidyverse.org>.
- [50] Guido VanRossum. *The Python language reference*. Hampton, NHRedwood City, Calif: Python Software FoundationSoHo Books, 2010. isbn: 9781441412690.

- [51] Charles R. Harris et al. “Array programming with NumPy”. In: *Nature* 585.7825 (Sept. 2020), pp. 357–362. doi: [10.1038/s41586-020-2649-2](https://doi.org/10.1038/s41586-020-2649-2).
- [52] Endre Bakken Stovner and Pål Sætrom. “PyRanges: efficient comparison of genomic intervals in Python”. In: *Bioinformatics* (Aug. 2019). Ed. by John Hancock. doi: [10.1093/bioinformatics/btz615](https://doi.org/10.1093/bioinformatics/btz615).
- [53] Andreas Heger, Kevin Jacobs, et al. *pysam: htslib interface for python*. Oct. 25, 2021. url: <https://github.com/pysam-developers/pysam> (visited on 10/26/2021).
- [54] James K Bonfield et al. “HTSlib: C library for reading/writing high-throughput sequencing data”. In: *GigaScience* 10.2 (Jan. 2021). doi: [10.1093/gigascience/giaboo7](https://doi.org/10.1093/gigascience/giaboo7).
- [55] Petr Danecek et al. “Twelve years of SAMtools and BCFtools”. In: *GigaScience* 10.2 (Jan. 2021). doi: [10.1093/gigascience/giaboo8](https://doi.org/10.1093/gigascience/giaboo8).
- [56] Alistair Miles et al. *zarr-developers/zarr-python: v2.10.2*. 2021. doi: [10.5281/ZENODO.5579625](https://doi.org/10.5281/ZENODO.5579625).
- [57] Wes McKinney et al. “Data Structures for Statistical Computing in Python”. In: *Proceedings of the 9th Python in Science Conference*. Vol. 445. Austin, TX. SciPy, 2010, pp. 51–56. doi: [10.25080/majora-92bf1922-00a](https://doi.org/10.25080/majora-92bf1922-00a).
- [58] Jeff Reback et al. *pandas-dev/pandas: Pandas 1.3.4*. 2021. doi: [10.5281/ZENODO.5574486](https://doi.org/10.5281/ZENODO.5574486).
- [59] Robert T. McGibbon et al. *quadprog: Quadratic Programming Solver (Python)* v0.1.10. Oct. 1, 2021. url: <https://github.com/quadprog/quadprog> (visited on 10/26/2021).
- [60] Pauli Virtanen et al. “SciPy 1.0: fundamental algorithms for scientific computing in Python”. In: *Nature Methods* 17.3 (Feb. 2020), pp. 261–272. doi: [10.1038/s41592-019-0686-2](https://doi.org/10.1038/s41592-019-0686-2).
- [61] J. D. Watson and F. H. C. Crick. “Molecular Structure of Nucleic Acids: A Structure for Deoxyribose Nucleic Acid”. In: *Nature* 171.4356 (Apr. 1953), pp. 737–738. doi: [10.1038/171737a0](https://doi.org/10.1038/171737a0).
- [62] F. Liang et al. “Homology-directed repair is a major double-strand break repair pathway in mammalian cells”. In: *Proceedings of the National Academy of Sciences* 95.9 (Apr. 1998), pp. 5172–5177. doi: [10.1073/pnas.95.9.5172](https://doi.org/10.1073/pnas.95.9.5172).
- [63] Richard R. Sinden. “Introduction to the Structure, Properties, and Reactions of DNA”. In: *DNA Structure and Function*. Elsevier, 1994, pp. 1–57. doi: [10.1016/b978-0-08-057173-7.50006-7](https://doi.org/10.1016/b978-0-08-057173-7.50006-7).
- [64] J. Craig Venter et al. “The Sequence of the Human Genome”. In: *Science* 291.5507 (Feb. 2001), pp. 1304–1351. doi: [10.1126/science.1058040](https://doi.org/10.1126/science.1058040).

- [65] Colin M. Hammond et al. “Histone chaperone networks shaping chromatin function”. In: *Nature Reviews Molecular Cell Biology* 18.3 (Jan. 2017), pp. 141–158. doi: [10.1038/nrm.2016.159](https://doi.org/10.1038/nrm.2016.159).
- [66] Boyan Bonev and Giacomo Cavalli. “Organization and function of the 3D genome”. In: *Nature Reviews Genetics* 17.11 (Oct. 2016), pp. 661–678. doi: [10.1038/nrg.2016.112](https://doi.org/10.1038/nrg.2016.112).
- [67] Tuguo Tateoka. “A contribution to the taxonomy of the Agrostis mertensii-flaccida complex (Poaceae) in Japan”. In: *The Botanical Magazine Tokyo* 88.2 (June 1975), pp. 65–87. doi: [10.1007/bf02491243](https://doi.org/10.1007/bf02491243).
- [68] R. Trivers and H Hare. “Haplodiploidy and the evolution of the social insect”. In: *Science* 191.4224 (Jan. 1976), pp. 249–263. doi: [10.1126/science.1108197](https://doi.org/10.1126/science.1108197).
- [69] Sarah P. Otto. “The Evolutionary Consequences of Polyploidy”. In: *Cell* 131.3 (Nov. 2007), pp. 452–462. doi: [10.1016/j.cell.2007.10.022](https://doi.org/10.1016/j.cell.2007.10.022).
- [70] Marvin I. Gottlieb et al. “Trisomy-17 syndrome”. In: *The American Journal of Medicine* 33.5 (Nov. 1962), pp. 763–773. doi: [10.1016/0002-9343\(62\)90253-x](https://doi.org/10.1016/0002-9343(62)90253-x).
- [71] Anna Cereda and John C Carey. “Trisomy 18 Syndrome”. In: *Atlas of Genetic Diagnosis and Counseling*. Vol. 7. 1. Humana Press, 2012, pp. 990–996. doi: [10.1186/1750-1172-7-81](https://doi.org/10.1186/1750-1172-7-81).
- [72] Maj A Hultén et al. “On the origin of trisomy 21 Down syndrome”. In: *Molecular Cytogenetics* 1.1 (2008), p. 21. doi: [10.1186/1755-8166-1-21](https://doi.org/10.1186/1755-8166-1-21).
- [73] David M. J. Lilley. “Structures of helical junctions in nucleic acids”. In: *Quarterly Reviews of Biophysics* 33.2 (May 2000), pp. 109–159. doi: [10.1017/s0033583500003590](https://doi.org/10.1017/s0033583500003590).
- [74] William P. Hanage, Christophe Fraser, and Brian G. Spratt. “The impact of homologous recombination on the generation of diversity in bacteria”. In: *Journal of Theoretical Biology* 239.2 (Mar. 2006), pp. 210–219. doi: [10.1016/j.jtbi.2005.08.035](https://doi.org/10.1016/j.jtbi.2005.08.035).
- [75] Ying Kong et al. “Homologous Recombination Drives Both Sequence Diversity and Gene Content Variation in *Neisseria meningitidis*”. In: *Genome Biology and Evolution* 5.9 (July 2013), pp. 1611–1627. doi: [10.1093/gbe/evt116](https://doi.org/10.1093/gbe/evt116).
- [76] Charles Darwin. “On the Origin of Species by Means of Natural Selection, or the Preservation of Favoured Races in the Struggle for Life”. In: *Evolutionary Writings*. Oxford University Press, May 2010. doi: [10.1093/owc/9780199580149.003.0005](https://doi.org/10.1093/owc/9780199580149.003.0005).
- [77] Kathleen Sprouffske et al. “High mutation rates limit evolutionary adaptation in *Escherichia coli*”. In: *PLOS Genetics* 14.4 (Apr. 2018). Ed. by Ivan Matic, e1007324. doi: [10.1371/journal.pgen.1007324](https://doi.org/10.1371/journal.pgen.1007324).

- [78] Ludmil B Alexandrov et al. “Clock-like mutational processes in human somatic cells”. In: *Nature Genetics* 47.12 (Nov. 2015), pp. 1402–1407. doi: [10.1038/ng.3441](https://doi.org/10.1038/ng.3441).
- [79] Luiza Moore et al. “The mutational landscape of human somatic and germline cells”. In: *Nature* (Aug. 2021). doi: [10.1038/s41586-021-03822-7](https://doi.org/10.1038/s41586-021-03822-7).
- [80] Hanan E. Shamseldin et al. “Identification of embryonic lethal genes in humans by autozygosity mapping and exome sequencing in consanguineous families”. In: *Genome Biology* 16.1 (June 2015). doi: [10.1186/s13059-015-0681-6](https://doi.org/10.1186/s13059-015-0681-6).
- [81] Laura Frey et al. “Mammalian VPS45 orchestrates trafficking through the endosomal system”. In: *Blood* 137.14 (Apr. 2021), pp. 1932–1944. doi: [10.1182/blood.2020006871](https://doi.org/10.1182/blood.2020006871).
- [82] Shan Dan et al. “Clinical application of massively parallel sequencing-based prenatal non-invasive fetal trisomy test for trisomies 21 and 18 in 11 105 pregnancies with mixed risk factors”. In: *Prenatal Diagnosis* 32.13 (Nov. 2012), pp. 1225–1232. doi: [10.1002/pd.4002](https://doi.org/10.1002/pd.4002).
- [83] Kypros H. Nicolaides et al. “Noninvasive Prenatal Testing for Fetal Trisomies in a Routinely Screened First-Trimester Population”. In: *Obstetrical & Gynecological Survey* 68.3 (Mar. 2013), pp. 173–175. doi: [10.1097/ogx.0b013e318285bf66](https://doi.org/10.1097/ogx.0b013e318285bf66).
- [84] Frank Diehl et al. “Circulating mutant DNA to assess tumor dynamics”. In: *Nature Medicine* 14.9 (July 2008), pp. 985–990. doi: [10.1038/nm.1789](https://doi.org/10.1038/nm.1789).
- [85] Heidi Schwarzenbach, Dave S. B. Hoon, and Klaus Pantel. “Cell-free nucleic acids as biomarkers in cancer patients”. In: *Nature Reviews Cancer* 11.6 (May 2011), pp. 426–437. doi: [10.1038/nrc3066](https://doi.org/10.1038/nrc3066).
- [86] R. Padmanabhan, E. Jay, and R. Wu. “Chemical Synthesis of a Primer and Its Use in the Sequence Analysis of the Lysozyme Gene of Bacteriophage T4”. In: *Proceedings of the National Academy of Sciences* 71.6 (June 1974), pp. 2510–2514. doi: [10.1073/pnas.71.6.2510](https://doi.org/10.1073/pnas.71.6.2510).
- [87] F. Sanger and A.R. Coulson. “A rapid method for determining sequences in DNA by primed synthesis with DNA polymerase”. In: *Journal of Molecular Biology* 94.3 (May 1975), pp. 441–448. doi: [10.1016/0022-2836\(75\)90213-2](https://doi.org/10.1016/0022-2836(75)90213-2).
- [88] F. Sanger, S. Nicklen, and A. R. Coulson. “DNA sequencing with chain-terminating inhibitors”. In: *Proceedings of the National Academy of Sciences* 74.12 (Dec. 1977), pp. 5463–5467. doi: [10.1073/pnas.74.12.5463](https://doi.org/10.1073/pnas.74.12.5463).
- [89] Eric S. Lander et al. “Initial sequencing and analysis of the human genome”. In: *Nature* 409.6822 (Feb. 2001), pp. 860–921. doi: [10.1038/35057062](https://doi.org/10.1038/35057062).

- [90] Inc Illumina. *How short inserts affect sequencing performance*. Sept. 2020. url: <https://sapac.support.illumina.com/bulletins/2020/12/how-short-inserts-affect-sequencing-performance.html> (visited on 09/08/2021).
- [91] Elaine R. Mardis. “Next-Generation DNA Sequencing Methods”. In: *Annual Review of Genomics and Human Genetics* 9.1 (Sept. 2008), pp. 387–402. doi: [10.1146/annurev.genom.9.081307.164359](https://doi.org/10.1146/annurev.genom.9.081307.164359).
- [92] Jenny Straiton et al. “From Sanger sequencing to genome databases and beyond”. In: *BioTechniques* 66.2 (Feb. 2019), pp. 60–63. doi: [10.2144/btn-2019-0011](https://doi.org/10.2144/btn-2019-0011).
- [93] G. M. Church and W. Gilbert. “Genomic sequencing.” In: *Proceedings of the National Academy of Sciences* 81.7 (Apr. 1984), pp. 1991–1995. doi: [10.1073/pnas.81.7.1991](https://doi.org/10.1073/pnas.81.7.1991).
- [94] G. Church and S Kieffer-Higgins. “Multiplex DNA sequencing”. In: *Science* 240 (Apr. 1988), pp. 185–188. doi: [10.1126/science.3353714](https://doi.org/10.1126/science.3353714).
- [95] Alexander Payne et al. “BulkVis: a graphical viewer for Oxford nanopore bulk FAST5 files”. In: *Bioinformatics* 35.13 (Nov. 2018). Ed. by Inanc Birol, pp. 2193–2198. doi: [10.1093/bioinformatics/bty841](https://doi.org/10.1093/bioinformatics/bty841).
- [96] Ploy N. Pratanwanich et al. “Identification of differential RNA modifications from nanopore direct RNA sequencing with xPore”. In: *Nature Biotechnology* (July 2021). doi: [10.1038/s41587-021-00949-w](https://doi.org/10.1038/s41587-021-00949-w).
- [97] Nicolas L Bray et al. “Near-optimal probabilistic RNA-seq quantification”. In: *Nature Biotechnology* 34.5 (Apr. 2016), pp. 525–527. doi: [10.1038/nbt.3519](https://doi.org/10.1038/nbt.3519).
- [98] Rob Patro et al. “Salmon provides fast and bias-aware quantification of transcript expression”. In: *Nature Methods* 14.4 (Mar. 2017), pp. 417–419. doi: [10.1038/nmeth.4197](https://doi.org/10.1038/nmeth.4197).
- [99] Brian D. Ondov et al. “Mash: fast genome and metagenome distance estimation using MinHash”. In: *Genome Biology* 17.1 (June 2016). doi: [10.1186/s13059-016-0997-x](https://doi.org/10.1186/s13059-016-0997-x).
- [100] Brian B Luczak, Benjamin T James, and Hani Z Girgis. “A survey and evaluations of histogram-based statistics in alignment-free sequence comparison”. In: *Briefings in Bioinformatics* 20.4 (Dec. 2017), pp. 1222–1237. doi: [10.1093/bib/bbx161](https://doi.org/10.1093/bib/bbx161).
- [101] Heng Li. “Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM”. In: (Mar. 16, 2013). arXiv: [1303.3997 \[q-bio.GN\]](https://arxiv.org/abs/1303.3997).

- [102] Ben Langmead et al. “Scaling read aligners to hundreds of threads on general-purpose processors”. In: *Bioinformatics* 35.3 (July 2018). Ed. by John Hancock, pp. 421–432. doi: [10.1093/bioinformatics/bty648](https://doi.org/10.1093/bioinformatics/bty648).
- [103] Klaus F. X. Mayer et al. “A chromosome-based draft sequence of the hexaploid bread wheat (*Triticum aestivum*) genome”. In: *Science* 345.6194 (July 2014), pp. 1251788–1251788. doi: [10.1126/science.1251788](https://doi.org/10.1126/science.1251788).
- [104] Erik Garrison and Gabor Marth. “Haplotype-based variant detection from short-read sequencing”. In: *arXiv preprint arXiv:1207.3907 [q-bio.GN]* (July 17, 2012). arXiv: <http://arxiv.org/abs/1207.3907v2> [q-bio.GN].
- [105] Zhongwu Lai et al. “VarDict: a novel and versatile variant caller for next-generation sequencing in cancer research”. In: *Nucleic Acids Research* 44.11 (Apr. 2016), e108–e108. doi: [10.1093/nar/gkw227](https://doi.org/10.1093/nar/gkw227).
- [106] Sangtae Kim et al. “Strelka2: fast and accurate calling of germline and somatic variants”. In: *Nature Methods* 15.8 (July 2018), pp. 591–594. doi: [10.1038/s41592-018-0051-x](https://doi.org/10.1038/s41592-018-0051-x).
- [107] David Benjamin et al. “Calling Somatic SNVs and Indels with Mutect2”. In: *bioRxiv* (Dec. 2019). doi: [10.1101/861054](https://doi.org/10.1101/861054). url: <https://doi.org/10.1101/861054>.
- [108] Daniel P. Cooke, David C. Wedge, and Gerton Lunter. “A unified haplotype-based method for accurate and comprehensive variant calling”. In: *Nature Biotechnology* 39.7 (Mar. 2021), pp. 885–892. doi: [10.1038/s41587-021-00861-3](https://doi.org/10.1038/s41587-021-00861-3).
- [109] GATK Team. *Somatic calling is NOT simply a difference between two callsets*. Sept. 15, 2021. url: <https://gatk.broadinstitute.org/hc/en-us/articles/360035890491> (visited on 09/23/2021).
- [110] Amaro Taylor-Weiner et al. “DeTiN: overcoming tumor-in-normal contamination”. In: *Nature Methods* 15.7 (June 2018), pp. 531–534. doi: [10.1038/s41592-018-0036-9](https://doi.org/10.1038/s41592-018-0036-9).
- [111] Konrad J. Karczewski et al. “The mutational constraint spectrum quantified from variation in 141,456 humans”. In: *Nature* 581.7809 (May 2020), pp. 434–443. doi: [10.1038/s41586-020-2308-7](https://doi.org/10.1038/s41586-020-2308-7).
- [112] Ali Karimnezhad et al. “Accuracy and reproducibility of somatic point mutation calling in clinical-type targeted sequencing data”. In: *BMC Med Genomics* 13.1 (Oct. 2020). doi: [10.1186/s12920-020-00803-z](https://doi.org/10.1186/s12920-020-00803-z).
- [113] Douglas Hanahan and Robert A Weinberg. “The Hallmarks of Cancer”. In: *Cell* 100.1 (Jan. 2000), pp. 57–70. doi: [10.1016/s0092-8674\(00\)81683-9](https://doi.org/10.1016/s0092-8674(00)81683-9).

- [114] Douglas Hanahan and Robert A. Weinberg. "Hallmarks of Cancer: The Next Generation". In: *Cell* 144.5 (Mar. 2011), pp. 646–674. doi: [10.1016/j.cell.2011.02.013](https://doi.org/10.1016/j.cell.2011.02.013).
- [115] Yousef Ahmed Fouad and Carmen Aanei. "Revisiting the hallmarks of cancer." In: *American journal of cancer research* 7 (5 2017), pp. 1016–1036. issn: 2156-6976. epublish.
- [116] Rebecca L. Siegel, Kimberly D. Miller, and Ahmedin Jemal. "Cancer statistics, 2018". In: *CA: A Cancer Journal for Clinicians* 68.1 (Jan. 2018), pp. 7–30. doi: [10.3322/caac.21442](https://doi.org/10.3322/caac.21442).
- [117] Julian R. Molina et al. "Non-Small Cell Lung Cancer: Epidemiology, Risk Factors, Treatment, and Survivorship". In: *Mayo Clinic Proceedings* 83.5 (May 2008), pp. 584–594. doi: [10.4065/83.5.584](https://doi.org/10.4065/83.5.584).
- [118] Sophie Sun, Joan H. Schiller, and Adi F. Gazdar. "Lung cancer in never smokers — a different disease". In: *Nature Reviews Cancer* 7.10 (Oct. 2007), pp. 778–790. doi: [10.1038/nrc2190](https://doi.org/10.1038/nrc2190).
- [119] Ryan Poplin et al. "Scaling accurate genetic variant discovery to tens of thousands of samples". In: *bioRxiv* (Nov. 2017). doi: [10.1101/201178](https://doi.org/10.1101/201178).
- [120] Neil A. Miller et al. "A 26-hour system of highly sensitive whole genome sequencing for emergency management of genetic diseases". In: *Genome Medicine* 7.1 (Sept. 2015). doi: [10.1186/s13073-015-0221-8](https://doi.org/10.1186/s13073-015-0221-8).
- [121] Ryan Poplin et al. "A universal SNP and small-indel variant caller using deep neural networks". In: *Nature Biotechnology* 36.10 (Sept. 2018), pp. 983–987. doi: [10.1038/nbt.4235](https://doi.org/10.1038/nbt.4235).
- [122] Melanie Schirmer et al. "Illumina error profiles: resolving fine-scale variation in metagenomic sequencing data". In: *BMC Bioinformatics* 17.1 (Mar. 2016). doi: [10.1186/s12859-016-0976-y](https://doi.org/10.1186/s12859-016-0976-y).
- [123] Nicholas Stoler and Anton Nekrutenko. "Sequencing error profiles of Illumina sequencing instruments". In: *NAR Genomics and Bioinformatics* 3.1 (Jan. 2021). doi: [10.1093/nargab/lqab019](https://doi.org/10.1093/nargab/lqab019).
- [124] Geraldine van der Auwera Brian O'Connor. *Genomics in the Cloud*. O'Reilly UK Ltd., May 1, 2020. 467 pp. isbn: 1491975199. url: <https://www.oreilly.com/library/view/genomics-in-the/9781491975183/>.
- [125] GATK Team. *Panel of Normals (PON)*. July 23, 2021. url: <https://gatk.broadinstitute.org/hc/en-us/articles/360035890631> (visited on 09/23/2021).
- [126] GATK Team. *Mutect2 multi-sample*. Sept. 25, 2020. url: <https://gatk.broadinstitute.org/hc/en-us/community/posts/360062528691> (visited on 10/23/2020).

- [127] Malvina Josephidou, Andy G. Lynch, and Simon Tavaré. “multiSNV: a probabilistic approach for improving detection of somatic point mutations from multiple related tumour samples”. In: *Nucleic Acids Research* 43.9 (Feb. 2015), e61–e61. doi: [10.1093/nar/gkv135](https://doi.org/10.1093/nar/gkv135).
- [128] Christoffer Flensburg et al. “SuperFreq: Integrated mutation detection and clonal tracking in cancer”. In: *PLOS Computational Biology* 16.2 (Feb. 2020). Ed. by Florian Markowetz, e1007603. doi: [10.1371/journal.pcbi.1007603](https://doi.org/10.1371/journal.pcbi.1007603).
- [129] Colby Chiang et al. “SpeedSeq: ultra-fast personal genome analysis and interpretation”. In: *Nature Methods* 12.10 (Aug. 2015), pp. 966–968. doi: [10.1038/nmeth.3505](https://doi.org/10.1038/nmeth.3505).
- [130] Brad Chapman et al. *bcbio/bcbio-nextgen: v1.2.4*. 2021. doi: [10.5281/ZENODO.3564938](https://doi.org/10.5281/ZENODO.3564938).
- [131] Brendan A. Veeneman et al. “Two-pass alignment improves novel splice junction quantification”. In: *Bioinformatics* 32 (Oct. 2015), pp. 43–49. doi: [10.1093/bioinformatics/btv642](https://doi.org/10.1093/bioinformatics/btv642).
- [132] Weichun Huang et al. “ART: a next-generation sequencing read simulator”. In: *Bioinformatics* 28.4 (Dec. 2011), pp. 593–594. doi: [10.1093/bioinformatics/btr708](https://doi.org/10.1093/bioinformatics/btr708).
- [133] Adam D Ewing et al. “Combining tumor genome simulation with crowdsourcing to benchmark somatic single-nucleotide-variant detection”. In: *Nature Methods* 12.7 (May 2015), pp. 623–630. doi: [10.1038/nmeth.3407](https://doi.org/10.1038/nmeth.3407).
- [134] Benjamin J. Solomon et al. “RET Solvent Front Mutations Mediate Acquired Resistance to Selective RET Inhibition in RET-Driven Malignancies”. In: *Journal of Thoracic Oncology* 15.4 (Apr. 2020), pp. 541–549. doi: [10.1016/j.jtho.2020.01.006](https://doi.org/10.1016/j.jtho.2020.01.006).
- [135] Ismael A. Vergara et al. “Evolution of late-stage metastatic melanoma is dominated by aneuploidy and whole genome doubling”. In: *Nature Communications* 12.1 (Mar. 2021). doi: [10.1038/s41467-021-21576-8](https://doi.org/10.1038/s41467-021-21576-8).
- [136] Zheng Hu et al. “Quantitative evidence for early metastatic seeding in colorectal cancer”. In: *Nature Genetics* 51.7 (June 2019), pp. 1113–1122. doi: [10.1038/s41588-019-0423-x](https://doi.org/10.1038/s41588-019-0423-x).
- [137] N Saitou and M Nei. “The neighbor-joining method: a new method for reconstructing phylogenetic trees.” In: *Molecular Biology and Evolution* (July 1987). doi: [10.1093/oxfordjournals.molbev.ao40454](https://doi.org/10.1093/oxfordjournals.molbev.ao40454).
- [138] Radu Mihaescu, Dan Levy, and Lior Pachter. “Why Neighbor-Joining Works”. In: *Algorithmica* 54.1 (Dec. 2007), pp. 1–24. doi: [10.1007/s00453-007-9116-4](https://doi.org/10.1007/s00453-007-9116-4).

- [139] Robert Reuven Sokal and Charles Duncan Michener. *A Statistical Method for Evaluating Systematic Relationships*. Vol. 38.2. University of Kansas science bulletin 22. University of Kansas, 1958. 30 pp.
- [140] Emile Zuckerkandl and Linus Pauling. “Molecular Disease, Evolution, and Genic Heterogeneity”. In: *Horizons in biochemistry* (1962), pp. 189–225.
- [141] D. Shibata. “Mutation and epigenetic molecular clocks in cancer”. In: *Carcinogenesis* 32.2 (Nov. 2010), pp. 123–128. doi: [10.1093/carcin/bgq239](https://doi.org/10.1093/carcin/bgq239).
- [142] Joseph Felsenstein. “Evolutionary trees from DNA sequences: A maximum likelihood approach”. In: *Journal of Molecular Evolution* 17.6 (Nov. 1981), pp. 368–376. doi: [10.1007/bf01734359](https://doi.org/10.1007/bf01734359).
- [143] Masami Hasegawa, Hirohisa Kishino, and Taka-aki Yano. “Dating of the human-ape splitting by a molecular clock of mitochondrial DNA”. In: *Journal of Molecular Evolution* 22.2 (Oct. 1985), pp. 160–174. doi: [10.1007/bf02101694](https://doi.org/10.1007/bf02101694).
- [144] R. W. Hamming. “Error Detecting and Error Correcting Codes”. In: *Bell System Technical Journal* 29.2 (Apr. 1950), pp. 147–160. doi: [10.1002/j.1538-7305.1950.tb00463.x](https://doi.org/10.1002/j.1538-7305.1950.tb00463.x).
- [145] Benjamin Werner et al. “Measuring single cell divisions in human tissues from multi-region sequencing data”. In: *Nature Communications* 11.1 (Feb. 2020). doi: [10.1038/s41467-020-14844-6](https://doi.org/10.1038/s41467-020-14844-6).
- [146] T. Arai et al. “Tumor doubling time and prognosis in lung cancer patients: evaluation from chest films and clinical follow-up study”. In: *Japanese journal of clinical oncology* 24 (4 Aug. 1994), pp. 199–204. issn: 0368-2811. ppublish.
- [147] Damien M de Vienne. “Tanglegrams Are Misleading for Visual Evaluation of Tree Congruence”. In: 36.1 (Oct. 2018). Ed. by Jeffrey Townsend, pp. 174–176. doi: [10.1093/molbev/msy196](https://doi.org/10.1093/molbev/msy196).
- [148] L. Tan et al. “Prediction and monitoring of relapse in stage III melanoma using circulating tumor DNA”. In: *Annals of Oncology* 30.5 (May 2019), pp. 804–814. doi: [10.1093/annonc/mdz048](https://doi.org/10.1093/annonc/mdz048).
- [149] “ctDNA is a specific and sensitive biomarker in multiple human cancers.” In: *Cancer discovery* 4.4 (4 Apr. 2014), OF8. issn: 2159-8290. doi: [10.1158/2159-8290.CD-RW2014-051](https://doi.org/10.1158/2159-8290.CD-RW2014-051). ppublish.
- [150] Amit G Deshwar et al. “PhyloWGS: Reconstructing subclonal composition and evolution from whole-genome sequencing of tumors”. In: *Genome Biology* 16.1 (Feb. 2015). doi: [10.1186/s13059-015-0602-8](https://doi.org/10.1186/s13059-015-0602-8).

- [151] Yuchao Jiang et al. “Assessing intratumor heterogeneity and tracking longitudinal and spatial clonal evolutionary history by next-generation sequencing”. In: 113.37 (Aug. 2016), E5528–E5537. doi: [10.1073/pnas.1522203113](https://doi.org/10.1073/pnas.1522203113).
- [152] Francesco Marass et al. “A phylogenetic latent feature model for clonal deconvolution”. In: *The Annals of Applied Statistics* 10.4 (Dec. 2016). doi: [10.1214/16-aoas986](https://doi.org/10.1214/16-aoas986).
- [153] Sayaka Miura et al. “Predicting clone genotypes from tumor bulk sequencing of multiple samples”. In: (June 2018). Ed. by John Hancock. doi: [10.1093/bioinformatics/bty469](https://doi.org/10.1093/bioinformatics/bty469).
- [154] Mohammed El-Kebir, Gryte Satas, and Benjamin J. Raphael. “Inferring parsimonious migration histories for metastatic cancers”. In: 50.5 (Apr. 2018), pp. 718–726. doi: [10.1038/s41588-018-0106-z](https://doi.org/10.1038/s41588-018-0106-z).
- [155] Giulio Caravagna et al. “The MOBSTER R package for tumour subclonal deconvolution from bulk DNA whole-genome sequencing data”. In: *BMC Bioinformatics* 21.1 (Nov. 2020). doi: [10.1186/s12859-020-03863-1](https://doi.org/10.1186/s12859-020-03863-1).
- [156] Maxime Tarabichi et al. “A practical guide to cancer subclonal reconstruction from DNA sequencing.” In: *Nature methods* 18.2 (2 Feb. 2021), pp. 144–155. issn: 1548-7105. doi: [10.1038/s41592-020-01013-2](https://doi.org/10.1038/s41592-020-01013-2). ppublish.
- [157] Sayaka Miura et al. “Power and pitfalls of computational methods for inferring clone phylogenies and mutation orders from bulk sequencing data”. In: *Scientific Reports* 10.1 (Feb. 2020). doi: [10.1038/s41598-020-59006-2](https://doi.org/10.1038/s41598-020-59006-2).
- [158] Ignaty Leshchiner et al. “Comprehensive analysis of tumour initiation, spatial and temporal progression under multiple lines of treatment”. In: *bioRxiv* (Dec. 2018). doi: [10.1101/508127](https://doi.org/10.1101/508127).
- [159] Moritz Gerstung et al. “The evolutionary history of 2,658 cancers”. In: *Nature* 578.7793 (Feb. 2020), pp. 122–128. doi: [10.1038/s41586-019-1907-7](https://doi.org/10.1038/s41586-019-1907-7).
- [160] Marian L. Burr et al. “An Evolutionarily Conserved Function of Polycomb Silences the MHC Class I Antigen Presentation Pathway and Enables Immune Evasion in Cancer”. In: *Cancer Cell* 36.4 (Oct. 2019), 385–401.e8. doi: [10.1016/j.ccr.2019.08.008](https://doi.org/10.1016/j.ccr.2019.08.008).
- [161] Mark A DePristo et al. “A framework for variation discovery and genotyping using next-generation DNA sequencing data”. In: *Nature Genetics* 43.5 (Apr. 2011), pp. 491–498. doi: [10.1038/ng.806](https://doi.org/10.1038/ng.806).
- [162] Berke Ç Toptaş et al. “Comparing complex variants in family trios”. In: *Bioinformatics* (June 2018). Ed. by Oliver Stegle. doi: [10.1093/bioinformatics/bty443](https://doi.org/10.1093/bioinformatics/bty443).

- [163] Di Wang et al. “Multiregion Sequencing Reveals the Genetic Heterogeneity and Evolutionary History of Osteosarcoma and Matched Pulmonary Metastases”. In: *Cancer Research* 79.1 (Nov. 2018), pp. 7–20. doi: [10.1158/0008-5472.can-18-1086](https://doi.org/10.1158/0008-5472.can-18-1086).
- [164] Zixi Chen et al. “Systematic comparison of somatic variant calling performance among different sequencing depth and mutation frequency”. In: *Scientific Reports* 10.1 (Feb. 2020). doi: [10.1038/s41598-020-60559-5](https://doi.org/10.1038/s41598-020-60559-5).
- [165] Richard Lupat et al. *Janis: A Python framework for Portable Pipelines*. en. 2021. doi: [10.5281/ZENODO.4427231](https://doi.org/10.5281/ZENODO.4427231).
- [166] H. Li and R. Durbin. “Fast and accurate short read alignment with Burrows-Wheeler transform”. In: *Bioinformatics* 25.14 (May 2009), pp. 1754–1760. doi: [10.1093/bioinformatics/btp324](https://doi.org/10.1093/bioinformatics/btp324).
- [167] Kenneth D. Doig et al. “Canary: an atomic pipeline for clinical amplicon assays”. In: *BMC Bioinformatics* 18.1 (Dec. 2017). doi: [10.1186/s12859-017-1950-z](https://doi.org/10.1186/s12859-017-1950-z).