

1 **Clonal phylogenies inferred from bulk, single**
2 **cell, and spatial transcriptomic analysis of**
3 **cancer**

4 **Authors: Andrew Erickson¹, Sandy Figiel¹, Timothy Rajakumar¹, Srinivasa Rao¹,**
5 **Wencheng Yin¹, Dimitrios Doultinos¹, Anette Magnussen¹, Reema Singh¹, Ninu**
6 **Poulose¹, Richard J Bryant^{1, 2}, Olivier Cussenot¹, Freddie C Hamdy^{1, 2}, Dan Woodcock¹,**
7 **Ian G Mills¹, Alastair D Lamb^{1, 2}**
8

9 **¹ Nuffield Department of Surgical Sciences, University of Oxford, Oxford, United**
10 **Kingdom, ² Department of Urology, Oxford University Hospitals NHS Foundation**
11 **Trust, Oxford, UK**

12 **Corresponding Author: Dr Alastair D Lamb, alastair.lamb@nds.ox.ac.uk**
13 **CRUK Clinician Scientist Fellow & Honorary Consultant Urological Surgeon**
14 **Nuffield Department of Surgical Sciences**
15 **University of Oxford**
16 **Old Road Campus Research Building**
17 **Oxford OX3 7DQ**
18 **alastair.lamb@nds.ox.ac.uk**
19

20 **Running title: Using transcript data for clonal tumor phylogenies**
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

1 **Abstract**

2

3 Epithelial cancers are typically heterogeneous with primary prostate cancer being a typical

4 example of histological and genomic variation. Prostate cancer is the second most common

5 male cancer in western industrialized countries. Prior studies of primary prostate cancer tumor

6 genetics revealed extensive inter and intra-patient tumor heterogeneity. Recent advances

7 have enabled extensive single-cell and spatial transcriptomic profiling of tissue specimens.

8 The ability to resolve accurate prostate cancer tumor phylogenies at high spatial resolution

9 would provide tools to address questions in tumorigenesis, disease progression, and

10 metastasis. Recent advances in machine learning have enabled the inference of ground-truth

11 genomic single-nucleotide and copy number variant status from transcript data. The inferred

12 SNV and CNV states can be used to resolve clonal phylogenies, however, it is still unknown

13 how faithfully transcript-based tumor phylogenies reconstruct ground truth DNA-based tumor

14 phylogenies. We sought to study the accuracy of inferred-transcript to recapitulate DNA-based

15 tumor phylogenies.

16

17 We first performed in-silico comparisons of inferred and directly resolved SNV and CNV status,

18 from single cancer cells, from three different cell lines. We found that inferred SNV phylogenies

19 accurately recapitulate DNA phylogenies (entanglement = 0.097). We observed similar results

20 in iCNV and CNV based phylogenies (entanglement = 0.11). Analysis of published prostate

21 cancer DNA phylogenies and inferred CNV, SNV and transcript based phylogenies

22 demonstrated phylogenetic concordance. Finally, a comparison of pseudo-bulked spatial

23 transcriptomic data to adjacent sections with WGS data also demonstrated recapitulation of

24 ground truth (entanglement = 0.35). These results suggest that transcript-based inferred

25 phylogenies recapitulate conventional genomic phylogenies. Further work will need to be done

26 to increase accuracy, genomic, and spatial resolution.

27

28

1 Introduction

2 It is generally accepted that cancers develop and evolve by adaptive genetic and
3 molecular changes over time (Nowell 1976; Greaves and Maley 2012; Black and McGranahan
4 2021). Sequential selection from this process of evolution leads to clones and subclones with
5 altered phenotype leading to more aggressive behaviour. Ultimately, these phenotypic
6 changes lead to metastatic spread and drug resistance, which is responsible for the majority
7 of cancer-related deaths (Gupta and Massagué 2006).

8 It is necessary to distinguish accurately tumour heterogeneity and determine clonal
9 evolution by identifying the clonal source of metastatic disease. This not only has an impact
10 on the understanding of tumour progression but the relationship between clonal composition
11 and the index lesion is also important and clinically relevant for both molecular diagnostics
12 and focal therapy (Lamb et al. 2017; Reiter et al. 2019; Erickson et al. 2021). Indeed, it would
13 help and support treatment decision-making by using new markers allowing to determine
14 whether cells are indicative of aggressive disease or to predict sensitivity to treatment.

15 One of the challenges to understand the tumour heterogeneity is that origin of
16 mutations occurring in cancer can be hereditary or somatic. Although identification of inherited
17 mutations is relatively straightforward, these are only responsible for 5 to 10% of all cancer
18 (Nagy et al. 2004; Garber and Offit 2005; Leon et al. 2021). By contrast, post-developmental
19 somatic genetic alterations are usually only present in a small fraction of clonally-expanding
20 cells but constitute the most common cause of cancer (Milholland et al. 2017). To identify
21 these somatic mutations *in situ*, techniques such as laser capture microdissection have been
22 employed, but this requires pre-knowledge to isolate a specific cell type or region of interest
23 from a tissue section (Asp et al. 2020) and so limits the ability to undertake a *de novo* spatial
24 clonal analysis. Recently, these limitations have been overcome by spatial transcriptomics,
25 which allows the analysis of gene expression profiles in a tissue sample while preserving
26 spatial tissue architecture. This approach captures transcripts *in situ*, with sequencing of

1 barcoded reads carried out *ex situ* and then mapped back to the cells of origin (Larsson et al.
2 2021; Ståhl et al. 2016). This cutting-edge technology permits visualisation and in-depth
3 analysis of intra-tumoral heterogeneity and could permit spatial analysis of clonal evolution.

4 Clonal evolution and, more precisely, the relationship between clones and subclones
5 is often represented and visualised by phylogenetic trees (Beerenwinkel et al. 2015; Schwartz
6 and Schäffer 2017). These phylogenetic trees have been used mainly in recent years to study
7 data derived from DNA sequencing (Schwartz and Schäffer 2017). However, to use spatial
8 transcriptomics to study clonal evolution, it is necessary to know whether RNA can also be
9 used to determine clonal phylogenetic hierarchies. In this meta-analysis, we investigate the
10 correlation between DNA sequencing data and RNA sequencing data using phylogenies
11 derived from inferred single-nucleotide variants (SNV) and copy-number variants (CNV) in
12 order to determine whether transcriptome-derived phylogenies can accurately reflect genome-
13 based phylogenies.

14 **Results**

15

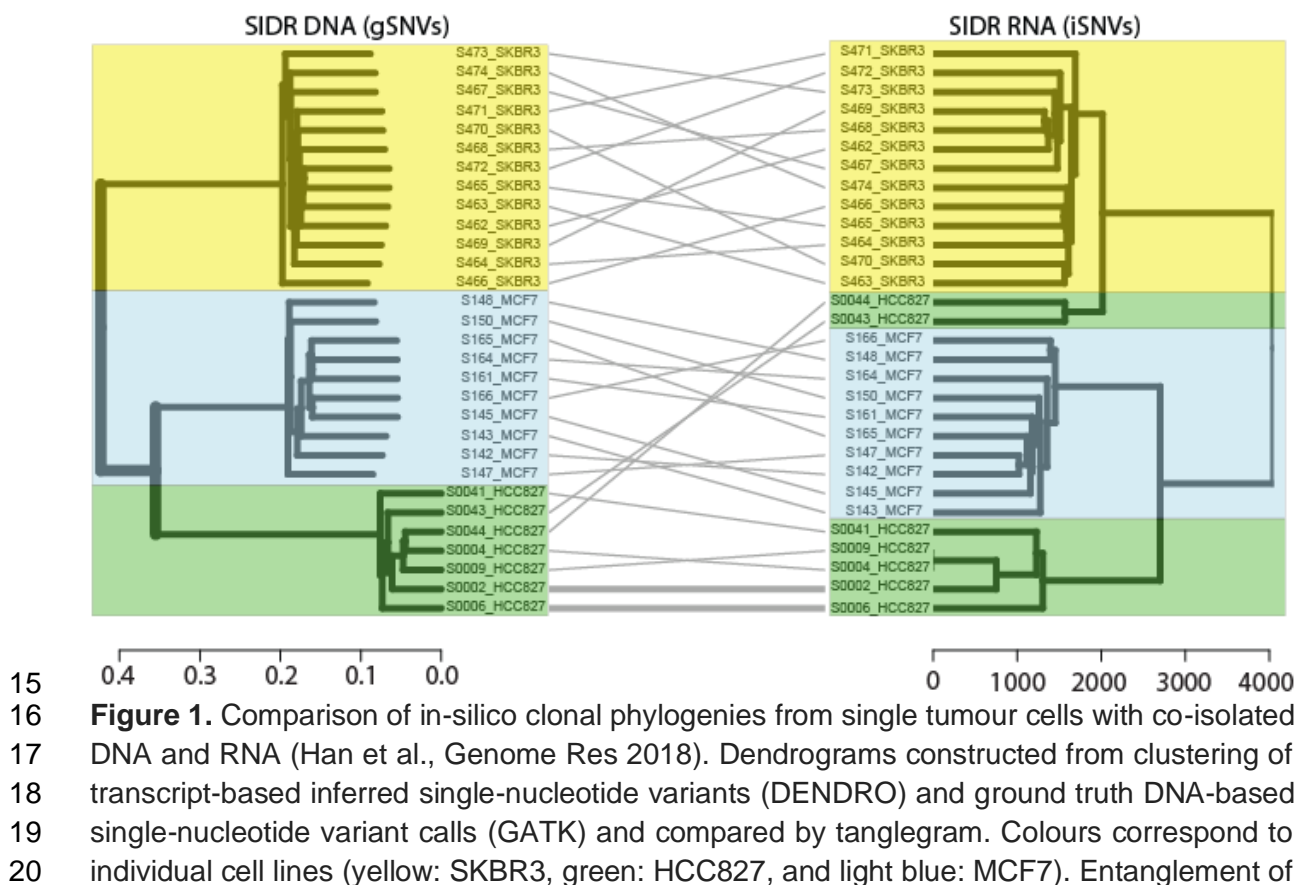
16 ***Transcriptome and Genome Derived clonal phylogenies from single cancer cells***

17

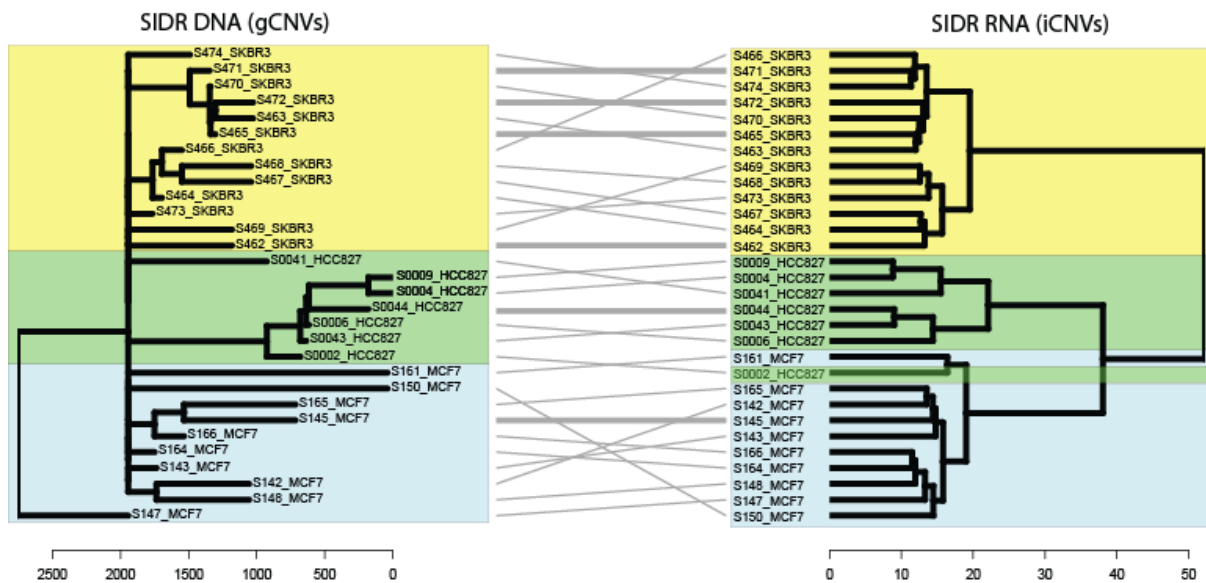
18 In order to benchmark performance of transcriptome-derived phylogenies, we first
19 identified an individual cancer cell dataset with simultaneously isolated DNA and RNA (SIDR)
20 from single cells (Han et al. 2018). The authors performed SIDR resulting in paired DNA and
21 RNA nucleic acid extractions from isolated single cells of three different cancer cell lines:
22 HCC827, MCF7 and SKBR3. They then performed whole-genome sequencing (WGS) and
23 RNA-sequencing on the extracted nucleic acids. Given the cell purity, we hypothesized that
24 WGS and RNA sequencing data from these individual cancer cells could be analyzed in an
25 “in-silico” experiment to benchmark performance of transcriptome and genome-derived
26 phylogenies.

27

1 We performed secondary analyses of the published, publicly available DNA and RNA
2 sequencing data from Han et al (Han et al. 2018). After quality control (Han et al. 2018), we
3 identified a total of 30 cells that had both sufficient quality DNA and RNA sequencing data,
4 resulting in a dataset of a total of 10 MCF7 cells, 7 HCC827 cells, and 13 SKBR3 cells for
5 analysis. We performed genomic SNV (gSNV) and inferred RNA-based SNV (iSNV) analyses
6 from all cells, derived dendrograms, and performed tanglegram analysis to compare gSNV
7 and iSNV dendrograms. In analysis of gSNVs and iSNVs, we observed a high concordance
8 of transcriptome and genomic phylogenies (**Figure 1**, entanglement = 0.097). Next, we
9 performed genomic CNV (gCNV) and inferred RNA-based CNV (iCNV) analyses from all cells,
10 derived dendrograms, and performed tanglegram analysis to compare gCNV and iCNV
11 dendrograms. In analysis of gCNVs and iCNVs, we also observed a high concordance of
12 transcriptome and genomic phylogenies (**Figure 2**, entanglement = 0.11). We therefore
13 concluded that RNA-derived inference of genomic SNVs and CNVs in three purified single cell
14 populations generated strong phylogenetic concordance.



1 the phylograms was 0.097 (an entanglement value of 1 corresponds with full entanglement of
2 two phylograms, whereas an entanglement value of 0 corresponds with no entanglement).
3



4
5
6 **Figure 2.** Comparison of in-silico clonal phylogenies from single tumour cells with co-isolated
7 DNA and RNA (Han et al., Genome Res 2018). Dendrograms constructed from clustering of
8 transcript-based inferred copy-number variants (inferCNV) and ground truth DNA-based copy
9 number variant calls (WGS-Ginkgo) and compared by tanglegram. Colours correspond to
10 individual cell lines (yellow: SKBR3, green: HCC827, and light blue: MCF7). Entanglement of
11 the phylograms was 0.11 (an entanglement value of 1 corresponds with full entanglement of
12 two phylograms, whereas an entanglement value of 0 corresponds with no entanglement). As
13 adapted from Erickson et al., Nature, 2022, Extended Data Fig. 1a.

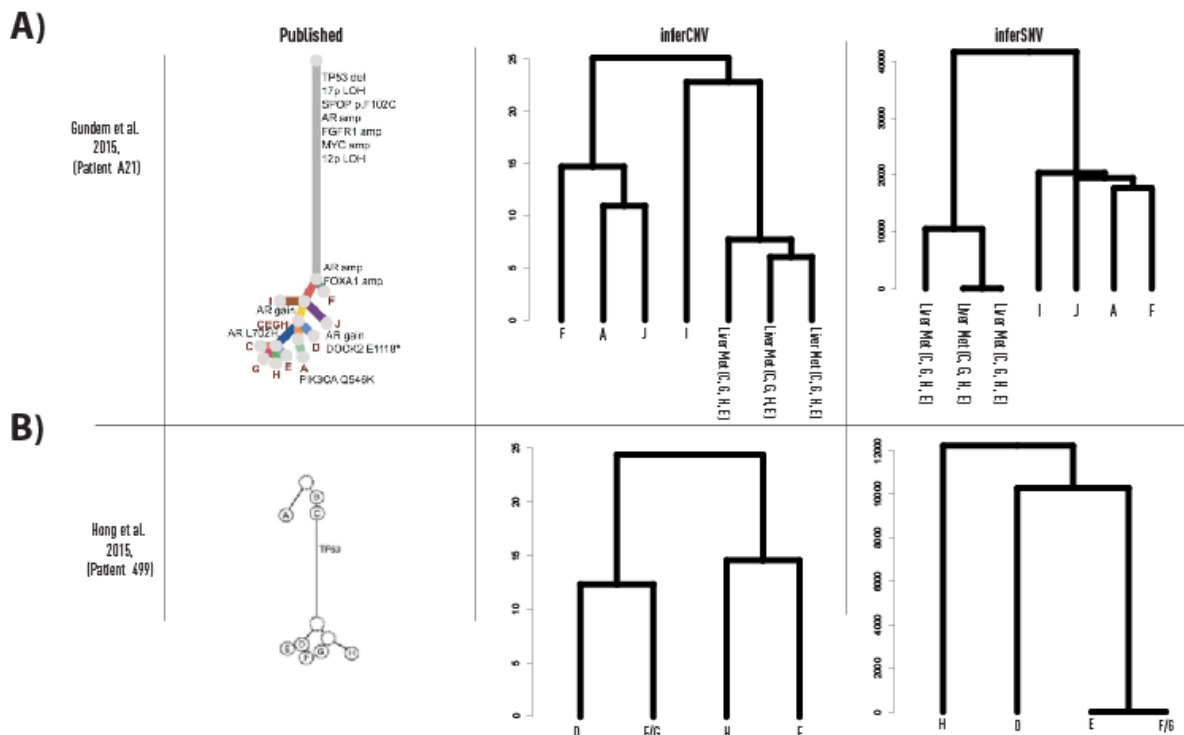
14 15 **Transcriptome and Genome Derived clonal phylogenies from bulk prostate cancer** 16 **sequencing**

17
18 Having established high *in-silico* concordance of transcriptome and genome-derived
19 phylogenies, we then sought to study prostate cancer sequencing data from patients with
20 paired DNA and RNA extracted from the same tumors. Gundem and colleagues reported
21 WGS data from 55 disseminated tumor samples, from 10 patients that underwent rapid-
22 autopsy after death due to prostate cancer (Gundem et al. 2015). A subset of $n = 7$ tumor
23 specimens from patient A21 also underwent RNA-sequencing (Bova et al. 2016).

24
25 We performed secondary analyses of RNA sequencing data from Bova et al. and
26 obtained iSNV and iCNV calls. From the iSNV and iCNV calls, we separately performed

1 phylogenetic analyses through hierarchical clustering, resulting in iSNV and iCNV derived
 2 dendrograms (**Figure 3a**). In both iSNV and iCNV analyses, liver metastases (C, G, H, E)
 3 clustered together. In both iSNV and iCNV analyses, Clones F, A and J also clustered
 4 together. Clone I, clustered together with the liver metastases in iCNV analyses, but not in the
 5 iSNV analyses. Taken together, the iSNV and iCNV dendrograms reflect the manually
 6 assembled clonal phylogeny published by Gundem et al, (Gundem et al. 2015).

7



8

9 **Figure 3.** Comparison of published DNA-based prostate cancer clonal phylogenies and
 10 transcript-based inferred single-nucleotide and copy-number variant derived dendrograms. a)
 11 Phylogeny from patient A21, as published and reproduced from Gundem et al., Nature, 2015.
 12 Transcript data were available only for a subset of specimens. b, Phylogeny from patient 498,
 13 as published and reproduced from Hong et al., Nat. Comms, 2015. Transcript data available
 14 for a subset of specimens. inferCNV-based clonal phylogenies adapted from Erickson et al.,
 15 Nature, 2022, Extended Data Fig. 1b.

16

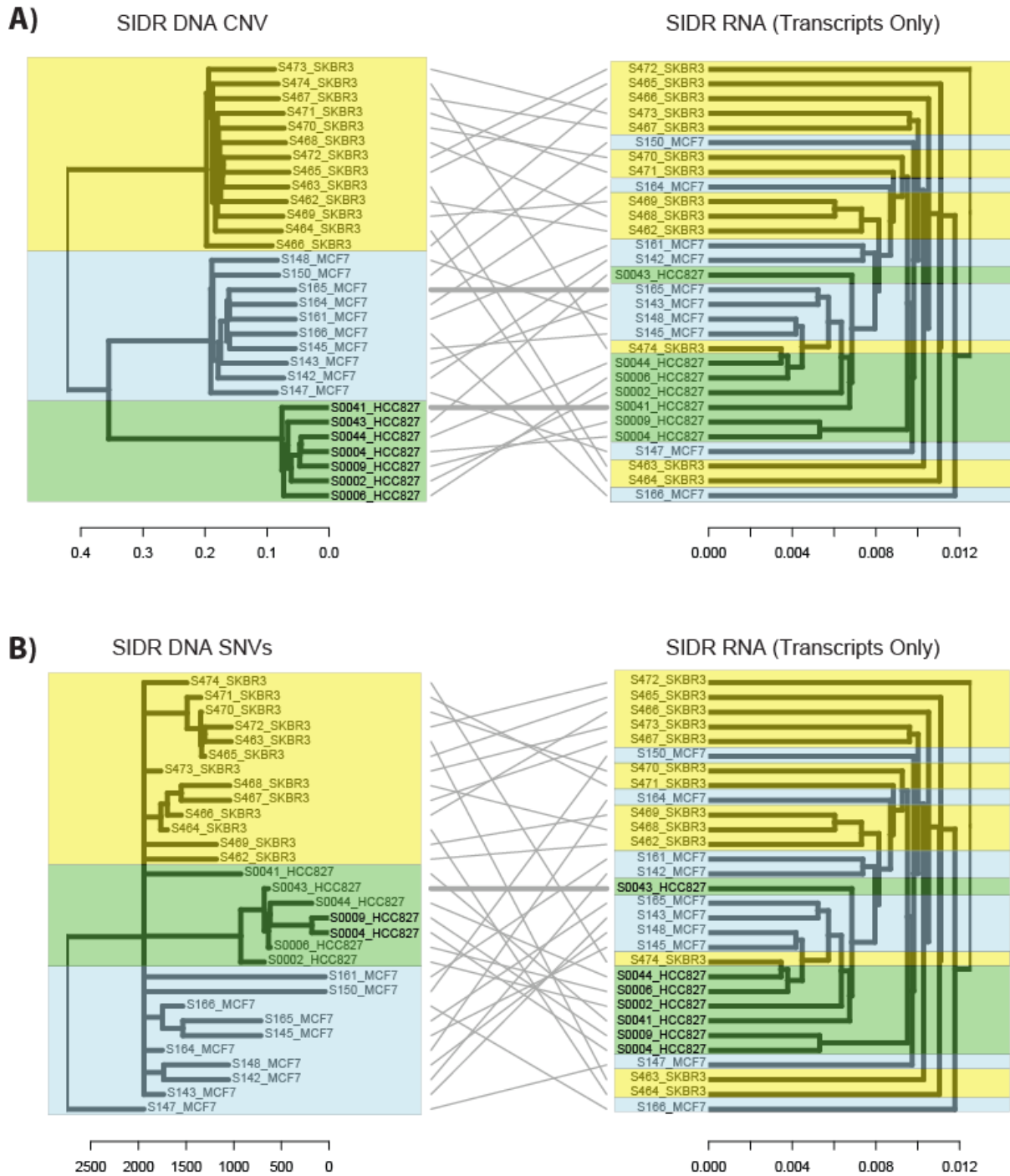
17 Next, we analyzed data from patient 498, analyzed by Hong et al. (Hong et al. 2015).

18 This patient's primary prostate cancer progressed to distant skeletal metastases, which then
 19 further re-seeded the prostatic bed. Of the n = 7 reported specimens, a total of n = 4 also
 20 underwent RNA sequencing. We performed secondary analyses of the RNA sequencing data
 21 and obtained iSNV and iCNV calls. From the iSNV and iCNV calls, we separately performed

1 phylogenetic analyses through hierarchical clustering, resulting in iSNV and iCNV derived
2 dendrograms (**Figure 3b**). In contrast to the results from Gundem et al., both iSNV and iCNV
3 presenting differing tree patterns as compared to one another.

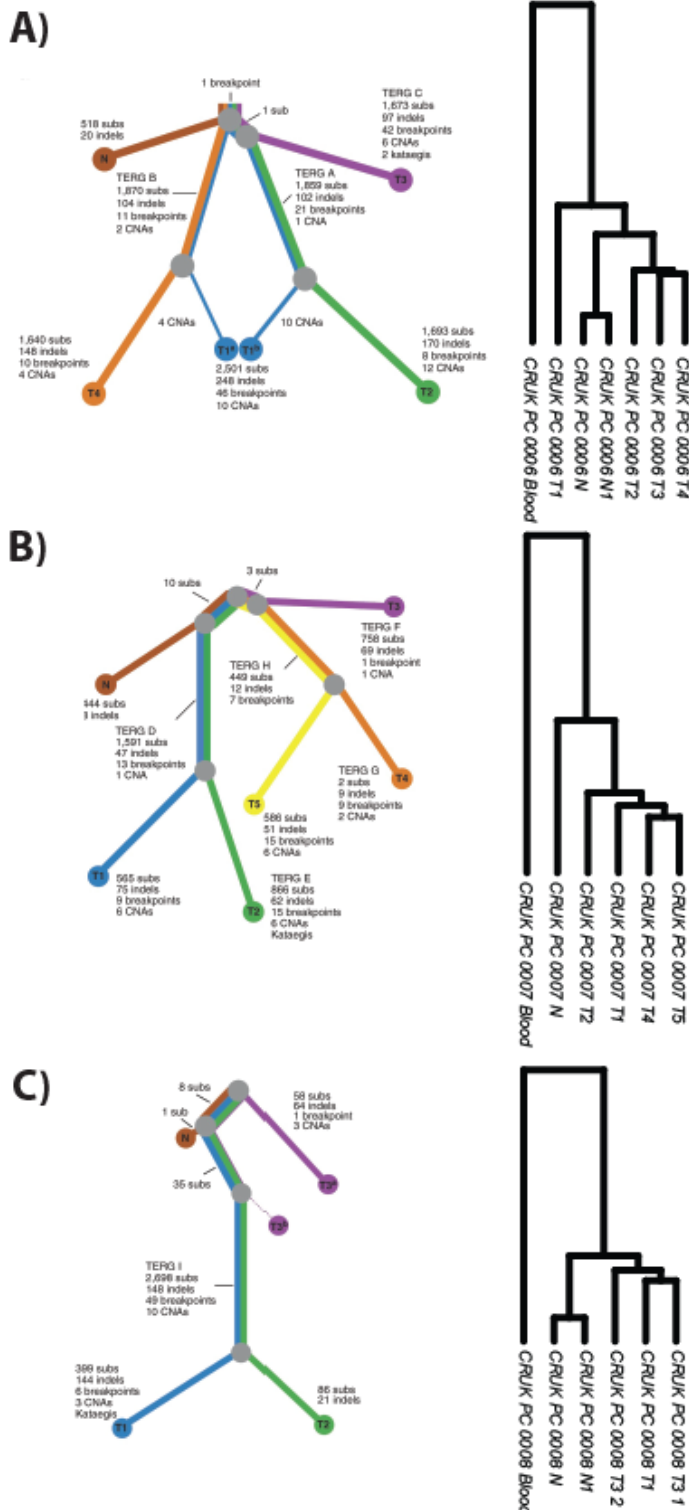
4
5 We then analyzed data from primary prostate cancer cases 6, 7 and 8, analyzed by
6 Cooper et al., who each underwent radical prostatectomy, from which multiple tissue punches
7 of both normal and tumor regions were sampled (Cooper et al. 2015). The samples then
8 underwent WGS, which were subsequently analyzed and tumor phylogenies were manually
9 produced. From a subset of the same specimens, adjacent tissue sections were taken and
10 subjected to RNA microarray analysis. Additionally, each patient had a blood sample taken,
11 that also underwent RNA microarray analysis. Being microarray data, we were unable to
12 derive iSNV and iCNVs. Therefore, we built a custom pipeline to analyze and cluster the RNA
13 microarray data directly, to generate hierarchical clustering represented as a dendrogram. To
14 benchmark this pipeline, we first compared gCNV and gSNV to SIDR data (Supplementary
15 Figure) and observed entanglement values of 0.21 and 0.16 respectively. Having established
16 this pipeline, we then applied it to the microarray data from Cooper et al to generate
17 dendrograms. These dendrograms were then analyzed in comparison to the published WGS-
18 based gDNA phylogenies (**Figure 4**). In all three patients, the blood specimen clustered
19 separately from the prostate tumor and normal tissue specimens. In cases 7 and 8, the
20 (multiple) normal tissue specimens clustered together and distinctly clustered separately from
21 the tumors, whereas in case 6 the two normals clustered with T2, T3 and T4, separate from
22 T1. Taken together, RNA-microarray derived dendrograms were able to recapitulate manually
23 assembled WGS-derived gDNA phylogenies.

24



1
2 **Supplementary Figure 1.** Comparison of in-silico clonal phylogenies from single tumour cells
3 with co-isolated DNA and RNA (Han et al., Genome Res 2018). **A)** Dendrograms constructed
4 from ground truth DNA-based copy number variant calls (WGS-Ginkgo) and direct transcripts
5 (hierarchical clustering) and compared by tanglegram. Colours correspond to individual cell
6 lines (yellow: SKBR3, green: HCC827, and light blue: MCF7). Entanglement of the phylograms
7 was 0.21 (an entanglement value of 1 corresponds with full entanglement of two phylograms,
8 whereas an entanglement value of 0 corresponds with no entanglement). **A)** Dendrograms
9 constructed from ground truth DNA-based single-nucleotide variant calls (DENDRO) and
10 direct transcripts (hierarchical clustering) and compared by tanglegram. Colours correspond
11 to individual cell lines (yellow: SKBR3, green: HCC827, and light blue: MCF7). Entanglement

1 of the phylograms was 0.16 (an entanglement value of 1 corresponds with full entanglement
 2 of two phylograms, whereas an entanglement value of 0 corresponds with no entanglement).
 3

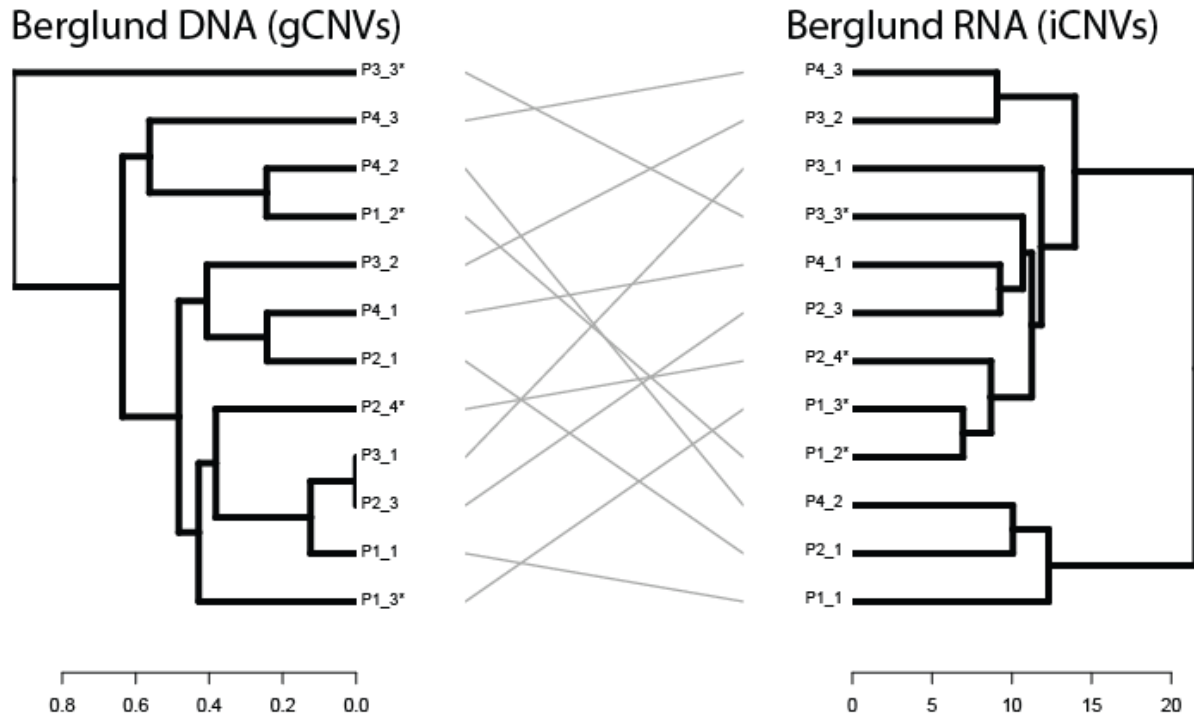


4
 5 **Figure 4.** Comparison of published DNA-based (WGS) phylogenetic trees (left) as compared
 6 to novel RNA-based (RNA Microarray) phylogenies (right) from Cooper et al., 2015. A)
 7 Phylogenies from patient CRUK0006, B) Phylogenies from patient CRUK0007, C)
 8 Phylogenies from patient CRUK0008. RNA phylogenies include blood samples not presented
 9 in DNA-based phylogenetic trees.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

Transcriptome and Genome Derived clonal phylogenies from bulk WGS and spatial transcriptomics from multi-region prostate cancer sequencing data

Next, we then sought to determine the ability of spatial transcriptome derived tumor phylogenies to recapitulate gDNA based phylogenies. Spatial transcriptomics generates transcriptome signal from poly-A captured short 3' RNA sequences of up to 200 bp length, sufficient for hg38 alignment and, we deduced, sufficient to enable iCNV analysis. Berglund and colleagues performed spatial transcriptomics (ST) (Ståhl et al. 2016) on a total of n = 12 prostate tissue regions from a patient that underwent radical prostatectomy (Berglund et al. 2018). Of these sections, a total of n = 4 were detected to have prostate cancer. The authors also performed WGS on adjacent serial sections from each of these 12 tissue sections, as well as a matched blood specimen from the same patient. Given that WGS is not spatially resolved, we performed 'pseudo-bulked' iCNV analyses on ST data from all 12 sections, and generated a clonal phylogeny in the form of a dendrogram. We also performed gDNA CNV calling from each of the 12 sections to generate a clonal phylogeny which was represented as a dendrogram. We then compared the iCNV and gCNV derived dendrograms using a tanglegram and observed a degree of concordance consistent with the resolution of the data (**Figure 5**, entanglement = 0.35). Interestingly, three of the tumor regions (P2_4, P1_3, P1_2) clustered together in the iCNV analysis, whereas they were represented on different subclusters in the gCNV phylogeny, suggesting that the iCNV approach may have generated a more accurate clustering in this case.



1
2 **Figure 5.** Comparison of DNA-based (WGS) phylogenetic trees (left) as compared to
3 transcript-based inferCNV clonal phylogenies (right) from Berglund et al., 2018. DNA
4 dendrogram constructed using patient-matched blood sample as a reference: such data were
5 not available for inferCNV. Entanglement of the phylograms was 0.35 (an entanglement value
6 of 1 corresponds with full entanglement of two phylograms, whereas an entanglement value
7 of 0 corresponds with no entanglement). A label with the ending of * represents a section
8 containing histologically detected cancer.
9

10 Discussion

11 Results from single-cancer cells demonstrate that transcriptome-derived iCNV and
12 iSNV phylogenies are highly concordant with ground truth gDNA based phylogenies. In our
13 in-silico analyses, the analyzed data represent a highly selected and well controlled set of
14 cells, with a 1:1 pairing of data resulting in extremely low entanglement values of the resultant
15 tanglegrams. These results are in line with findings by Han et al., where they reported positive
16 correlations for all three cell lines between gCNV and mRNA expression levels that were
17 binned across the genome (Han et al. 2018). Our quantitative results in single-cells were
18 supported by qualitative comparisons in prostate cancer cells where we did not have access
19 to all ground truth data to enable a true like-to-like comparison.

1 There are limitations to consider in the construction of transcriptome-derived inferred
2 phylogenies. First, the design and resolution of the genetic sequencing technologies can
3 greatly affect the 'resolved signal'. For example, only 2% of the entire genome is translated
4 into proteins (International Human Genome Sequencing Consortium 2004), and thus the
5 genomic coverage of the transcriptome represents a sub-fraction of potential data for mapping
6 tumor phylogenies. This is further compounded by variable coverage within transcripts
7 themselves: many modern scRNAseq and spatial transcriptomics techniques, such as
8 Chromium and Visium offered by 10x Genomics, perform polyA capture, resulting in
9 sequencing of 75-300 bp near the end of transcripts. Further, for iSNV approaches (Petti et
10 al. 2019; Zhou et al. 2020), the coverage of transcribed SNV loci can be extremely low being
11 confined to the exome. Potential issues with iSNVs seem to be mitigated in iCNV approaches
12 (Patel et al. 2014; Gao et al. 2021; Elyanow et al. 2021), which incorporate machine learning
13 algorithms to bin genomically adjacent transcripts. Additionally transcriptional regulation
14 programs (Lee and Young 2013; Bradner et al. 2017; Davies et al. 2020) can affect
15 transcription without any changes to copy-number status: these may result in false positive or
16 negatives in iCNV analyses. Indeed, Han et al observed a discrepancy in Chromosome 3
17 gCNV calls and expression profiles (Han et al. 2018). Finally, one key factor affecting the
18 ability of iCNV/iSNV (as well as gCNV and gSNV) approaches is use of well annotated
19 references. All of the patient-derived WGS analyses in the data used in this publication had
20 access to reference blood controls for calling gCNVs and gSNVs. Such data are not often
21 taken or obtained for RNA sequencing, and thus are unavailable for iCNV and iSNV calling.
22 This can also be further compounded by tissue or cell-of-origin transcriptional programs
23 unrelated to copy-number alterations. Spatial transcriptomic data offers the opportunity to
24 compensate for this through selection of histologically normal regions as control references.

25 As the tumor evolution community moves increasingly to single cell and spatial
26 resolution, our ability to resolve clonal and subclonal tumor evolution patterns will greatly
27 increase. Our results underscore the need for proper reference sets when calling iCNV and

1 iSNV derived clonal phylogenies. These issues may be partly mitigated by next-generation
2 iCNV and iSNV algorithms that incorporate both into combined iSNV+iCNV phylogenies (Gao
3 et al. 2022). Other approaches incorporating evolutionary game theory through mathematical
4 models could aid in resolving clonal phylogenies (Wöflfl et al. 2022). Further work will also
5 need to be done to identify and control for non copy-number alteration derived transcriptional
6 regulation leading to further refinements in the ability of transcript-based clonal phylogenies to
7 resolve ground truth.

8 **Methods**

9

10 ***Data Acquisition***

11

12 In order to benchmark and validate methods to generate phylogenies derived from inferred
13 single-nucleotide variants and copy-number variants, we reviewed the literature and found a
14 recent publication which simultaneously extracted both DNA and RNA, from the same exact
15 single tumor cells, and performed whole genome and whole transcriptome sequencing (Han
16 et al. 2018). These public datasets contained data from 38 single cells that had been subject
17 to simultaneous WGS and RNAseq using the SIDR methodology. Han et al describe a quality
18 control process to determine which cells were satisfactorily sequenced for downstream
19 analysis, leaving a total of 30 paired samples that passed all qc metrics.

20

21 Next, we reviewed the literature for publications and available data from patients with prostate
22 cancer, who had both conventional bulk DNA and RNA sequencing applied to the same
23 specimen, and from patients that had 3 or more total specimens. We identified patient A21
24 (Gundem et al. 2015; Bova et al. 2016), patient 498 (Hong et al. 2015). For further validation
25 and comparison, WGS and RNA-microarray data were obtained from cases 6, 7 and 8 from
26 Cooper et al. (Cooper et al. 2015).

27

1 Lastly, we obtained paired WGS sequencing data and paired Spatial Transcriptomics data
2 from the n = 12 regions from a single patient in a recent publication (Berglund et al. 2018).

3 4 ***Analysis of Single Cell Data***

5 6 ***Quality Control of Single-Cell Whole Genome Sequencing Data***

7
8 Only 38 paired cells were available with both scWGS and scRNAseq (Han et al. 2018). After
9 removing the individual cells that failed either scWGS or scRNAseq QC left only 30 in
10 common.

11 12 ***DNA Sequencing Preprocessing of Single-Cell Whole Genome Sequencing Data***

13
14 Paired end sequencing data was aligned against the GRCh38 reference genome with the
15 Burrow-Wheeler Aligner (0.7.17).

16 17 ***iSNV Calling from Single-Cell Whole Genome Sequencing Data***

18
19 WGS variants were called using a pipeline broadly based on the GATK best practice Germline
20 short variant discovery (SNPs + Indels) workflow using Picard (2.23.0) and GATK (4.1.7.0).
21 This consisted of pre-processing the raw alignment to mark duplicate reads and perform base
22 recalibration. Raw variants were called using GATK HaplotypeCaller in GVCF mode followed
23 by GATK GenotypeGVCFs. Finally the raw variants were filtered to generate a downstream
24 analysis ready cell by variant dataset.

25
26 The processed variants were converted to an Identity by State matrix, clustered and converted
27 to dendrogram format in R using the SNPrelate package(Zheng et al. 2012)(Zheng et al.
28 2017).

29 30 ***gCNV Calling from Single-Cell Whole Genome Sequencing Data***

31
32 After preprocessing and QCing, n = 30 cells remained, and were then analyzed by Gingko
33 (Garvin et al. 2015). BAM files were converted to .BED files using bamToBed in BedTools.

1 We utilized a variable bin size of 50 kb, with 101 bp reads (Han et al. 2018). The clustering of
2 CNV's was performed using ward linkage and Euclidean distance as the distance metric.
3 Copy-Number tree results were downloaded in Newick format for further downstream analysis.

4 5 ***RNA Sequencing Preprocessing of Single-Cell Whole Transcriptome Sequencing Data***

6
7 Paired end sequencing data was aligned against the GRCh38 reference genome with STAR
8 (2.7.3a) with per-sample 2-pass mapping and annotation with comprehensive gene annotation
9 data from GENCODE GRCh38. Gene counts per cell were tabulated from aligned data using
10 the featureCounts function from the Subread (1.6.4) package.

11 12 ***iSNV Calling from Single-Cell Whole Transcriptome Sequencing Data***

13
14 iSNV calling from RNAseq data was performed according to the pipeline outlined by Zhou et
15 al and based on GATK best practices (Zhou et al. 2020). The STAR aligned data underwent
16 sorting, annotation with read group information, deduplication, SplitNCigarReads,
17 realignment, and base recalibration, before variant calling with GATK (3.8.0) HaplotypeCaller.
18 Raw iSNVs were processed by DENDRO to calculate a genetic divergence matrix between
19 cells and to generate a phylogeny using hierarchical clustering (ward.D method).

20 21 ***iCNV Calling from Single-Cell Whole Transcriptome Sequencing Data***

22
23 Data were analyzed using R version 4.0.1, and inferCNV (version 1.4.0) ([CSL STYLE
24 ERROR: reference with no printed form.]). A merged file from the previously described pre-
25 processing steps, containing feature counts for each cell, as well as a gene position file, and
26 an annotation file were generated for input to inferCNV. An inferCNV object was created with
27 no defined reference group. After creation of the InferCNV object, inferCNV was ran with the
28 following parameters: cutoff = 0.1, cluster_by_groups = FALSE, denoise = TRUE, HMM =
29 TRUE.

30 31 ***Comparison of Dendrograms from Single-Cells***

32

1 For comparison of dendrograms created by WGS-CNVs (Gingko) and inferred CNV's from
2 RNA (InferCNV), the `clust2.newick` and `infercnv.21_denoised.observations_dendrogram.txt`
3 files were imported into R and analyzed with packages *dendextend* and *phylogram*.

4

5 ***Analysis of transcript derived phylogenies***

6

7 RNA counts were analyzed, by comparing individual gene count values to the median (MED)
8 and standard deviation (SD) values of global RNA count values per sample: if the count value
9 was less than MED-SD, then it was assigned a value of -1, else if the count value was greater
10 than MED+SD, then it was assigned a value of +1, else it was assigned 0. The resultant values
11 from each sample or cell were converted into a phydat object using *phangorn*'s function
12 `phyDat()`, with the parameters `type="USER"`, `levels = c('-1', '0', '1')`. Pairwise distances
13 between cells or tissue samples were calculated using the *phangorn* `dist.ml()` function with
14 previously described `phyDat()` object as input. UPGMA clustering was applied using the
15 *phangorn* `upgma()` function and converted to a dendrogram using the *dendextend* function
16 `as.dendrogram()`.

17

18 ***Analysis of Spatial Transcriptomics Data***

19

20 ***CNV Calling from Spatial Transcriptomics Data***

21

22 Data were analyzed as previously described (Erickson et al. 2022) with the following
23 exceptions. Original 1k array Spatial Transcriptomics data were obtained. As gCNV
24 comparison data were from whole sections, all ST count data were 'pseudo-bulked' within
25 sections, resulting in 12 pseudobulked count matrices for analyses. InferCNV was ran using
26 standard parameters with no reference set. The resultant
27 `infercnv.observations_dendrogram.txt` dendrogram was used for downstream tanglegram
28 analysis.

29

30 ***Comparison of Dendrograms from WGS and ST***

1 The original outputs for CNV calling from Berglund et al., were not available, and the
2 ReadDepth package used to generate the calls has since been deprecated by the author
3 (Miller). Thus, we ran a new pipeline using the WGS data from Berglund et al (Berglund et al.
4 2018). FASTQ files were obtained and aligned to HG38. Battenberg CNV analyses (Nik-Zainal
5 et al. 2012) were performed using the matched reference blood FASTQ data as the reference.

6 ***Copy number calling with Battenberg***

7 The Battenberg package (v2.2.10) was used to determine copy number, and estimate tumour
8 purity and ploidy from WGS data. Impute2 (v2.3.0) was used with GRCh38 loci for phasing
9 germline heterozygous SNPs. The Battenberg pipeline was run with the following parameters:
10 segmentation_gamma = 10, phasing_gamma = 10, platform_gamma = 1, min_ploidy = 1.6,
11 max_ploidy = 4.8, min_rho = 0.13, max_rho = 1.02.

12 The *recal_subclones.txt* text files were downloaded for each of the 12 prostate tissues, and
13 processed through a custom pipeline as follows. Battenberg CNV segments were binned into
14 1200 bp segments and aligned, generating n = 2439447 bins across the genome. CN
15 amplifications and deletions were called at thresholded values of -1.5 and 2.5 respectively.
16 Next, the processed bins from all samples were merged to create a CN bin matrix. CN calls
17 for segments that were shared for all samples were dropped, resulting in a final matrix
18 containing n = 28 discordant CN calls.

19 This CN matrix was then used similarly as described by Berglund et al., with the R package
20 *pvclust*, and n = 1000 bootstraps. The structure of the cluster was converted to a dendrogram
21 using the R package *dendrogram* for comparison to the inferCNV dendrogram via a
22 tanglegram using the *dendextend* package (step2side).

23 24 **Data Access**

25 Data from single cell experiments (Han et al. 2018) were previously deposited to ENA:
26 PRJEB20144 (WGS) and PRJEB20143 (RNA). All sequence data from patient 499 (Hong et

1 al. 2015) samples were previously deposited into the EGA Sequence Read Archive under
2 accession number EGAS00001000942. RNA sequencing data from patient A21 (Bova et al.
3 2016) were previously deposited into the EGA Sequence Read Archive under accession
4 number EGAS00001001659. Sequencing data from patient 1 (Berglund et al. 2018) were
5 previously deposited at the European Genome–Phenome Archive (EGA), hosted by the
6 European Bioinformatics Institute (EBI), under the accession number EGAS0000100300.

7 **Competing Interest Statement**

8 The authors have no conflicts of interest to declare.

9 **Acknowledgements**

10 Computation used the Oxford Biomedical Research Computing (BMRC) facility, a joint
11 development between the Wellcome Centre for Human Genetics and the Big Data Institute
12 supported by Health Data Research UK and the NIHR Oxford Biomedical Research Centre.

13 The views expressed are those of the author(s) and not necessarily those of the NHS, the
14 NIHR or the Department of Health.

15 Author contributions: A.E., A.L., and I.M. conceived the study. A.E., T.R., and S.R. performed
16 computational experiments. All authors interpreted the data and wrote the manuscript.

17
18
19
20
21

1 References

- 2
- 3 Asp M, Bergenstråhle J, Lundeberg J. 2020. Spatially Resolved Transcriptomes-Next
4 Generation Tools for Tissue Exploration. *Bioessays* **42**: e1900221.
- 5 Beerenwinkel N, Schwarz RF, Gerstung M, Markowitz F. 2015. Cancer evolution:
6 mathematical models and computational inference. *Syst Biol* **64**: e1–25.
- 7 Berglund E, Maaskola J, Schultz N, Friedrich S, Marklund M, Bergenstråhle J, Tarish F,
8 Tanoglidi A, Vickovic S, Larsson L, et al. 2018. Spatial maps of prostate cancer
9 transcriptomes reveal an unexplored landscape of heterogeneity. *Nat Commun* **9**: 2419.
- 10 Black JRM, McGranahan N. 2021. Genetic and non-genetic clonal diversity in cancer
11 evolution. *Nat Rev Cancer* **21**: 379–392.
- 12 Bova GS, Kallio HML, Annala M, Kivinummi K, Högnäs G, Häyrynen S, Rantapero T,
13 Kivinen V, Isaacs WB, Tolonen T, et al. 2016. Integrated clinical, whole-genome, and
14 transcriptome analysis of multisampled lethal metastatic prostate cancer. *Cold Spring*
15 *Harb Mol Case Stud* **2**: a000752.
- 16 Bradner JE, Hnisz D, Young RA. 2017. Transcriptional Addiction in Cancer. *Cell* **168**: 629–
17 643.
- 18 Cooper CS, Eeles R, Wedge DC, Van Loo P, Gundem G, Alexandrov LB, Kremeyer B,
19 Butler A, Lynch AG, Camacho N, et al. 2015. Analysis of the genetic phylogeny of
20 multifocal prostate cancer identifies multiple independent clonal expansions in
21 neoplastic and morphologically normal prostate tissue. *Nat Genet* **47**: 367–372.
- 22 Davies A, Zoubeidi A, Selth LA. 2020. The epigenetic and transcriptional landscape of
23 neuroendocrine prostate cancer. *Endocr Relat Cancer* **27**: R35–R50.
- 24 Elyanow R, Zeira R, Land M, Raphael BJ. 2021. STARCH: copy number and clone inference
25 from spatial transcriptomics data. *Phys Biol* **18**: 035001.
- 26 Erickson A, Hayes A, Rajakumar T, Verrill C, Bryant RJ, Hamdy FC, Wedge DC, Woodcock
27 DJ, Mills IG, Lamb AD. 2021. A Systematic Review of Prostate Cancer Heterogeneity:
28 Understanding the Clonal Ancestry of Multifocal Disease. *Eur Urol Oncol* **4**: 358–369.
- 29 Erickson A, He M, Berglund E, Marklund M, Mirzazadeh R, Schultz N, Kvastad L, Andersson
30 A, Bergenstråhle L, Bergenstråhle J, et al. 2022. Spatially resolved clonal copy number
31 alterations in benign and malignant tissue. *Nature* **608**: 360–367.
- 32 Gao R, Bai S, Henderson YC, Lin Y, Schalck A, Yan Y, Kumar T, Hu M, Sei E, Davis A, et
33 al. 2021. Delineating copy number and clonal substructure in human tumors from single-
34 cell transcriptomes. *Nat Biotechnol* **39**: 599–608.
- 35 Gao T, Soldatov R, Sarkar H, Kurkiewicz A, Biederstedt E, Loh P-R, Kharchenko PV. 2022.
36 Haplotype-aware analysis of somatic copy number variations from single-cell
37 transcriptomes. *Nat Biotechnol*. <http://dx.doi.org/10.1038/s41587-022-01468-y>.
- 38 Garber JE, Offit K. 2005. Hereditary cancer predisposition syndromes. *J Clin Oncol* **23**: 276–
39 292.
- 40 Garvin T, Aboukhalil R, Kendall J, Baslan T, Atwal GS, Hicks J, Wigler M, Schatz MC. 2015.
41 Interactive analysis and assessment of single-cell copy-number variations. *Nat Methods*

- 1 **12**: 1058–1060.
- 2 Greaves M, Maley CC. 2012. Clonal evolution in cancer. *Nature* **481**: 306–313.
- 3 Gundem G, Van Loo P, Kremeyer B, Alexandrov LB, Tubio JMC, Papaemmanuil E, Brewer
4 DS, Kallio HML, Högnäs G, Annala M, et al. 2015. The evolutionary history of lethal
5 metastatic prostate cancer. *Nature* **520**: 353–357.
- 6 Gupta GP, Massagué J. 2006. Cancer metastasis: building a framework. *Cell* **127**: 679–695.
- 7 Han KY, Kim K-T, Joung J-G, Son D-S, Kim YJ, Jo A, Jeon H-J, Moon H-S, Yoo CE, Chung
8 W, et al. 2018. SIDR: simultaneous isolation and parallel sequencing of genomic DNA
9 and total RNA from single cells. *Genome Res* **28**: 75–87.
- 10 Hong MKH, Macintyre G, Wedge DC, Van Loo P, Patel K, Lunke S, Alexandrov LB, Sloggett
11 C, Cmero M, Marass F, et al. 2015. Tracking the origins and drivers of subclonal
12 metastatic expansion in prostate cancer. *Nat Commun* **6**: 6605.
- 13 International Human Genome Sequencing Consortium. 2004. Finishing the euchromatic
14 sequence of the human genome. *Nature* **431**: 931–945.
- 15 Lamb AD, Zargar H, Murphy DG, Corcoran NM, Hovens CM. 2017. Disrupting the Status
16 Quo in Prostate Cancer Diagnosis. *Eur Urol* **71**: 193–194.
- 17 Larsson L, Frisé J, Lundeberg J. 2021. Spatially resolved transcriptomics adds a new
18 dimension to genomics. *Nat Methods* **18**: 15–18.
- 19 Lee TI, Young RA. 2013. Transcriptional regulation and its misregulation in disease. *Cell*
20 **152**: 1237–1251.
- 21 Leon P, Cancel-Tassin G, Bourdon V, Buecher B, Oudard S, Brureau L, Jouffe L, Blanchet
22 P, Stoppa-Lyonnet D, Coulet F, et al. 2021. Bayesian predictive model to assess
23 BRCA2 mutational status according to clinical history: Early onset, metastatic phenotype
24 or family history of breast/ovary cancer. *Prostate* **81**: 318–325.
- 25 Milholland B, Dong X, Zhang L, Hao X, Suh Y, Vijg J. 2017. Differences between germline
26 and somatic mutation rates in humans and mice. *Nat Commun* **8**: 15183.
- 27 Miller C. *readDepth*. Github <https://github.com/chrisamiller/readDepth> (Accessed October 9,
28 2022).
- 29 Nagy R, Sweet K, Eng C. 2004. Highly penetrant hereditary cancer syndromes. *Oncogene*
30 **23**: 6445–6470.
- 31 Nik-Zainal S, Van Loo P, Wedge DC, Alexandrov LB, Greenman CD, Lau KW, Raine K,
32 Jones D, Marshall J, Ramakrishna M, et al. 2012. The life history of 21 breast cancers.
33 *Cell* **149**: 994–1007.
- 34 Nowell PC. 1976. The clonal evolution of tumor cell populations. *Science* **194**: 23–28.
- 35 Patel AP, Tirosh I, Trombetta JJ, Shalek AK, Gillespie SM, Wakimoto H, Cahill DP, Nahed
36 BV, Curry WT, Martuza RL, et al. 2014. Single-cell RNA-seq highlights intratumoral
37 heterogeneity in primary glioblastoma. *Science* **344**: 1396–1401.
- 38 Petti AA, Williams SR, Miller CA, Fiddes IT, Srivatsan SN, Chen DY, Fronick CC, Fulton RS,
39 Church DM, Ley TJ. 2019. A general approach for detecting expressed mutations in
40 AML cells using single cell RNA-sequencing. *Nat Commun* **10**: 3660.

- 1 Reiter JG, Baretta M, Gerold JM, Makohon-Moore AP, Daud A, Iacobuzio-Donahue CA, Azad
2 NS, Kinzler KW, Nowak MA, Vogelstein B. 2019. An analysis of genetic heterogeneity in
3 untreated cancers. *Nat Rev Cancer* **19**: 639–650.
- 4 Schwartz R, Schäffer AA. 2017. The evolution of tumour phylogenetics: principles and
5 practice. *Nat Rev Genet* **18**: 213–229.
- 6 Ståhl PL, Salmén F, Vickovic S, Lundmark A, Navarro JF, Magnusson J, Giacomello S, Asp
7 M, Westholm JO, Huss M, et al. 2016. Visualization and analysis of gene expression in
8 tissue sections by spatial transcriptomics. *Science* **353**: 78–82.
- 9 Wölfl B, Te Rietmole H, Salvioli M, Kaznatcheev A, Thuijsman F, Brown JS, Burgering B,
10 Staňková K. 2022. The Contribution of Evolutionary Game Theory to Understanding and
11 Treating Cancer. *Dyn Games Appl* **12**: 313–342.
- 12 Zheng X, Gogarten SM, Lawrence M, Stilp A, Conomos MP, Weir BS, Laurie C, Levine D.
13 2017. SeqArray—a storage-efficient high-performance data format for WGS variant
14 calls. *Bioinformatics* **33**: 2251–2257.
- 15 Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS. 2012. A high-performance
16 computing toolset for relatedness and principal component analysis of SNP data.
17 *Bioinformatics* **28**: 3326–3328.
- 18 Zhou Z, Xu B, Minn A, Zhang NR. 2020. DENDRO: genetic heterogeneity profiling and
19 subclone detection by single-cell RNA sequencing. *Genome Biol* **21**: 10.
- 20 *infercnv*. Github <https://github.com/broadinstitute/infercnv> (Accessed August 19, 2020).
- 21