

TESTING FOR PHYLOGENETIC SIGNAL IN SINGLE-CELL RNA-SEQ DATA

JIRÍ C. MORAVEC¹, ROB LANFEAR², DAVID L. SPECTOR³, SARAH D. DIERMEIER⁴,
AND ALEX GAVRYUSHKIN^{1,5} ✉

ABSTRACT. Phylogenetic methods are emerging as a useful tool to understand cancer evolutionary dynamics, including tumor structure, heterogeneity, and progression. Most currently used approaches utilize either bulk whole genome sequencing (WGS) or single-cell DNA sequencing (scDNA-seq) and are based on calling copy number alterations and single nucleotide variants (SNVs). Here we explore the potential of single-cell RNA sequencing (scRNA-seq) to reconstruct cancer evolutionary dynamics. scRNA-seq is commonly applied to explore differential gene expression of cancer cells throughout tumor progression. The method exacerbates the single-cell sequencing problem of low yield per cell with uneven expression levels. This accounts for low and uneven sequencing coverage and makes SNV detection and phylogenetic analysis challenging. In this paper, we demonstrate for the first time that scRNA-seq data contains sufficient evolutionary signal and can be utilized in phylogenetic analyses. We explore and compare results of such analyses based on both expression levels and SNVs called from scRNA-seq data. Both techniques are shown to be useful for reconstructing phylogenetic relationships between cells, reflecting the clonal composition of a tumor. Both standardized expression values and SNVs appear to be equally capable of reconstructing a similar pattern of phylogenetic relationship. This pattern is stable even when phylogenetic uncertainty is taken in account. Our results open up a new direction of somatic phylogenetics based on scRNA-seq data. Further research is required to refine and improve these approaches to capture the full picture of somatic evolutionary dynamics in cancer.

INTRODUCTION

Phylogenetic analysis is an approach that relies on reconstructing evolutionary relationships between organisms to determine population genetics parameters such as population growth (Kingman 1982; Heled et al. 2008), structure (Müller et al. 2017a) or geographical distribution (Lemey et al. 2009; Lemey et al. 2010). Typically, the reconstructed phylogeny is not the end-goal. Using previously estimated trees, various evolutionary hypotheses can be explored, such as the evolutionary relationship of traits carried by individual taxa (Grafen et al. 1989; Pagel et al. 2004; Freckleton 2012).

Within-organism cancer evolution is increasingly being studied using population genetics approaches, including phylogenetics (Navin et al. 2011; Yuan et al. 2015; Alves et al. 2017; Schwartz et al. 2017; Caravagna et al. 2018; Singer et al. 2018; Alves et al. 2019; Caravagna et al. 2019; Detering et al. 2019; Malikic et al. 2019; Werner et al. 2019; Kuipers et al. 2020), to understand evolutionary dynamics of cancer cell populations. These approaches have shown promise to be developed into therapeutic applications in the personalized medicine framework (Gerlinger et al. 2012; Abbosh et al. 2017; Rao et al. 2020b). Specifically, the clonal composition of tumors, metastasis initiation, development, and timing can be reconstructed using phylogenetic methods (Yuan et al. 2015; Angelova et al. 2018; El-Kebir et al. 2018; Alves et al. 2019). Unlike other evolutionary

¹ DEPARTMENT OF COMPUTER SCIENCE, UNIVERSITY OF OTAGO, NEW ZEALAND

² ECOLOGY, EVOLUTION, AND GENETICS, THE AUSTRALIAN NATIONAL UNIVERSITY, AUSTRALIA

³ COLD SPRING HARBOR LABORATORY, NEW YORK, UNITED STATES OF AMERICA

⁴ DEPARTMENT OF BIOCHEMISTRY, UNIVERSITY OF OTAGO, NEW ZEALAND

⁵ SCHOOL OF MATHEMATICS AND STATISTICS, UNIVERSITY OF CANTERBURY, NEW ZEALAND

E-mail addresses: jiri.moravec@otago.ac.nz, rob.lanfear@anu.edu.au, spector@cshl.edu, sarah.diermeier@otago.ac.nz, ✉ alex@biods.org.

processes prone to events such as hybridization or horizontal gene transfer, population dynamics of somatic cells is underpinned by a strictly bifurcating clonal process driven by cell division. This is in perfect agreement with theoretical assumptions routinely applied in stochastic phylogenetic models such as coalescent (Kingman 1982; Hudson et al. 1990; Posada 2020) or birth-death processes (Aldous 1996; Aldous 2001; Komarova 2006).

From the methodological perspective, however, cancer is an evolutionary process with unique characteristics which are not modeled in conventional phylogenetic approaches. These include a high level of genomic instability with structural changes (gene losses and duplications) which accumulate along with point mutations during the course of growth and evolution (Beerenwinkel et al. 2015; Posada 2015).

Traditional Whole Genome Sequencing (WGS) methods have been instrumental in understanding cancer mutational profiles and oncogene detection (Mardis et al. 2009; Nakagawa et al. 2018). DNA from a tissue sample is isolated and sequenced “in bulk”. This increases the total amount of DNA which improves coverage and reduces amplification errors. To establish the presence or absence of mutations, a variant allele frequency (VAF) is calculated and compared to a threshold, typically 10 – 20% (Strom 2016). This filters out rare mutations present only in a few reads that are likely to be false positives or sequencing errors (Petrackova et al. 2019). More recently, bulk sequencing is used to study cancer evolution using phylogenetic methods, either by comparing VAF (Zhai et al. 2017; Zhao et al. 2016; Ling et al. 2015) or estimating copy number variants (CNV) (Desper et al. 1999; Demeulemeester et al. 2016; Tarabichi et al. 2021). However, the usage of bulk sequencing in this context is problematic. Bulk samples contain cells from multiple cell lineages including non-tumor cells, such as immune or blood vessel cells (Racle et al. 2017), and there is strong evidence for a constant migration of metastatic cells between tumors (Aguirre-Ghiso 2010; Cheung et al. 2016; Reiter et al. 2017; Casasent et al. 2018). High VAF thresholds ignore tumor heterogeneity, but by lowering the threshold, mutations in non-tumor cells or clonal lineages are retained instead. Sequences or mutational profiles derived from bulk samples thus have a chimeric origin (Alves et al. 2017).

A typical assumption in classical phylogenetics is that the sequences or mutational profiles represent individual taxonomic units, either individuals or populations of closely related individuals. If these methods are used on the data from bulk samples, the reconstructed trees are not phylogenies describing an evolutionary history, but evolutionarily meaningless sample similarity trees (Alves et al. 2017). To address this issue, phylogenetic trees are reconstructed by estimating the sequential order of somatic mutations using VAF from one or multiple tumor samples (Deshwar et al. 2015; El-Kebir et al. 2018; Miura et al. 2018). Given the tumor heterogeneity and insufficient read depth to reliably estimate VAF, this is not a simple problem and the performance of current methods is limited (Miura et al. 2020).

Single-cell DNA sequencing (scDNA-seq) does not suffer from the chimeric DNA origin of bulk-sequencing as each DNA segment is barcoded to guarantee its known cell of origin. Recent progress in WGS technology made sequencing individual cells cost-efficient (Gawad et al. 2016) and this approach is now regularly used for the phylogenetic reconstruction of metastatic cancer or the subclonal structure of a single tumor (Potter et al. 2013; Roth et al. 2016; Leung et al. 2017; Myers et al. 2019). However, this increased resolution comes with additional complications. Current methods are not sensitive enough to sequence DNA from a single cell and DNA amplification is required (Gawad et al. 2016). This process suffers from a random bias with different parts of the genome amplified in different quantities or not at all (Satas et al. 2018). In addition, polymerase does not replicate DNA without error, this can have a significant impact if the replication errors occur early in the amplification process (Gawad et al. 2016). This does not only increase the error rate for identified SNVs, but a large proportion of SNVs might be simply missing (Hicks et al. 2018). The advantages associated with scDNA-seq led to the development of novel approaches that tackle

these challenges using an error model to correct for amplification errors and false-positive SNV calls (Zafar et al. 2016; Zafar et al. 2018; Luquette et al. 2019; Kozlov et al. 2020).

Similar technological development led to proliferation of single-cell RNA sequencing (scRNA-seq) which, compared to traditional bulk RNA sequencing, enabled detection of gene expression profiles for individual cells in the tissue sample (Müller et al. 2017b; Olsen et al. 2018; Jerby-Arnon et al. 2018; González-Silva et al. 2020). This allows understanding tumor heterogeneity by identifying different cell populations (Andrews et al. 2018), estimating immune cell content within a tumor (Yu et al. 2019), or even identifying individual clones and subclones, as they can differ in their behavior (Fan et al. 2020). However, as the levels of RNA expression vary between genes and cells, the amplification problems of scDNA-seq that cause unequal expression and drop-out effects are more pronounced in scRNA-seq. There is an increased interest for SNV calling on scRNA-seq data using bulk-SNV callers (Chen et al. 2016; Poirion et al. 2018; Liu et al. 2019; Schnepf et al. 2019) and specialized CNV callers (Kuipers et al. 2020; Harmanci et al. 2020b; Harmanci et al. 2020a; Gao et al. 2021) as this allows for identification of mutations in actively expressed genes.

In this work, we test if expression values and SNVs inferred from scRNA-seq contain phylogenetic information to reconstruct a population history of cancer. We perform an experiment to guarantee a known population history, and then try to reconstruct this history using computational phylogenetics from both expression values and identified SNVs derived from the same scRNA-seq data set. We then compare phylogenies obtained from these methods against the known population history to evaluate the strength of the phylogenetic signal contained in the scRNA-seq data sets.

METHODS

Experimental design. To guarantee known population history, immunosuppressed mice were injected with human breast cancer cells. The tumors that develop are derived from the same population and thus share a common ancestor, but evolved independently in each mouse and should form separate clades on reconstructed phylogenetic trees when analyzed together. As each tumor was seeded by a population of cancer cells, a number of small sample-specific clades representing subclonal diversity of the population sample should be observed. We would expect clustering of each tumor and CTC sample as well as clustering of tumor and CTC samples isolated from a single individual. To test for the presence of these sample-specific clades, as well as the strength of the phylogenetic relationship between cells from each tumor, we employ phylogenetic clustering tests. If the phylogenetic tests confirm sample-specific clustering of cells, then the scRNA-seq data contains sufficient phylogenetic signal. Due to the lack of a specialized scRNA-seq caller or error model to account for the uncertainty in the data, some intermixing is possible, but heavy intermixing would demonstrate an insufficiency of scRNA-seq for phylogenetic analyses.

Sample preparation and scRNA sequencing. MDA-MB-231-LM2 (GFP+) (Minn et al. 2005) cells were injected into the R4 mammary fat pad of Nu/J mice (250,000 cells per mouse, 3 mice), and tumor growth was monitored for 8 weeks. Mice were euthanized when tumor size approached the endpoint (2 cm). Tumors were resected and dissociated into single cells. To extract circulating tumor cells (CTC), up to 1 ml of blood was drawn immediately post euthanasia using cardiac puncture. Red blood cells were removed using RBC lysis buffer. All cells (tumor derived and circulating tumor cells) were stained with DAPI and sorted for DAPI and GFP using a BD FACSAria cell sorter. Libraries were generated using the 10x Chromium single cell gene expression system immediately after cell sorting, and sequenced on an Illumina NextSeq platform together to eliminate batch effect.

Mapping and expression analysis. Reads were mapped with the Cellranger v5.0 software to the GRCh38 v15 from the Genome Reference Consortium using the analysis-ready assembly without alternative locus scaffolds (no_alt_analysis_set) and associated GTF annotation file.

The Cellranger software performs mapping, demultiplexing, cell detection, and gene quantification for the 10x Genomics scRNA-seq data.

Postprocessing expression data.

Standardizing expression values. The filtered feature-barcode expression values from Cellranger were processed using the R Seurat v4.0.4 package (Stuart et al. 2019) and according to the Seurat’s standard pre-processing workflow. However, low-quality cells, such as cells with small number of unique reads or small number of represented genes, were not removed at this stage and no normalization was performed. The expression values for each gene were centered ($\mu = 0$) and rescaled ($\sigma^2 = 1$).

Discretizing expression values. The rescaled expression values were then categorized into a 5 level ordinal scale ranging from 1 (low level of expression) to 5 (high level of expression). The five-level scale was chosen to capture the data distribution of the rescaled expression values and represent a compromise between introducing data noise with too many levels or artificial similarity with only a few categories.

Interval ranges, according to which the values were categorized, were chosen according to the 60% and 90% Highest Density Intervals (HDI), the shortest intervals containing 60% or 90% of values respectively. The values inside the 60% HDI were categorized as normal, values inside the 90% HDI, but outside the 60% as increased/decreased expression and values outside the 90% HDI as a extremely increased/decreased expression.

Genes that contain only a single categorized value for each cell were removed as phylogenetically irrelevant and the discretized values were then transformed into fasta format.

Recording unexpressed genes as unknown data. The amount of coverage in a standard bulk RNA-seq expression analysis is usually sufficient to conclude that genes for which no molecule was detected are not expressed (Lähnemann et al. 2020). In scRNA-seq however, the sequencing coverage is very small, drop out effect is likely, and thus this assumption does not hold. This is especially a problem for non-UMI based technologies (Cao et al. 2021), but not entirely absent from the UMI-based technologies as well due to biological and technological processes (Townes et al. 2020; Hsiao et al. 2020).

According to the standard expression pipeline, these values are commonly treated as biological zeros, i.e., no detected expression of a particular gene, and have a significant impact on the data distribution during the normalization and rescaling steps (Hicks et al. 2018; Townes et al. 2020). Without an explicit model of drop out effect to account for technical or biological variation, these values might be more accurately represented as unknown values rather than true biological zeros (Van den Berge et al. 2018). We have modified the Seurat code to treat these values as unknown values (`NA` in R) and included modified functions in the `phyloRNA` package.

We will further use *data density* to describe the number of unknown values in both expression and SNV datasets, with 100% representing data set without unknown values, while 0% would represent a dataset formed entirely of unknown values.

SNV.

Pre-processing reads for SNV detection. The BAM files from Cellranger were processed using the Broad Institute’s Genome Analysis ToolKit (GATK) v4.2.3.0 (Poplin et al. 2018) according to GATK best practices of somatic short variant discovery.

SNV detection and filtering. To obtain SNVs for individual cells of the scRNA-seq data, first a list of SNVs were obtained by running Mutect2 (Benjamin et al. 2019), treating the data set as a pseudo-bulk sample, and retaining only the SNVs that passed all filters. Mutect2 was run in the tumor with matched normal sample using the parental cell lineage MDA-MB-231 from Kidwell et al. (2021) and Panel of Normals derived from the same source, see supplementary materials for details.

SNVs for individual cells were then obtained by individually summarizing reads belonging to each single cell at the positions of the SNVs obtained beforehand using the pysam library, which is built

on htlib (Li et al. 2009). The most common base for every cell and every position was retained, base heterogeneity or CNVs was ignored. This SNV table was then transformed into fasta format.

Finding a well-represented subset of data. When treating the potentially unexpressed genes as unknown values, only a small proportion of the expression count values was known, with the data set derived from SNV suffering from the same problem due to the low number of reads for each cell.

While model-based phylogenetic methods can process missing data by treating the missing data as phylogenetically neutral, this significantly flattens the likelihood space which can cause artifacts, convergence problems or increase computational time (Wiens 2006; Jiang et al. 2014; Xi et al. 2016).

Published phylogenetic tools designed for single-cell DNA data sets rang from 47 cells and 40 SNVs (Jahn et al. 2016) to 370 cells and 50 SNVs (Singer et al. 2018) or in an extreme case 18 cells and 50,000 SNVs (Singer et al. 2018) with at most 58% of missing data across these data sets. In comparison, scRNA-seq technology can produce up to tens of thousand of cells with tens of thousand detected genes (Chen et al. 2019) and data reduction is often required.

To alleviate these issues, we employ two different filtering strategies to reduce the size of the datasets, while preserving as much information as possible, a selection strategy, where a set of high-quality cells is selected, and a stepwise filtration algorithm, where a subset of data with the highest data density is selected. Under the selection filtering strategy, a set of cells is selected, either cells of interest from the expression analysis, or a fixed subset of cells with the highest data density. This allows for a construction of datasets of specific size.

The stepwise filtering algorithm aims to find a well-represented subset of the data. By iteratively cutting out cells and genes/SNVs with the smallest number of known values, we increase the data density until a local maximum or desired data density is reached. This is equivalent to the gene/cell quality filtering during the scRNA-seq post-processing pipeline, such as Seurat's standard pre-processing workflow described above, where low-quality cells and genes are removed. The advantage of this method is that a desired density can be reached with the least amount of data removed.

Phylogenetic analysis. To reconstruct phylogenetic trees from the categorized expression values and identified SNVs, we used IQ-TREE v2.1.4 (Minh et al. 2020) and BEAST2 v2.6.3 (Bouckaert et al. 2019).

The IQ-TREE analysis was performed with an ordinal model and an ascertainment bias correction (`-m ORDINAL+ASC`) for the expression data, and a standard model selection was performed for the SNV data (`-m TEST`). Where the size of the dataset allowed, tree support was evaluated using the standard non-parametric bootstrap (Felsenstein 1985) with 100 replicates (`-b 100`).

The BEAST2 analysis was performed with a birth-death tree prior (Kingman 1982) with an exponential population growth (Kuhner et al. 1998), as these models most closely mimic the biological conditions of tumor growth. For the expression data, the BEAST2 was run using ordinal model available in the Morph-Models package, while the SNV values were analyzed using the Generalized Time-Reversible model (Tavaré 1986). For both the expression and the SNV data sets, BEAST2 was set to not ignore ambiguous states.

Phylogenetic clustering tests. To test if the phylogenetic methods were able to recover expected population history, we employ Mean Pairwise Distance (MPD) (Webb 2000) and Mean Nearest Taxon Distance (MNTD) (Webb 2000). MPD is calculated as a mean distance between each pair of taxa from the same group, while MNTD is calculated as a mean distance to the nearest taxon from the same group. For each sample and samples isolated from a single individual, MPD and MNTD are calculated and compared to a null distribution obtained by permuting sample labels on a tree and calculating MPD and MNTD for these permutations. The p-value is then calculated as a rank of the observed MPD/MNTD in the null distribution normalized by the number of permutations. The MPD and MNTD are calculated using the `ses.mpd` and `ses.mntd` functions implemented in the package `picante` (Kembel et al. 2010) For the Bayesian phylogenies, MPD and MNTD were

calculated for a sample of 1000 trees from the posterior distribution and then summarized with mean and 95% confidence interval.

Application to other datasets. To further evaluate our approach, we analyze two previously published datasets, a UMI based dataset of small intestinal neuroendocrine cancer (Rao et al. 2020a) and non-UMI based dataset of gastric cancer (Wang et al. 2021). These datasets contain primary and metastatic cells from two regional samples, allowing us to assess the performance of the phylogenetic analysis using the phylogenetic clustering tests. We assume that primary, metastatic, and cells from regional samples will each cluster together, forming a cell-type and region-specific clades.

Intestinal neuroendocrine cancer. The small intestinal neuroendocrine cancer dataset from Rao et al. (2020a) consisted of a primary tumor and a paired liver metastatic sample. Both samples contained a mixture of cancerous and non-cancerous cells (Fibroblasts, Endothelial cells, and Immune cells). The expression values for both samples from Rao et al. (2020a) were processed as per the methodology section, with zeros recoded as unknown data. To obtain the SNVs, the raw reads were mapped using the Cellranger v5 and processed as per the methodology section. We have used Mutect2 in a tumor-only mode using the Panel of Normals and the GNOMAD germline data from the GATK best practices resource bundle. Cells were labeled according to their sample of origin (primary tumor and metastasis) and their cell type, which was determined by replicating the analysis from Rao et al. (2020a). Two subsets for both data types were then derived, a subset with all cell types and a subset with only cancer cells. To reduce the computational burden, 1000 cells with the least amount of missing data were selected, 500 from the primary tumor and 500 from the metastatic sample. To derive subsets from the SNVs, the cells from the expression subsets were used. Maximum Likelihood trees were then reconstructed and the relationship between cells of the same type and sample of origin were then tested using the phylogenetic clustering tests.

Gastric cancer. The gastric cancer dataset from Wang et al. (2021) consisted of 94 cells from a primary tumor and a lymph node of three patients (GC1, GC2, and GC3). We would expect that for each patient, the lymph node cells would form a monophyletic lineage derived from the primary tumor cell, but due to the small number of cells, clustering of the primary tumor cells is also interpreted as a success. The expression values were split into patient-specific datasets and analyzed separately as per the methodology section and the discretized expression values were analyzed using the Maximum Likelihood and the Bayesian phylogenetic methods. To obtain the SNV values, the raw reads were mapped using the STAR v2.7.9a (Dobin et al. 2013) and mapped reads were then processed as per the methodology section. As with the Intestinal neuroendocrine cancer, Mutect2 was used in a tumor-only mode using the Panel of Normals and the GNOMAD germline data from the GATK best practices resource bundle. Bayesian and Maximum Likelihood trees were constructed and the clustering of primary and lymph node cells was then explored using the phylogenetic clustering tests.

Code and data availability. Code required to replicate the data processing steps is available at <https://github.com/bioDS/phyloRNAanalysis>.

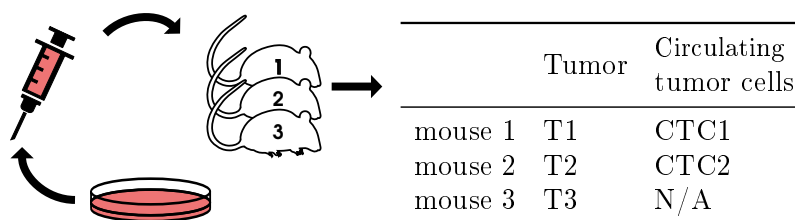
To aid in creating pipelines for phylogenetic analysis of scRNA-seq data, we have integrated a number of common tools in the R `phyloRNA` package, which is available at <https://github.com/bioDS/phyloRNA>.

All the data is available in the NCBI GEO under the accession number GSE163210.

RESULTS

Sample overview. In total, five samples were used in this analysis, three tumor samples (T1, T2, T3) and two CTC samples (CTC1, CTC2). The number of cells isolated from the CTC3 sample

TABLE 1. An overview of data set used in this work. In total, five samples were isolated from three individuals (Table 1a): 3 tumor samples (T1, T2, T3) and 2 circulating tumor cell samples (CTC1, CTC2). For each sample, the number of cells from fluorescent-activated cell sorting (FACS), the number of identified by Cellranger, the number of detected genes, the number of unique molecular identifiers (UMIs), UMI/Cell ratio, and the data density are reported (Table 1b).



(a) Experiment overview

Sample	Cells (FACS)	Cells (Cellranger)	Genes	UMI	UMI/Cell	data density
T1	11,258	701	17k	3,167k	4,518	2.97 %
T2	20,233	2,794	5k	69k	25	0.04 %
T3	13,865	806	18k	2,876k	3,569	2.57 %
CTC1	605	3,125	8k	129k	41	0.06 %
CTC2	415	3,161	9k	155k	49	0.06 %
total:	46,376	10,587	20k	6.4M	604	0.44 %

(b) Sample overview

was too small for scRNA sequencing and the sample was removed from the study. The number of detected cells in the tumor samples was generally smaller than in the CTC samples, but the reverse was true for the total number of detected unique molecular identifiers (UMIs) – the number of unique mRNA transcripts (see Table 1). In the T2 sample, a large number of cells but a small number of UMIs were detected in a similar pattern to the CTC samples.

Compared to the fluorescent-activated cell sorting (FACS), Cellranger detected fewer cells for tumor samples, but more cells for the CTC samples. Cellranger classifies barcodes as cells based on the amount of UMI detected to distinguish real cells from a background noise (Lun et al. 2019). The large number of detected cells in the CTC samples is likely a result of lysed cells or cell-free RNA (Fleming et al. 2019). In all cases, the number of expression values across data sets was relatively low, with the best sample T3 amounting to about 3% of known expression values.

SNV identification. To identify SNVs in scRNA-seq data, we first identified a list of SNVs by treating the single-cell reads as a pseudo-bulk sample. The total of 21,261 SNVs that passed all quality filters were identified this way. When these SNVs were called for each individual cell, the resulting data set had data density of less than 0.13%. The expression data is expected to have higher data density than SNV because for expression quantification a presence or absence of a molecule is sufficient while for SNV, knowledge of each position is required. This expectation is confirmed in Table 1, where data density of the expression data is summarized. About 40% of the 10,587 cells represented in this data set did not contain any positively identified SNV after filtering out false-positives, these were relatively equally distributed among the T2 (1487), CTC1 (1379) and CTC2 (1324) samples. This represents a challenge from a data analysis perspective given the large sample size and its small data density.

Data reduction. With over 10,000 cells and more than 20,000 genes and SNVs, the unfiltered datasets would require substantial computational resources. An additional issue we have encountered in our data was a significant difference in the quality between individual samples, only five CTC1 and six CTC2 cells passed the quality filtering criteria of a minimum of 250 represented genes and a minimum of 500 UMI per cell, with no T2 cells passing the quality filtering. This contrasts with the T1 and T3 samples, where 701 and 806 passed the quality filtering criteria respectively. Due to this varied quality of samples, filtering data to a higher data density using the stepwise filtering algorithm leads to the removal of the low-quality samples (T1, CTC1 and CTC2), which bar us from testing the phylogenetic structure using the phylogenetic clustering tests. For this reason, we have selected a small number of cells with the least amount of missing data from each sample using the selection filtering method. The small number of cells is not sufficient to represent the full diversity of the tumor, but allow us to test the phylogenetic relationship between individual samples without introducing a bias due to an unequal size of the samples.

A total of 58 cells were retained for both the expression and SNV datasets: 20 cells for T1 and T3 samples and six cells for T2, CTC1 and CTC2 samples. In these reduced datasets, genes that were not present in any of the cells or present only in a single cell, are removed. The reduced expression data set contained 30% of known data distributed across 7,520 genes. The SNV data set contained 10% of known data distributed across 1,058 SNVs. These reduced data sets are analyzed using Maximum Likelihood and Bayesian method to further explore the topological uncertainty.

Reconstructed trees and phylogenetic tests for the data filtered to the 20%, 50% and 90% data density using the stepwise filtering algorithm are provided in the supplementary materials.

Phylogenetic reconstruction from expression data. The Maximum Likelihood tree reconstructed from the reduced expression data set showed significant clustering of all samples (Figure 1a). This is confirmed by the phylogenetic clustering tests where all but CTC2 cells had a significant MPD p-value (Table 2). Four out of six CTC2 cells clustered together, but on the opposite side of the tree with phylogenetic proximity to the T1 cells. This close phylogenetic relationship suggests that T1 and CTC2 were isolated from a single individual. This pattern is further reinforced as T2 cells clustered in a single compact clade with phylogenetic proximity to the CTC1 sample. When this relationship was tested with phylogenetic clustering methods, both MPD and MNTD confirmed the strong clustering signal between T2 and CTC1. The same tests were not significant for the T1-CTC2 grouping likely due to the presence of two non-clustering CTC2 cells.

The phylogenies reconstructed from the same data using the Bayesian inference show a similar pattern of clustering (Figure 1b, Table 2), although neither CTC1 nor CTC2 formed a compact cluster. The T2 and CTC1 connection is not supported, but about half of the CTC1 cells were placed in a group with the T2 samples. Similarly to the Maximum Likelihood tree, this group was not closely related to the T1 and T3 cells, instead it formed a distantly related sister group. The relationship between T1 and CTC2 is supported by the MNTD statistics on the Bayesian phylogeny.

Neither MNTD nor MPD statistics on the Maximum Likelihood and Bayesian phylogeny supported the clustering of CTC2 cells. This might suggest that the CTC2 cells are polyphyletic, with their origin in the seeding population before the injection. This is not unlikely given that the cell lineage used (MDA-MB-231-LM2) is highly metastatic (Minn et al. 2005).

In addition to testing on the best phylogeny, we have integrated the topological uncertainty of the reconstructed phylogenies by performing the phylogenetic clustering tests on the 100 bootstrap replicates from the Maximum Likelihood analysis and a sample of 1000 trees from the Bayesian posterior tree sample. The distribution of MPD and MNTD p-values calculated on each tree were then summarized using mean and 95% confidence interval. The majority of relationships from the best tree were also supported by the tree samples (Supplementary Table 1). This suggests that while there is high uncertainty in the data and reconstructed topologies, we can reconstruct broad topological patterns with relatively high certainty.

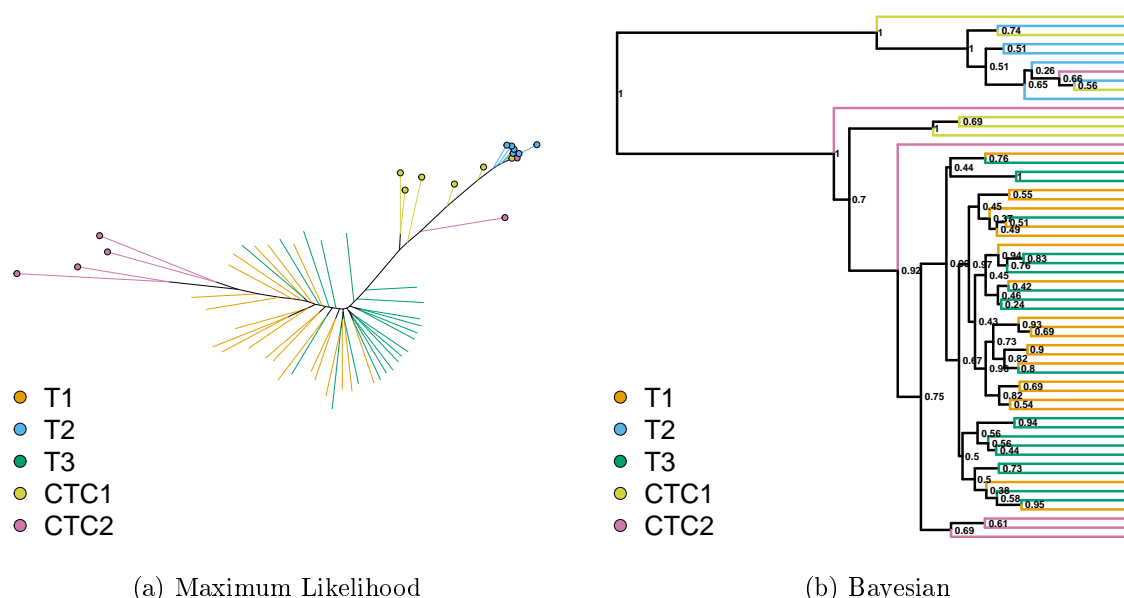


FIGURE 1. Maximum Likelihood and Bayesian trees reconstructed from the expression data for the 58 selected cells. Terminal branches are colored according to cell's sample of origin (T1, T2, T3, CTC1, CTC2). In the Maximum Likelihood tree, the T2, CTC1 and CTC2 samples are also marked with colored circles. For the Bayesian tree, Bayesian posterior values show the topology uncertainty.

TABLE 2. Test of phylogenetic clustering for the reduced dataset of the 58 selected cells. Mean Pairwise Distance (MPD) and Mean Nearest Taxon Distance (MNTD) calculated for the Maximum Likelihood (ML) and Bayesian (BI) trees from the expression and SNV data. P-values for MPD and MNTD were calculated for each sample (T1, T2, T3, CTC1, CTC2) and expected clustering for cells isolated from a single individual (T1 with CTC1, and T2 with CTC2) and to test a possible mislabeling between CTC1 and CTC2 samples (T1 with CTC2, and T2 with CTC1). Significant p-values at $\alpha = 0.05$ after correcting for multiple comparisons using the False Discovery Rate method (Benjamini et al. 1995) are marked with an asterisk. Values of MNTD and MPD calculated for the Maximum Likelihood bootstrap sample and Bayesian posterior tree sample are available in the Supplementary Table 1.

Groups	Cells	Expression (ML)		Expression (BI)		SNV (ML)		SNV (BI)	
		MPD	MNTD	MPD	MNTD	MPD	MNTD	MPD	MNTD
T1	20	*0.003	0.755	*0.001	*0.006	*0.001	*0.008	*0.001	0.434
T2	6	*0.001	*0.001	*0.001	*0.003	*0.001	*0.001	*0.001	*0.001
T3	20	*0.001	0.423	*0.001	*0.017	*0.002	0.660	*0.001	0.184
CTC1	6	*0.002	*0.004	0.992	0.479	*0.001	*0.001	*0.001	*0.002
CTC2	6	1.000	0.988	0.724	0.966	0.247	0.595	0.223	0.594
T1 & CTC1	26	0.059	0.407	0.139	0.261	0.475	*0.008	0.466	0.305
T2 & CTC2	12	0.997	0.027	0.970	0.165	0.998	0.800	0.977	0.317
T1 & CTC2	26	0.637	0.999	*0.005	0.932	*0.001	0.630	*0.001	0.943
T2 & CTC1	12	*0.001	*0.001	0.674	0.034	*0.001	*0.001	*0.001	*0.001

* significant support

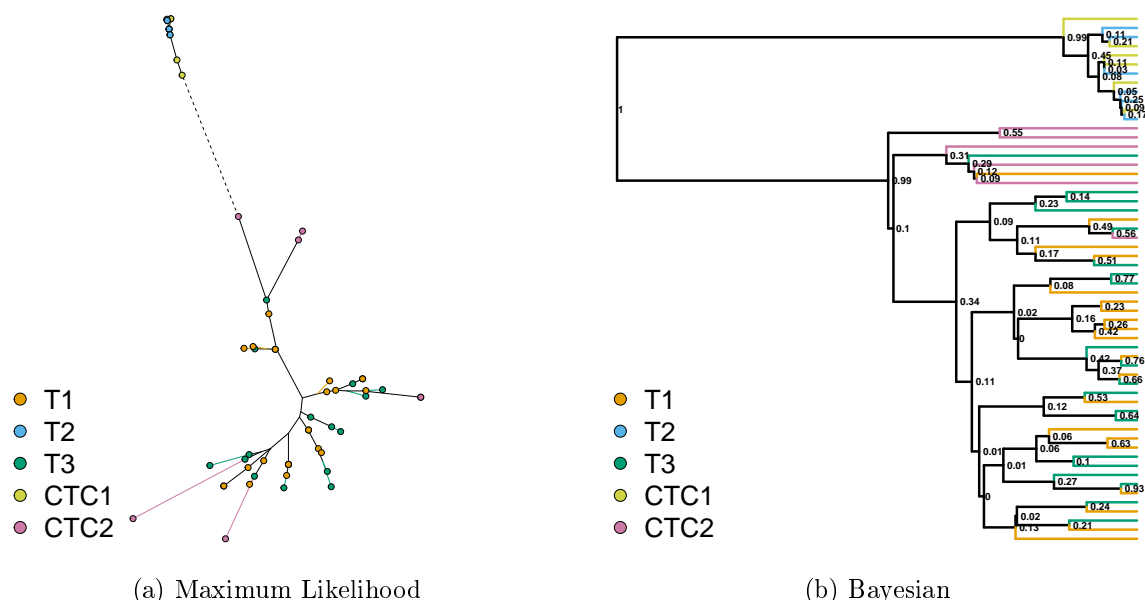


FIGURE 2. Maximum Likelihood and Bayesian trees reconstructed from the SNV data for the 58 selected cells. Terminal branches are colored according to cell's sample of origin (T1, T2, T3, CTC1, CTC2). In the Maximum Likelihood tree, cells are also marked with colored circles. For the Bayesian tree, Bayesian posterior values show the topology uncertainty.

Phylogenetic reconstruction from the SNV data. The Maximum Likelihood tree reconstructed from the reduced SNV dataset (Figure 2) displayed similar but weaker patterns to the one reconstructed from the expression data. The CTC2 cells no longer formed two compact clusters and were dispersed along the tree. Similarly to the expression data, the T2 and CTC1 cells were placed together on a long branch suggesting a long shared evolutionary history. However, unlike the expression data, the T1 and T3 were more interspersed with very short branches. The phylogenetic clustering tests confirm the grouping of all samples (Table 2), except for the CTC2 sample, in addition to the putative relationship between T1 and CTC2, and T2 and CTC1 samples. This reinforces the hypothesis about possible mislabeling between CTC1 and CTC2 samples.

A similar pattern of sample clustering can be observed on the Bayesian phylogeny reconstructed from the same data (Figure 2b), with T2 and CTC1 cells placed on a distantly related sister branch to all other samples. The T1 and T3 cells are still interspersed, but the CTC2 cells seem to cluster together more closely. Like with the expression analysis, when these relationships are stable when the topological uncertainty is integrated into the phylogenetic clustering tests (Supplementary Table 1).

Biological zero or unknown value. To test the assumption if the zero expression values should be treated as unknown data rather than biological zeros, i.e., no expression of a particular gene, we have reconstructed the phylogenies from the scRNA-seq expression by treating the zeros in the dataset as biological zeros. Data were processed as per the standard methodology to get the alignments, but instead of treating the zeros as an unknown position, they were treated as a category 0 in addition to the 5 level ordinal scale. Phylogenies were then reconstructed using both Maximum Likelihood and Bayesian methods with sample clustering explored using the phylogenetic clustering tests.

In the phylogenies reconstructed from the expression data when zero is treated as a biological zero (Figure 3), the CTC2 cells did not form a cluster but clustered closely with the T1 and CTC2 cluster. This cluster was no longer placed as a sister branch to the T1 and T3 cells but was deeply

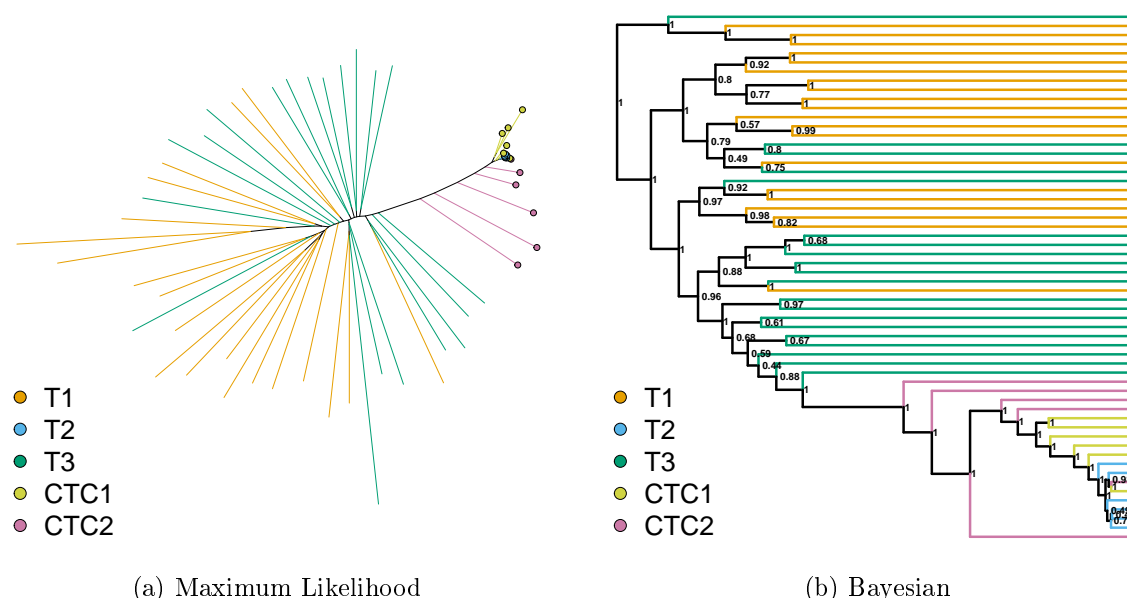


FIGURE 3. Maximum Likelihood and Bayesian trees reconstructed from the expression data for the 58 selected cells. Terminal branches are colored according to cell's sample of origin (T1, T2, T3, CTC1, CTC2). In the Maximum Likelihood tree, the T2, CTC1 and CTC2 samples are also marked with colored circles. For the Bayesian tree, Bayesian posterior values show the topology uncertainty.

TABLE 3. Test of phylogenetic clustering for the expression data when zero expression level is treated as biological zero. Mean Pairwise Distance (MPD) and Mean Nearest Taxon Distance (MNTD) calculated for the Maximum Likelihood (ML) and Bayesian (BI) trees from the expression data, with zeros treated as biological zeros. P-values for MPD and MNTD were calculated for each sample (T1, T2, T3, CTC1, CTC2) and expected clustering for cells isolated from a single individual (T1 with CTC1, and T2 with CTC2) and to test a possible mislabeling between CTC1 and CTC2 samples (T1 with CTC2, and T2 with CTC1). Significant p-values at $\alpha = 0.05$ after correcting for multiple comparisons using the False Discovery Rate method (Benjamini et al. 1995) are marked with an asterisk.

Groups	Cells	Expression (ML)		Expression (BI)	
		MPD	MNTD	MPD	MNTD
T1	20	1.000	1.000	0.998	0.992
T2	6	*0.001	*0.001	*0.001	*0.001
T3	20	0.543	0.994	0.804	0.998
CTC1	6	*0.001	*0.001	*0.001	*0.001
CTC2	6	*0.004	*0.013	*0.001	*0.011
T1 & CTC1	26	0.992	0.928	0.968	0.560
T2 & CTC2	12	*0.001	*0.001	*0.001	*0.001
T1 & CTC2	26	0.997	0.994	0.991	0.888
T2 & CTC1	12	*0.001	*0.001	*0.001	*0.001

* significant support

nested. The T1 and T3 samples were less interspersed than when zero is treated as unknown data. This change in the phylogenetic structure is supported by the phylogenetic clustering tests, with T1 and T3 no longer being supported and instead, the clustering of CTC2 cells is being supported in both the Maximum Likelihood and Bayesian phylogenies. Likewise, the T1 and CTC2 grouping is not supported, as the CTC2 cells group together with the CTC1 and T2 samples.

These results do not provide a conclusive answer on which assumption should be preferred. Assuming all zeros to be biological zeros will bias the model as many of those might be technical zeros instead. At the same time, the pattern of expression and non-expression seems to carry information. This information is lost when all zeros are assumed to be technical zeros and thus unknown data. For our datasets, the assumption of zeros as technical zeros and thus unknown data seems to create better agreement in the phylogenetic structure between the expression and SNVs and thus should be preferred. However, our datasets also suffered from unequal data quality issues (Table 1), and under different conditions, assuming zeros as biological zeros might be preferred.

Application to other datasets.

Intestinal neuroendocrine cancer. We have derived two subsets from the expression and SNV data for the small intestinal neuroendocrine cancer dataset from Rao et al. (2020a), a subset with all cell types and a subset with cancer cells only. However, not all cells found in the expression subsets were found in the SNV data. This is likely due to a different version of the Cellranger software used in this work compared to the Rao et al. (2020a). In both derived subsets from the expression data, metastatic cells showed a strong clustering tendency ($p = 0.001$) into several large clades (Figure 4). This suggests a strong phylogenetic relationship with several well-preserved lineages. In addition, in the derived subset containing all cell types, the cancer cells showed a significant clustering ($p = 0.001$), while other cell types showed the opposite tendency (Figure 4). However, the cancer clade contained deeply nested clades of Endothelial cells and Immune cells. A similar albeit significantly weaker pattern of cancer cell clustering can be observed on the trees derived from the SNV data (Figure 4, Figure 4). In both subsets derived from the SNV data, the primary cells clustered together, but the pattern was less consistent and confirmed only by one of the two tested statistics.

Gastric cancer. For both the expression and SNV data from the gastric cancer dataset published by Wang et al. (2021), only a single patient showed significant clustering of lymph nodes (Figure 5). Poor separation of primary and lymph node cells from the expression levels was pointed out in the original study (Wang et al. 2021). Additionally, non-UMI based methods suffer from an increased error rate through zero-count inflation (Cao et al. 2021) and amplification variability (Townes et al. 2020). In the absence of a strong phylogenetic signal shared by a large percentage of genes, this additional noise is making a phylogenetic reconstruction difficult, if not impossible. At the same time, the typically higher coverage in the non-UMI based sequencing compared to the UMI should improve the identification of SNVs and decrease the misspecification error. This might suggest that different strategies for the phylogenetic reconstruction should be applied to UMI and non-UMI based sequencing.

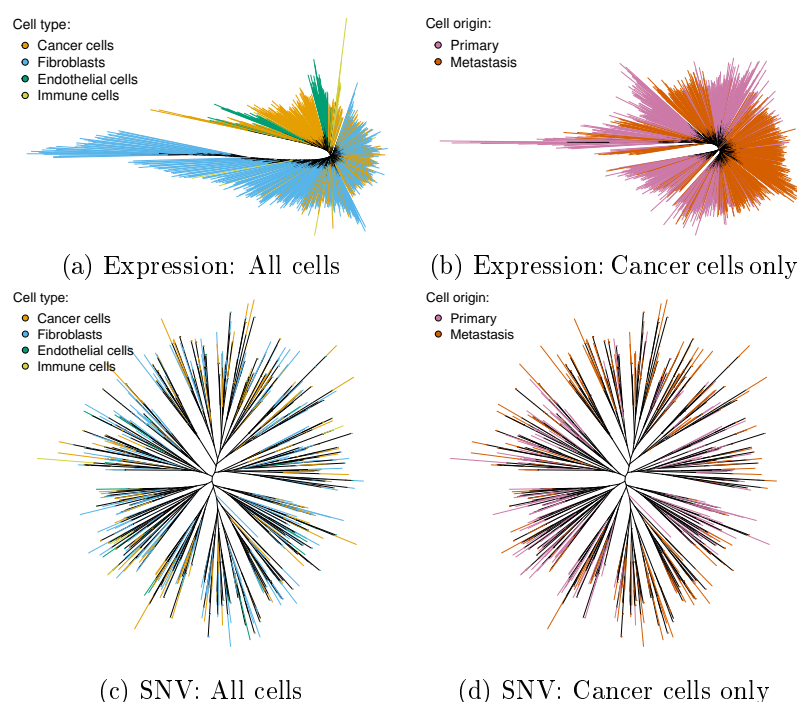


FIGURE 4. Maximum Likelihood trees constructed from the expression and SNV data published by Rao et al. (2020a). Terminal branches are colored according to cell's type or sample of origin. In the tree reconstructed from expression data for all cells (Figure 4a), the vast majority of cancer cells cluster in a single clade. The tree reconstructed from expression data for cancer cells only (Figure 4b) shows a strong clustering of primary and metastatic cells. While the metastatic cells are not clustered in a single clade, multiple metastatic events are biologically plausible. In the trees reconstructed from the SNV data (Figure 4c, Figure 4d), primary and metastatic cells, as well as cells of different type, are relatively evenly distributed without any apparent clustering.

TABLE 4. Test of phylogenetic clustering on the Maximum Likelihood trees from Rao et al. (2020a). Mean Pairwise Distance (MPD) and Mean Nearest Taxon Distance (MNTD) calculated for the phylogeny reconstructed from the dataset containing only cancer cells and from the dataset containing all cell types. P-values for MPD and MNTD were calculated for the sample of origin and cell types where applicable. Significant p-values at $\alpha = 0.05$ after correcting for multiple comparisons using the False Discovery Rate method (Benjamini et al. 1995) are marked with an asterisk.

Data	Groups	Cancer only			All cell types		
		Cells	MPD	MNTD	Cells	MPD	MNTD
Expression	Cancer cells	1000	—	—	355	*0.001	*0.001
	Fibroblasts	0	—	—	552	1.000	0.989
	Endothelial cells	0	—	—	71	0.860	0.903
	Immune cells	0	—	—	22	0.651	0.565
	Metastasis	500	*0.001	*0.001	500	*0.001	*0.001
	Primary	500	1.000	1.000	500	1.000	1.000
SNV	Cancer cells	981	—	—	355	0.362	*0.004
	Fibroblasts	0	—	—	552	0.402	0.997
	Endothelial cells	0	—	—	71	0.596	0.785
	Immune cells	0	—	—	22	0.949	0.968
	Metastasis	500	0.907	0.132	500	1.000	*0.001
	Primary	481	*0.004	0.076	500	*0.001	0.095

* significant support

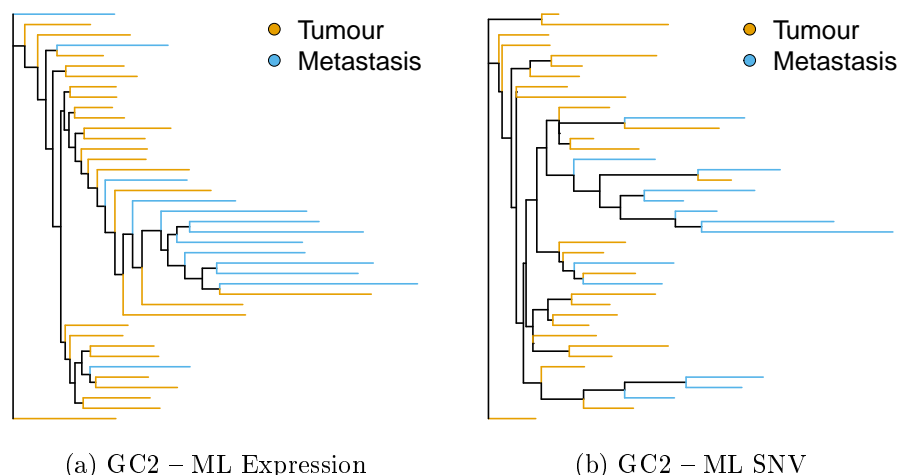


FIGURE 5. Maximum Likelihood trees for the patient G2 constructed from the expression and SNV data published by Wang et al. (2021). Terminal branches are colored according to cell's sample of origin. Only the patient G2 shows a significant clustering signal both on the trees from Expression and SNV data. For all trees, see Supplementary Figure 3 and Supplementary Figure 4

TABLE 5. Test of phylogenetic clustering on the Maximum Likelihood and Bayesian trees calculated from expression and SNV data published by Wang et al. (2021). Mean Pairwise Distance (MPD) and Mean Nearest Taxon Distance (MNTD) calculated for the Maximum Likelihood and Bayesian trees reconstructed from the expression and the SNV data for patients GC1, GC2 and GC2. P-values for MPD and MNTD were calculated for the sample of origin. Significant p-values at $\alpha = 0.05$ after correcting for multiple comparisons using the False Discovery Rate method (Benjamini et al. 1995) are marked with an asterisk.

Data	Type	Groups	GC1			GC2			GC3		
			Cells	MPD	MNTD	Cells	MPD	MNTD	Cells	MPD	MNTD
Expression	ML	Primary tumor	19	0.171	0.123	27	*0.001	*0.001	19	0.829	0.854
		Lymph node	4	0.830	0.869	13	1.000	1.000	12	0.138	0.093
	BI	Primary tumor	19	0.776	0.995	27	0.999	0.953	19	0.218	0.232
		Lymph node	4	0.086	0.035	13	*0.002	*0.001	12	0.205	0.333
SNV	ML	Primary tumor	19	0.102	0.111	27	*0.001	*0.001	19	0.552	0.231
		Lymph node	4	0.507	0.092	13	1.000	0.945	12	0.276	0.208
	BI	Primary tumor	19	0.117	0.025	27	*0.001	*0.014	19	0.531	0.372
		Lymph node	4	0.955	0.935	13	1.000	0.960	12	0.487	0.322

* significant support

DISCUSSION

Phylogenetic methods using scDNA-seq data are becoming increasingly common in tumor evolution studies. scRNA-seq is currently used for studying expression profiles of cancer cells and their behavior. However, while clustering approaches to identify cells with similar expression profiles are common and frequently used, scRNA-seq data are yet to be used in phylogenetic analyses to reconstruct the population history of somatic cells. To test if the scRNA-seq contains a phylogenetic signal to reliably reconstruct the population history of cancer, we have performed an experiment to produce a known history by infecting immunosuppressed mice with human cancer cells derived

from the same population. Then using two different forms of scRNA-seq data, expression values, and SNVs, we reconstructed phylogenies using Maximum Likelihood and Bayesian phylogenetic methods. By comparing the reconstructed trees to the known population history, we confirmed that scRNA-seq contains a phylogenetic signal to reconstruct the population history of cancer, with both the expression values and SNVs producing a similar phylogenetic pattern. However, this signal is burdened by uncertainty in both the source data as well as reconstructed phylogeny. Accurate phylogenies might thus need an explicit error model to account for this increased uncertainty (Hicks et al. 2018). Still, by taking this topological uncertainty into account, we can make a conclusion about the structural relationship of individual cells. This highlights that scRNA-seq can be utilized to explore both the physiological behavior of cancer cells and their population history using a single source of data.

Without any specialized phylogenetic or error models for the scRNA-seq data, conventional methods and software tools developed for systematic biology are able to reconstruct population history from this data, potentially at low computational cost. This implies that more accurate inference will be possible when and if specialized models and software are developed, and serious computational resources are employed. For example, computationally more intensive standard non-parametric bootstrap or Bayesian methods on the unfiltered data sets are certainly within the reach of modern computing clusters. This is a future direction for research.

In this work, we tested for phylogenetic signal on three data sets, a new data set consisting of 5 tumor samples seeded using a population sample, and two previously published data sets consisting of a primary tumor with a paired lymph node or a metastatic samples. Due to the nature of the experiment and the amount of uncertainty in the scRNA-seq data, this barred us from a more detailed exploration of the tree topology as only broad patterns, the phylogenetic clustering of cells according to sample and individual of origin, could be considered. Our clustering analyses show that the phylogenetic trees conform broadly to the expected shapes under different experimental conditions, and thus that expression and SNV data can both be used to infer phylogenetic trees from scRNA-seq. Nevertheless, our results also demonstrate that all such trees contain significant uncertainty, so new datasets and methods will be required to extend this work.

The degree to which low and uneven gene expression plays a role in scRNA-seq requires special attention, especially for non-UMI based data sets, as this causes not only a large proportion of missing data, but also burdens the known values with a significant error rate. Research should aim at trying to quantify this expression-specific error rate and build specialized models to include the uncertainty about the observed data in the phylogenetic reconstruction itself. This could potentially include removing a large proportion of low-coverage data in favor of robust analysis and proper uncertainty estimation of the inferred topology.

The estimation of the topological uncertainty, be it the Bootstrap branch support or the Bayesian posterior clade probabilities, is a staple for phylogenetic analyses. Currently existing methods for the phylogenetic analysis of scDNA-seq, such as SCITE (Jahn et al. 2016), SiFit (Zafar et al. 2017), or SCIΦ (Singer et al. 2018), do not provide this uncertainty estimate. This makes interpretation of the estimated topology difficult because a single topology can only be marginally more accurate than a number of alternative topologies. Of packages we are aware of, only CellPhy, through its integration in the phylogenetic software RAXML-NG (Kozlov et al. 2019), provides an estimate of topological uncertainty through the bootstrap method. Bayesian methods could be a solution as they provide an uncertainty estimate through the posterior distribution. However, they are significantly more computationally intensive than Maximum Likelihood methods. Instead, as the size of single-cell data sets will only increase, bootstrap approximations optimized for a large amount of missing data need to be developed to provide a fast and accurate estimate of topological uncertainty.

An aspect of scRNA-seq expression data that was not considered here is correlated gene expression (Wang et al. 2004; Bageritz et al. 2019). A single somatic mutation could thus induce a change of expression of multiple genes. This might be problematic given that phylogenetic methods assume

that individual sites are independent and this would cause an overestimation of a mutation rate. However, phylogenetic methods are generally rather robust to a wide range of model violations (Huelsenbeck 1995a; Huelsenbeck 1995b; Song et al. 2010; Philippe et al. 2011). In addition, by randomly sampling sites, the bootstrap analysis does explore solutions that would arise from this model violation. An investigation of the effect of correlated gene expression on the estimated phylogeny provides an interesting direction for further research.

Multimic approaches are increasingly popular as they integrate information from multiple biological layers (Bock et al. 2016; Hasin et al. 2017; Nam et al. 2020). While CNVs were ignored in this paper, it is possible to detect large-scale CNVs from scRNA-seq data (Müller et al. 2018; Kuipers et al. 2020; Harmanci et al. 2020b; Harmanci et al. 2020a; Gao et al. 2021). Combined with the SNVs and expression data as analyzed in this paper, this enables a multimic approach using just a single scRNA-seq data source, without the additional cost of DNA sequencing.

ACKNOWLEDGEMENT

We thank Dr. Jon Preall and the Genomics Technology Development Core (CSHL) for scRNA-seq library preparation, and Pamela Moody and the Flow Cytometry Facility (CSHL) for support with single-cell sorting. We acknowledge Suzanne Russo for technical assistance with animal experiments.

AG and JCM acknowledge support from the Royal Society te Apārangi through a Rutherford Discovery Fellowship (RDF-U001702), AG, RL, and SDD acknowledge support of from an Endeavour Smart Ideas grant (U00X1912), AG acknowledges support from a Data Science Programmes grant (UOAX1932), SDD acknowledges support from a Rutherford Discovery Fellowship (RDF-U001802) and the NHI/NCI grant (1K99CA215362-01), and DLS acknowledges support from the NCI grant (5P01CA013106-Project 3). We would also like to acknowledge the CSNL Next-Gen Sequencing Core (NCI-2P30CA45508).

REFERENCES

- [1] Christopher Abbosh, Nicolai J Birkbak, Gareth A Wilson, Mariam Jamal-Hanjani, Tudor Constantin, Raheleh Salari, John Le Quesne, David A Moore, Selvaraju Veeriah, et al. “Phylogenetic ctDNA analysis depicts early-stage lung cancer evolution”. *Nature* 545.7655 (Apr. 2017), pp. 446–451.
- [2] Julio A Aguirre-Ghiso. “On the theory of tumor self-seeding: implications for metastasis progression in humans”. *Breast Cancer Res.* 12.2 (Apr. 2010), p. 304.
- [3] David Aldous. “Probability Distributions on Cladograms”. *Random Discrete Structures*. Springer New York, 1996, pp. 1–18.
- [4] David J Aldous. “Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today”. *Stat. Sci.* 16.1 (Feb. 2001), pp. 23–34.
- [5] João M Alves, Sonia Prado-Lopez, Jose Manuel Cameselle-Teijeiro, and David Posada. “Rapid evolution and biogeographic spread in a colorectal cancer”. May 2019.
- [6] João M Alves, Tamara Prieto, and David Posada. “Multiregional Tumor Trees Are Not Phylogenies”. *Trends Cancer Res.* 3.8 (Aug. 2017), pp. 546–550.
- [7] Tallulah S Andrews and Martin Hemberg. “Identifying cell populations with scRNASeq”. *Mol. Aspects Med.* 59 (Feb. 2018), pp. 114–122.
- [8] Mihaela Angelova, Bernhard Mlecnik, Angela Vasaturo, Gabriela Bindea, Tessa Fredriksen, Lucie Lafontaine, Bénédicte Buttard, Erwan Morgand, Daniela Bruni, et al. “Evolution of Metastases in Space and Time under Immune Selection”. *Cell* (Oct. 2018).
- [9] Josephine Bageritz, Philipp Willnow, Erica Valentini, Svenja Leible, Michael Boutros, and Aurelio A Teleman. “Gene expression atlas of a developing tissue by single cell expression correlation analysis”. *Nat. Methods* 16.8 (Aug. 2019), pp. 750–756.

- [10] Niko Beerenwinkel, Roland F Schwarz, Moritz Gerstung, and Florian Markowetz. “Cancer evolution: mathematical models and computational inference”. *Syst. Biol.* 64.1 (Jan. 2015), e1–25.
- [11] David Benjamin, Takuto Sato, Kristian Cibulskis, Gad Getz, Chip Stewart, and Lee Lichtenstein. “Calling Somatic SNVs and Indels with Mutect2”. *bioRxiv* (2019). DOI: 10.1101/861054. eprint: <https://www.biorxiv.org/content/early/2019/12/02/861054.full.pdf>. URL: <https://www.biorxiv.org/content/early/2019/12/02/861054>.
- [12] Y Benjamini and Y Hochberg. “Controlling the false discovery rate: a practical and powerful approach to multiple testing”. *J. R. Stat. Soc.* (1995).
- [13] Christoph Bock, Matthias Farlik, and Nathan C Sheffield. “Multi-Omics of Single Cells: Strategies and Applications”. *Trends Biotechnol.* 34.8 (Aug. 2016), pp. 605–608.
- [14] Remco Bouckaert, Timothy G Vaughan, Joelle Barido-Sottani, Sebastian Duchene, Mathieu Fourment, Alexandra Gavryushkina, Joseph Heled, Graham Jones, Denise Kuhnert, et al. “BEAST 2.5: An Advanced Software Platform for Bayesian Evolutionary Analysis”. *PLoS Comput. Biol.* 15.4 (Apr. 2019), e1006650.
- [15] Yingying Cao, Simo Kitanovski, Ralf Küppers, and Daniel Hoffmann. “UMI or not UMI, that is the question for scRNA-seq zero-inflation”. *Nat. Biotechnol.* 39.2 (Feb. 2021), pp. 158–159.
- [16] Giulio Caravagna, Ylenia Giarratano, Daniele Ramazzotti, Ian Tomlinson, Trevor A Graham, Guido Sanguinetti, and Andrea Sottoriva. “Detecting repeated cancer evolution from multi-region tumor sequencing data”. *Nat. Methods* 15.9 (Sept. 2018), pp. 707–714.
- [17] Giulio Caravagna, Timon Heide, Marc Williams, Luis Zapata, Daniel Nichol, Ketevan Chkhaidze, William Cross, George D Cresswell, Benjamin Werner, et al. “Model-based tumor subclonal reconstruction”. Mar. 2019.
- [18] Anna K Casasent, Aislyn Schalck, Ruli Gao, Emi Sei, Annalyssa Long, William Pangburn, Tod Casasent, Funda Meric-Bernstam, Mary E Edgerton, et al. “Multiclonal Invasion in Breast Tumors Identified by Topographic Single Cell Sequencing”. *Cell* 172.1-2 (Jan. 2018), 205–217.e12.
- [19] Geng Chen, Baitang Ning, and Tielu Shi. “Single-Cell RNA-Seq Technologies and Related Computational Data Analysis”. *Front. Genet.* 10 (Apr. 2019), p. 317.
- [20] Jiahuan Chen, Qian Zhou, Yangfan Wang, and Kang Ning. “Single-cell SNP analyses and interpretations based on RNA-Seq data for colon cancer research”. *Sci. Rep.* 6.1 (Sept. 2016), p. 34420.
- [21] Kevin J Cheung and Andrew J Ewald. “A collective route to metastasis: Seeding by tumor cell clusters”. *Science* 352.6282 (Apr. 2016), pp. 167–169.
- [22] Jonas Demeulemeester, Parveen Kumar, Elen K Møller, Silje Nord, David C Wedge, April Peterson, Randi R Mathiesen, Renathe Fjellidal, Masoud Zamani Esteki, et al. “Tracing the origin of disseminated tumor cells in breast cancer using single-cell sequencing”. *Genome Biol.* 17.1 (Dec. 2016), pp. 1–15.
- [23] Amit G Deshwar, Shankar Vembu, Christina K Yung, Gun Ho Jang, Lincoln Stein, and Quaid Morris. “PhyloWGS: reconstructing subclonal composition and evolution from whole-genome sequencing of tumors”. *Genome Biol.* 16 (Feb. 2015), p. 35.
- [24] R Desper, F Jiang, O P Kallioniemi, H Moch, C H Papadimitriou, and A A Schäffer. “Inferring tree models for oncogenesis from comparative genome hybridization data”. *J. Comput. Biol.* 6.1 (1999), pp. 37–51.
- [25] Harald Detering, Laura Tomás, Tamara Prieto, and David Posada. “Accuracy of somatic variant detection in multiregional tumor sequencing data”. May 2019.
- [26] Alexander Dobin, Carrie A Davis, Felix Schlesinger, Jorg Drenkow, Chris Zaleski, Sonali Jha, Philippe Batut, Mark Chaisson, and Thomas R Gingeras. “STAR: ultrafast universal RNA-seq aligner”. *Bioinformatics* 29.1 (Jan. 2013), pp. 15–21.

- [27] Jean Fan, Kamil Slowikowski, and Fan Zhang. “Single-cell transcriptomics in cancer: computational challenges and opportunities”. *Exp. Mol. Med.* 52.9 (Sept. 2020), pp. 1452–1465.
- [28] Joseph Felsenstein. “Confidence Limits on Phylogenies: An Approach Using the Bootstrap”. *Evolution* 39.4 (1985), pp. 783–791.
- [29] Stephen J Fleming, John C Marioni, and Mehrtash Babadi. “CellBender remove-background: a deep generative model for unsupervised removal of background noise from scRNA-seq datasets”. Oct. 2019.
- [30] Robert P. Freckleton. “Fast likelihood calculations for comparative analyses”. *Methods in Ecology and Evolution* 3.5 (2012), pp. 940–947. DOI: 10.1111/j.2041-210X.2012.00220.x. eprint: <https://besjournals.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2041-210X.2012.00220.x>. URL: <https://besjournals.onlinelibrary.wiley.com/doi/abs/10.1111/j.2041-210X.2012.00220.x>.
- [31] Ruli Gao, Shanshan Bai, Ying C Henderson, Yiyun Lin, Aislyn Schalek, Yun Yan, Tapsi Kumar, Min Hu, Emi Sei, et al. “Delineating copy number and clonal substructure in human tumors from single-cell transcriptomes”. *Nat. Biotechnol.* (Jan. 2021), pp. 1–10.
- [32] Charles Gawad, Winston Koh, and Stephen R Quake. “Single-cell genome sequencing: current state of the science”. *Nat. Rev. Genet.* 17.3 (Mar. 2016), pp. 175–188.
- [33] Marco Gerlinger, Andrew J Rowan, Stuart Horswell, M Math, James Larkin, David Endesfelder, Eva Gronroos, Pierre Martinez, Nicholas Matthews, et al. “Intratumor heterogeneity and branched evolution revealed by multiregion sequencing”. *N. Engl. J. Med.* 366.10 (Mar. 2012), pp. 883–892.
- [34] Laura González-Silva, Laura Quevedo, and Ignacio Varela. “Tumor Functional Heterogeneity Unraveled by scRNA-seq Technologies”. *Trends Cancer Res.* 6.1 (Jan. 2020), pp. 13–19.
- [35] Alan Grafen and William Donald Hamilton. “The phylogenetic regression”. *Philosophical Transactions of the Royal Society of London. B, Biological Sciences* 326.1233 (1989), pp. 119–157. DOI: 10.1098/rstb.1989.0106. eprint: <https://royalsocietypublishing.org/doi/pdf/10.1098/rstb.1989.0106>. URL: <https://royalsocietypublishing.org/doi/abs/10.1098/rstb.1989.0106>.
- [36] Akdes Serin Harmanci, Arif O Harmanci, and Xiaobo Zhou. “CaSpER identifies and visualizes CNV events by integrative analysis of single-cell or bulk RNA-sequencing data”. *Nat. Commun.* 11.1 (Jan. 2020), pp. 1–16.
- [37] Akdes Serin Harmanci, Arif O Harmanci, and Xiaobo Zhou. “Inference of Clonal Copy Number Alterations from RNA-sequencing data”. *Journal of Cancer Immunology* 2.3 (2020).
- [38] Yehudit Hasin, Marcus Seldin, and Aldons Lusi. “Multi-omics approaches to disease”. *Genome Biol.* 18.1 (May 2017), p. 83.
- [39] Joseph Heled and Alexei J Drummond. “Bayesian inference of population size history from multiple loci”. *BMC Evol. Biol.* 8 (Oct. 2008), p. 289.
- [40] Stephanie C Hicks, F William Townes, Mingxiang Teng, and Rafael A Irizarry. “Missing data and technical variability in single-cell RNA-sequencing experiments”. *Biostatistics* 19.4 (Oct. 2018), pp. 562–578.
- [41] Chiaowen Joyce Hsiao, Poyuan Tung, John D Blischak, Jonathan E Burnett, Kenneth A Barr, Kushal K Dey, Matthew Stephens, and Yoav Gilad. “Characterizing and inferring quantitative cell cycle phase in single-cell RNA-seq data analysis”. *Genome Res.* 30.4 (Apr. 2020), pp. 611–621.
- [42] Richard R Hudson et al. “Gene genealogies and the coalescent process”. *Oxford surveys in evolutionary biology* 7.1 (1990), p. 44.
- [43] J P Huelsenbeck. “Performance of phylogenetic methods in simulation”. *Syst. Biol.* (1995).
- [44] J P Huelsenbeck. “The robustness of two phylogenetic methods: four-taxon simulations reveal a slight superiority of maximum likelihood over neighbor joining”. *Mol. Biol. Evol.* 12.5 (Sept. 1995), pp. 843–849.

- [45] Katharina Jahn, Jack Kuipers, and Niko Beerenwinkel. “Tree inference for single-cell data”. *Genome Biol.* 17 (May 2016), p. 86.
- [46] Livnat Jerby-Arnon, Parin Shah, Michael S Cuoco, Christopher Rodman, Mei-Ju Su, Johannes C Melms, Rachel Leeson, Abhay Kanodia, Shaolin Mei, et al. “A Cancer Cell Program Promotes T Cell Exclusion and Resistance to Checkpoint Blockade”. *Cell* 175.4 (Nov. 2018), 984–997.e24.
- [47] Wei Jiang, Si-Yun Chen, Hong Wang, De-Zhu Li, and John J Wiens. “Should genes with missing data be excluded from phylogenetic analyses?” *Mol. Phylogenet. Evol.* 80 (Nov. 2014), pp. 308–318.
- [48] Mohammed El-Kebir, Gryte Satas, and Benjamin J Raphael. “Inferring parsimonious migration histories for metastatic cancers”. *Nat. Genet.* 50.5 (May 2018), pp. 718–726.
- [49] S.W. Kembel, P.D. Cowan, M.R. Helmus, W.K. Cornwell, H. Morlon, D.D. Ackerly, S.P. Blomberg, and C.O. Webb. “Picante: R tools for integrating phylogenies and ecology”. *Bioinformatics* 26 (2010), pp. 1463–1464.
- [50] Chelsea U Kidwell, Joseph R Casalini, Soorya Pradeep, Sandra D Scherer, Daniel Greiner, Jarrod S Johnson, Gregory S Olson, Jared Rutter, Alana L Welm, et al. “Laterally transferred macrophage mitochondria act as a signaling source promoting cancer cell proliferation”. Aug. 2021.
- [51] J F C Kingman. “The coalescent”. *Stochastic Process. Appl.* 13.3 (Sept. 1982), pp. 235–248.
- [52] Natalia L Komarova. “Spatial stochastic models for cancer initiation and progression”. *Bull. Math. Biol.* 68.7 (Oct. 2006), pp. 1573–1599.
- [53] Alexey Kozlov, João Alves, Alexandros Stamatakis, and David Posada. “CellPhy: accurate and fast probabilistic inference of single-cell phylogenies from scDNA-seq data”. Aug. 2020.
- [54] Alexey M Kozlov, Diego Darriba, Tomáš Flouri, Benoit Morel, and Alexandros Stamatakis. “RAxML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference”. *Bioinformatics* 35.21 (Nov. 2019), pp. 4453–4455.
- [55] M K Kuhner, J Yamato, and J Felsenstein. “Maximum likelihood estimation of population growth rates based on the coalescent”. *Genetics* 149.1 (May 1998), pp. 429–434.
- [56] Jack Kuipers, Mustafa Anil Tuncel, Pedro Ferreira, Katharina Jahn, and Niko Beerenwinkel. “Single-cell copy number calling and event history reconstruction”. Apr. 2020.
- [57] David Lähnemann, Johannes Köster, Ewa Szczurek, Davis J McCarthy, Stephanie C Hicks, Mark D Robinson, Catalina A Vallejos, Kieran R Campbell, Niko Beerenwinkel, et al. “Eleven grand challenges in single-cell data science”. *Genome Biol.* 21.1 (Dec. 2020), p. 1.
- [58] Philippe Lemey, Andrew Rambaut, Alexei J Drummond, and Marc A Suchard. “Bayesian phylogeography finds its roots”. *PLoS Comput. Biol.* 5.9 (Sept. 2009), e1000520.
- [59] Philippe Lemey, Andrew Rambaut, John J Welch, and Marc A Suchard. “Phylogeography takes a relaxed random walk in continuous space and time”. *Mol. Biol. Evol.* 27.8 (Aug. 2010), pp. 1877–1885.
- [60] Marco L Leung, Alexander Davis, Ruli Gao, Anna Casasent, Yong Wang, Emi Sei, Eduardo Vilar, Dipen Maru, Scott Kopetz, et al. “Single-cell DNA sequencing reveals a late-dissemination model in metastatic colorectal cancer”. *Genome Res.* 27.8 (Aug. 2017), pp. 1287–1299.
- [61] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, Richard Durbin, et al. “The Sequence Alignment/Map format and SAMtools”. *Bioinformatics* 25.16 (June 2009), pp. 2078–2079. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btp352. eprint: <https://academic.oup.com/bioinformatics/article-pdf/25/16/2078/531810/btp352.pdf>. URL: <https://doi.org/10.1093/bioinformatics/btp352>.
- [62] Shaoping Ling, Zheng Hu, Zuyu Yang, Fang Yang, Yawei Li, Pei Lin, Ke Chen, Lili Dong, Lihua Cao, et al. “Extremely high genetic diversity in a single tumor points to prevalence of

- non-Darwinian cell evolution”. *Proc. Natl. Acad. Sci. U. S. A.* 112.47 (Nov. 2015), E6496–E6505.
- [63] Fenglin Liu, Yuanyuan Zhang, Lei Zhang, Ziyi Li, Qiao Fang, Ranran Gao, and Zemin Zhang. “Systematic comparative analysis of single-nucleotide variant detection methods from single-cell RNA sequencing data”. *Genome Biol.* 20.1 (Nov. 2019), p. 242.
- [64] Aaron T L Lun, Samantha Riesenfeld, Tallulah Andrews, The Phuong Dao, Tomas Gomes, participants in the 1st Human Cell Atlas Jamboree, and John C Marioni. “EmptyDrops: distinguishing cells from empty droplets in droplet-based single-cell RNA sequencing data”. *Genome Biol.* 20.1 (Mar. 2019), p. 63.
- [65] Lovelace J Luquette, Craig L Bohrsen, Max A Sherman, and Peter J Park. “Identification of somatic mutations in single cell DNA-seq using a spatial model of allelic imbalance”. *Nat. Commun.* 10.1 (Aug. 2019), p. 3908.
- [66] Salem Malikic, Katharina Jahn, Jack Kuipers, S Cenk Sahinalp, and Niko Beerenwinkel. “Integrative inference of subclonal tumour evolution from single-cell and bulk sequencing data”. *Nat. Commun.* 10.1 (June 2019), p. 2750.
- [67] Elaine R Mardis and Richard K Wilson. “Cancer genome sequencing: a review”. *Hum. Mol. Genet.* 18.R2 (Oct. 2009), R163–8.
- [68] Bui Quang Minh, Heiko A Schmidt, Olga Chernomor, Dominik Schrempf, Michael D Woodhams, Arndt von Haeseler, and Robert Lanfear. “IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era”. *Molecular Biology and Evolution* 37.5 (Feb. 2020), pp. 1530–1534. ISSN: 0737-4038. DOI: 10.1093/molbev/msaa015. eprint: <https://academic.oup.com/mbe/article-pdf/37/5/1530/33386032/msaa015.pdf>. URL: <https://doi.org/10.1093/molbev/msaa015>.
- [69] Andy J Minn, Gaorav P Gupta, Peter M Siegel, Paula D Bos, Weiping Shu, Dilip D Giri, Agnes Viale, Adam B Olshen, William L Gerald, et al. “Genes that mediate breast cancer metastasis to lung”. *Nature* 436.7050 (July 2005), pp. 518–524.
- [70] Sayaka Miura, Karen Gomez, Oscar Murillo, Louise A Huuki, Tracy Vu, Tiffany Buturla, and Sudhir Kumar. “Predicting clone genotypes from tumor bulk sequencing of multiple samples”. *Bioinformatics* 34.23 (June 2018), pp. 4017–4026.
- [71] Sayaka Miura, Tracy Vu, Jiamin Deng, Tiffany Buturla, Olumide Oladeinde, Jiyeong Choi, and Sudhir Kumar. “Power and pitfalls of computational methods for inferring clone phylogenies and mutation orders from bulk sequencing data”. *Sci. Rep.* 10.1 (Feb. 2020), p. 3498.
- [72] Nicola F Müller, David A Rasmussen, and Tanja Stadler. “The Structured Coalescent and Its Approximations”. *Mol. Biol. Evol.* 34.11 (Nov. 2017), pp. 2970–2981.
- [73] Sören Müller, Ara Cho, Siyuan J Liu, Daniel A Lim, and Aaron Diaz. “CONICS integrates scRNA-seq with DNA sequencing to map gene expression to tumor sub-clones”. *Bioinformatics* 34.18 (Sept. 2018), pp. 3217–3219.
- [74] Sören Müller, Gary Kohanbash, S John Liu, Beatriz Alvarado, Diego Carrera, Aparna Bhaduri, Payal B Watchmaker, Garima Yagnik, Elizabeth Di Lullo, et al. “Single-cell profiling of human gliomas reveals macrophage ontogeny as a basis for regional differences in macrophage activation in the tumor microenvironment”. *Genome Biol.* 18.1 (Dec. 2017), p. 234.
- [75] Matthew A Myers, Gryte Satas, and Benjamin J Raphael. “CALDER: Inferring Phylogenetic Trees from Longitudinal Tumor Samples”. *Cell Syst* 8.6 (June 2019), 514–522.e5.
- [76] Hidewaki Nakagawa and Masashi Fujita. “Whole genome sequencing analysis for cancer genomics and precision medicine”. *Cancer Sci.* 109.3 (Mar. 2018), pp. 513–522.
- [77] Anna S Nam, Ronan Chaligne, and Dan A Landau. “Integrating genetic and non-genetic determinants of cancer evolution by single-cell multi-omics”. *Nat. Rev. Genet.* (Aug. 2020).
- [78] Nicholas Navin, Jude Kendall, Jennifer Troge, Peter Andrews, Linda Rodgers, Jeanne McIndoo, Kerry Cook, Asya Stepansky, Dan Levy, et al. “Tumour evolution inferred by single-cell sequencing”. *Nature* 472.7341 (Apr. 2011), pp. 90–94.

- [79] Thale Kristin Olsen and Ninib Baryawno. “Introduction to Single-Cell RNA Sequencing”. *Current Protocols in Molecular Biology* 122.1 (2018), e57. DOI: 10.1002/cpmb.57. eprint: <https://currentprotocols.onlinelibrary.wiley.com/doi/pdf/10.1002/cpmb.57>. URL: <https://currentprotocols.onlinelibrary.wiley.com/doi/abs/10.1002/cpmb.57>.
- [80] Mark Pagel, Andrew Meade, and Daniel Barker. “Bayesian Estimation of Ancestral Character States on Phylogenies”. *Systematic Biology* 53.5 (Oct. 2004), pp. 673–684. ISSN: 1063-5157. DOI: 10.1080/10635150490522232. eprint: <https://academic.oup.com/sysbio/article-pdf/53/5/673/24197159/53-5-673.pdf>. URL: <https://doi.org/10.1080/10635150490522232>.
- [81] Anna Petrackova, Michal Vasinek, Lenka Sedlarikova, Tereza Dyskova, Petra Schneiderova, Tomas Novosad, Tomas Papajik, and Eva Kriegova. “Standardization of Sequencing Coverage Depth in NGS: Recommendation for Detection of Clonal and Subclonal Mutations in Cancer Diagnostics”. *Front. Oncol.* 9 (Sept. 2019), p. 851.
- [82] Hervé Philippe, Henner Brinkmann, Dennis V Lavrov, D Timothy J Littlewood, Michael Manuel, Gert Wörheide, and Denis Baurain. “Resolving difficult phylogenetic questions: why more sequences are not enough”. *PLoS Biol.* 9.3 (Mar. 2011), e1000602.
- [83] Olivier Poirion, Xun Zhu, Travers Ching, and Lana X Garmire. “Using single nucleotide variations in single-cell RNA-seq to identify subpopulations and genotype-phenotype linkage”. *Nat. Commun.* 9.1 (Nov. 2018), p. 4892.
- [84] Ryan Poplin, Valentin Ruano-Rubio, Mark A. DePristo, Tim J. Fennell, Mauricio O. Carneiro, Geraldine A. Van der Auwera, David E. Kling, Laura D. Gauthier, Ami Levy-Moonshine, et al. “Scaling accurate genetic variant discovery to tens of thousands of samples”. *bioRxiv* (2018). DOI: 10.1101/201178. eprint: <https://www.biorxiv.org/content/early/2018/07/24/201178.full.pdf>. URL: <https://www.biorxiv.org/content/early/2018/07/24/201178>.
- [85] David Posada. “Cancer Molecular Evolution”. *J. Mol. Evol.* 81.3-4 (Oct. 2015), pp. 81–83.
- [86] David Posada. “CellCoal: Coalescent Simulation of Single-Cell Sequencing Samples”. *Mol. Biol. Evol.* 37.5 (May 2020), pp. 1535–1542.
- [87] Nicola E Potter, Luca Ermini, Elli Papaemmanuil, Giovanni Cazzaniga, Gowri Vijayaraghavan, Ian Titley, Anthony Ford, Peter Campbell, Lyndal Kearney, et al. “Single-cell mutational profiling and clonal phylogeny in cancer”. *Genome Res.* 23.12 (Dec. 2013), pp. 2115–2125.
- [88] Julien Racle, Kaat de Jonge, Petra Baumgaertner, Daniel E Speiser, and David Gfeller. “Simultaneous enumeration of cancer and immune cell types from bulk tumor gene expression data”. *eLife* 6 (Nov. 2017).
- [89] Manisha Rao, Ki Oh, Richard Moffitt, Patricia Thompson, Jinyu Li, Jingxuan Liu, Aaron Sasson, George Georgakis, Joseph Kim, et al. “Comparative single-cell RNA sequencing (scRNA-seq) reveals liver metastasis-specific targets in a patient with small intestinal neuroendocrine cancer”. *Cold Spring Harb Mol Case Stud* 6.2 (Apr. 2020).
- [90] Sneha R Rao, Jason A Somarelli, Erdem Altunel, Laura E Selmic, Mark Byrum, Maya U Sheth, Serene Cheng, Kathryn E Ware, So Young Kim, et al. “From the Clinic to the Bench and Back Again in One Dog Year: How a Cross-Species Pipeline to Identify New Treatments for Sarcoma Illuminates the Path Forward in Precision Medicine”. *Front. Oncol.* 10 (Feb. 2020), p. 117.
- [91] Johannes G Reiter, Alvin P Makohon-Moore, Jeffrey M Gerold, Ivana Bozic, Krishnendu Chatterjee, Christine A Iacobuzio-Donahue, Bert Vogelstein, and Martin A Nowak. “Reconstructing metastatic seeding patterns of human cancers”. *Nat. Commun.* 8 (Jan. 2017), p. 14114.

- [92] Andrew Roth, Andrew McPherson, Emma Laks, Justina Biele, Damian Yap, Adrian Wan, Maia A Smith, Cydney B Nielsen, Jessica N McAlpine, et al. “Clonal genotype and population structure inference from single-cell tumor sequencing”. *Nat. Methods* 13.7 (July 2016), pp. 573–576.
- [93] Gryte Satas and Benjamin J Raphael. “Haplotype phasing in single-cell DNA-sequencing data”. *Bioinformatics* 34.13 (July 2018), pp. i211–i217.
- [94] Patricia M Schnepf, Mengjie Chen, Evan T Keller, and Xiang Zhou. “SNV identification from single-cell RNA sequencing data”. *Hum. Mol. Genet.* 28.21 (Nov. 2019), pp. 3569–3583.
- [95] Russell Schwartz and Alejandro A. Schäffer. “The evolution of tumour phylogenetics: principles and practice”. *Nature Reviews Genetics* 18.4 (Apr. 2017), pp. 213–229. ISSN: 1471-0064. DOI: 10.1038/nrg.2016.170. URL: <https://doi.org/10.1038/nrg.2016.170>.
- [96] Jochen Singer, Jack Kuipers, Katharina Jahn, and Niko Beerenwinkel. “Single-cell mutation identification via phylogenetic inference”. *Nat. Commun.* 9.1 (Dec. 2018), p. 5144.
- [97] Hojun Song, Nathan C Sheffield, Stephen L Cameron, Kelly B Miller, and Michael F Whiting. “When phylogenetic assumptions are violated: base compositional heterogeneity and among-site rate variation in beetle mitochondrial phylogenomics”. *Syst. Entomol.* 35.3 (July 2010), pp. 429–448.
- [98] Samuel P Strom. “Current practices and guidelines for clinical next-generation sequencing oncology testing”. *Cancer Biol Med* 13.1 (Mar. 2016), pp. 3–11.
- [99] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi, William M Mauck III, Yuhao Hao, Marlon Stoeckius, Peter Smibert, et al. “Comprehensive Integration of Single-Cell Data”. *Cell* 177 (2019), pp. 1888–1902. DOI: 10.1016/j.cell.2019.05.031. URL: <https://doi.org/10.1016/j.cell.2019.05.031>.
- [100] Maxime Tarabichi, Adriana Salcedo, Amit G Deshwar, Máire Ni Leathlobhair, Jeff Wintersinger, David C Wedge, Peter Van Loo, Quaid D Morris, and Paul C Boutros. “A practical guide to cancer subclonal reconstruction from DNA sequencing”. *Nat. Methods* 18.2 (Jan. 2021), pp. 144–155.
- [101] Simon Tavaré. “Some probabilistic and statistical problems in the analysis of DNA sequences”. *Lectures on mathematics in the life sciences* 17.2 (1986), pp. 57–86.
- [102] F William Townes and Rafael A Irizarry. “Quantile normalization of single-cell RNA-seq read counts without unique molecular identifiers”. *Genome Biol.* 21.1 (July 2020), pp. 1–17.
- [103] Koen Van den Berge, Fanny Perraudeau, Charlotte Soneson, Michael I Love, Davide Risso, Jean-Philippe Vert, Mark D Robinson, Sandrine Dudoit, and Lieven Clement. “Observation weights unlock bulk RNA-seq tools for zero inflation and single-cell applications”. *Genome Biol.* 19.1 (Feb. 2018), p. 24.
- [104] Bin Wang, Yingyi Zhang, Tao Qing, Kaichen Xing, Jie Li, Timing Zhen, Sibao Zhu, and Xianbao Zhan. “Comprehensive analysis of metastatic gastric cancer tumour cells using single-cell RNA-seq”. *Sci. Rep.* 11.1 (Jan. 2021), pp. 1–10.
- [105] Haiying Wang, Francisco Azuaje, Olivier Bodenreider, and Joaquín Dopazo. “Gene Expression Correlation and Gene Ontology-Based Similarity: An Assessment of Quantitative Relationships”. *Proc. IEEE Symp. Comput. Intell. Bioinforma. Comput. Biol.* 2004 (Oct. 2004), pp. 25–31.
- [106] Campbell O Webb. “Exploring the Phylogenetic Structure of Ecological Communities: An Example for Rain Forest Trees”. *Am. Nat.* 156.2 (Aug. 2000), pp. 145–155.
- [107] Benjamin Werner, Jack Case, Marc J Williams, Kate Chkhaidze, Daniel Temko, Javier Fernandez-Mateos, George D Cresswell, Daniel Nichol, William Cross, et al. “Measuring single cell divisions in human cancers from multi-region sequencing data”. Feb. 2019.
- [108] John J Wiens. “Missing data and the design of phylogenetic analyses”. *J. Biomed. Inform.* 39.1 (Feb. 2006), pp. 34–42.

- [109] Zhenxiang Xi, Liang Liu, and Charles C Davis. “The Impact of Missing Data on Species Tree Estimation”. *Mol. Biol. Evol.* 33.3 (Mar. 2016), pp. 838–860.
- [110] Xiaoqing Yu, Y Ann Chen, Jose R Conejo-Garcia, Christine H Chung, and Xuefeng Wang. “Estimation of immune cell content in tumor using single-cell RNA-seq reference data”. *BMC Cancer* 19.1 (July 2019), p. 715.
- [111] Ke Yuan, Thomas Sakoparnig, Florian Markowitz, and Niko Beerenwinkel. “BitPhylogeny: a probabilistic framework for reconstructing intra-tumor phylogenies”. *Genome Biol.* 16.1 (Feb. 2015), p. 36.
- [112] H Zafar, N Navin, K Chen, and L Nakhleh. “SiCloneFit: Bayesian inference of population structure, genotype, and phylogeny of tumor clones from single-cell genome sequencing data”. *bioRxiv* (2018).
- [113] Hamim Zafar, Anthony Tzen, Nicholas Navin, Ken Chen, and Luay Nakhleh. “SiFit: inferring tumor trees from single-cell sequencing data under finite-sites models”. *Genome Biol.* 18.1 (Sept. 2017), p. 178.
- [114] Hamim Zafar, Yong Wang, Luay Nakhleh, Nicholas Navin, and Ken Chen. “Monovar: single-nucleotide variant detection in single cells”. *Nat. Methods* 13.6 (June 2016), pp. 505–507.
- [115] Weiwei Zhai, Tony Kiat-Hon Lim, Tong Zhang, Su-Ting Phang, Zenia Tiang, Peiyong Guan, Ming-Hwee Ng, Jia Qi Lim, Fei Yao, et al. “The spatial organization of intra-tumour heterogeneity and evolutionary trajectories of metastases in hepatocellular carcinoma”. *Nat. Commun.* 8 (Feb. 2017), p. 4565.
- [116] Zi-Ming Zhao, Bixiao Zhao, Yalai Bai, Atila Iamarino, Stephen G Gaffney, Joseph Schlessinger, Richard P Lifton, David L Rimm, and Jeffrey P Townsend. “Early and multiple origins of metastatic lineages within primary tumors”. *Proc. Natl. Acad. Sci. U. S. A.* 113.8 (Feb. 2016), pp. 2140–2145.

Supplementary materials

Constructing the Panel of Normals and normal samples from the MDA-MB-231 cell lineage. The data from the Kidwell et al. (2021) contain a mixture of macrophages and MDA-MB-231 cells with and without mitochondrial transfer. To create the Panel of Normals, we have used all reads, but for the matched normal samples, only reads belonging to the MDA-MB-231 cells without mitochondrial transfer were used.

The fastq files were downloaded from the NCBI GEO database (ascension number GSE181410), mapped using the Cellranger and preprocessed using the GATK best practices as per the methodology section. The Panel of Normals was then constructed as per GATK instructions by first running Mutect2 in a tumor-only mode, merging the resulting variants into a database and then creating a Panel of Normals variant file from this database.

For the normal samples, we have used the preprocessed bam files from the previous steps. Only the bam files that contained MDA-MB-231 cells without mitochondrial transfer were retained (GSM5501832 and GSM5501833). To remove macrophages, the bam files were then filtered using the cell barcodes from the h5-Seurat expression analysis, that identified MDA-MB-231 cells without mitochondrial transfer (samples 2A and 2B in the Seurat object). These filtered bam files were used as matched normal samples in the main analysis.

Integrating topological uncertainty. To investigate how the topological uncertainty influences the phylogenetic relationship between samples, we perform the phylogenetic clustering tests over 100 bootstrap samples of the Maximum Likelihood trees and over sample of 1000 posterior trees from the Bayesian inference. First we have sub-sampled the posterior tree sample to gain a sample of 1000 trees using the `logcombiner` from the BEAST2 package. Then, we calculate the MPD and MNTD p-values for each tree in the Maximum Likelihood and Bayesian tree sample. Finally, we summarize each sample using the mean and the 95% confidence interval. The calculated values for the Maximum Likelihood and Bayesian trees reconstructed from the subset of 58 cells from the expression and SNV data are presented in Supplementary Table 1. The majority of relationship that were present on the best tree is stable when the topological uncertainty is taken in account. MNTD is less stable than MPD, likely due to a higher sensitivity of MNTD to patterns closer to the tips of a tree.

Data reduction using the stepwise filtering algorithm. Using the stepwise filtering algorithm, we iteratively remove cells and genes/SNVs with the smallest number of known values, until a desired data density is reached. By using this method, the least amount of data is removed. Here we investigate the effect of unknown data on the reconstructed topology by preparing datasets with different data densities using our stepwise filtering algorithm.

The expression data set filtered to 20% density contained 1,627 cells and 6,187 genes. Cells were mainly represented by T1 and T3 samples which form over 92% of the data set. In contrast, other samples (T2, CTC1, CTC2) were significantly underrepresented despite their larger amount of cells in an unfiltered data set. In filtering to 50% density, the numbers decreased to 1,454 cells and 1,634 genes. The sample diversity also decreased, with T2 dropping out entirely. When filtered to 90% density, the data set was reduced to 593 cells and 528 genes. The data diversity further decreased to T1, T3 and CTC2, with the CTC2 sample reduced to 2 cells.

When the SNV data set was filtered to 20% density, the numbers decreased significantly – to 870 cells and 317 SNVs, with only T1, T3 and CTC2 samples present. In subsequent filtering to 50% data density, the dataset was reduced to 254 cells and only 69 SNVs, with subsequent filtering to the 90% reducing the dataset further to 60 cells and 8 SNVs.

We inferred Maximum Likelihood trees of the expression data filtered to 20%, 50%, and 90% data density (Supplementary Figure 1). In the reconstructed phylogeny at the 20% density filtering (Supplementary Figure 1a), individual tumor samples did not form three separate clades,

SUPPLEMENTARY TABLE 1. Test of phylogenetic clustering for the reduced dataset of the 58 selected cells. Mean Pairwise Distance (MPD) and Mean Nearest Taxon Distance (MNTD) calculated for the Maximum Likelihood (ML) and Bayesian (BI) trees from the expression and SNV data. P-values for MPD and MNTD were calculated for each sample (T1, T2, T3, CTC1, CTC2) and expected clustering for cells isolated from a single individual (T1 with CTC1, and T2 with CTC2) and to test a possible mislabeling between CTC1 and CTC2 samples (T1 with CTC2, and T2 with CTC1). P-values were calculated for the sample of 100 bootstrap trees and 1000 posterior trees from the Maximum Likelihood and Bayesian analyses respectively, and this distribution of p-values is summarized with mean and 95% confidence interval. Significant p-values at $\alpha = 0.05$ after correcting for multiple comparisons using the False Discovery Rate method (Benjamini et al. 1995) are marked with an asterisk.

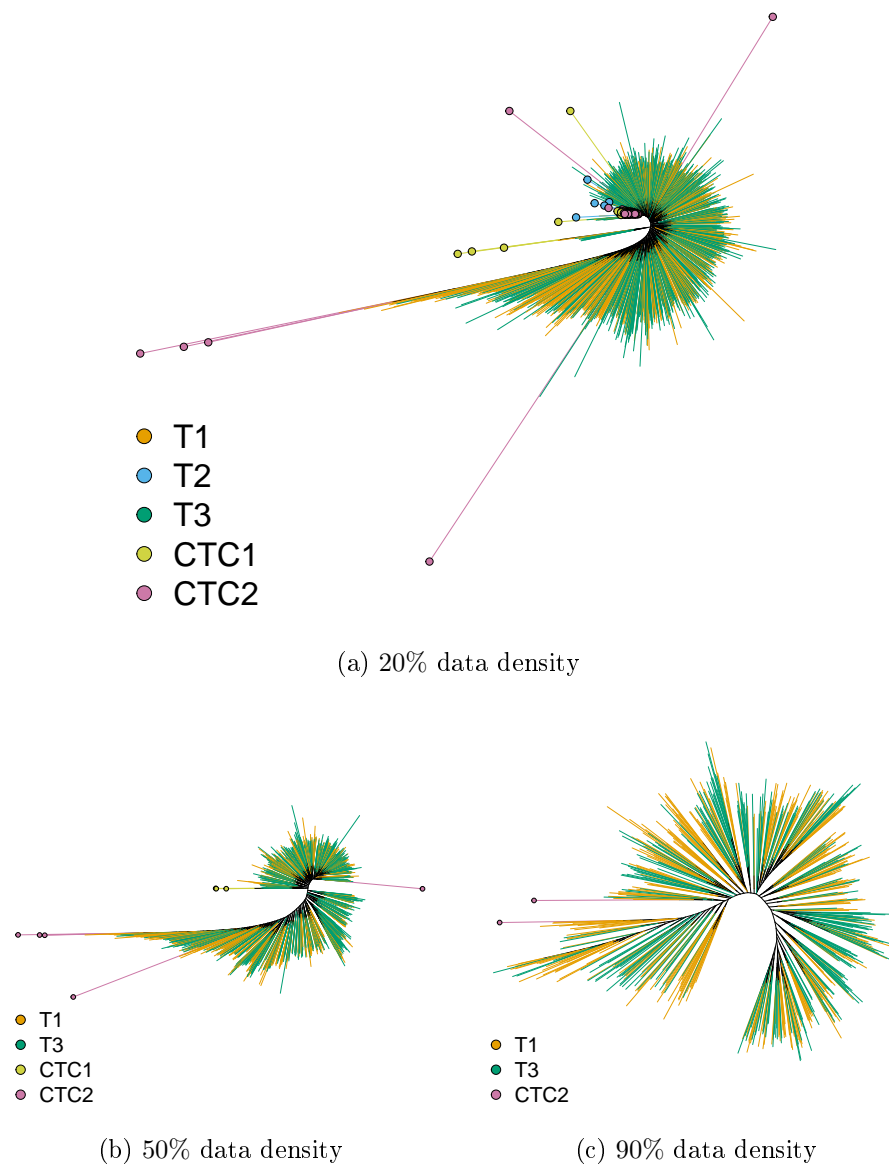
Groups	Cells	Expression (ML)		Expression (BI)		SNV (ML)		SNV (BI)	
		MPD	MNTD	MPD	MNTD	MPD	MNTD	MPD	MNTD
T1	20	*0.001 (0.001–0.002)	0.750 (0.617–0.862)	*0.001 (0.001–0.001)	*0.007 (0.001–0.020)	*0.001 (0.001–0.001)	0.058 (0.001–0.202)	*0.001 (0.001–0.002)	0.168 (0.001–0.567)
T2	6	*0.001 (0.001–0.001)	*0.001 (0.001–0.001)	*0.002 (0.001–0.006)	*0.004 (0.001–0.015)	*0.016 (0.001–0.002)	0.032 (0.001–0.008)	*0.003 (0.001–0.001)	*0.005 (0.001–0.006)
T3	20	*0.001 (0.001–0.001)	0.390 (0.280–0.521)	*0.001 (0.001–0.002)	0.039 (0.002–0.096)	*0.002 (0.001–0.004)	0.387 (0.018–0.781)	*0.001 (0.001–0.002)	0.267 (0.007–0.701)
CTC1	6	*0.003 (0.001–0.012)	*0.006 (0.001–0.014)	0.989 (0.972–0.999)	0.470 (0.263–0.602)	*0.011 (0.001–0.004)	0.024 (0.001–0.027)	*0.002 (0.001–0.002)	*0.007 (0.001–0.021)
CTC2	6	1.000 (1.000–1.000)	0.987 (0.934–1.000)	0.700 (0.644–0.743)	0.946 (0.841–1.000)	0.237 (0.110–0.318)	0.492 (0.233–0.740)	0.175 (0.092–0.249)	0.433 (0.178–0.617)
T1 & CTC1	26	0.052 (0.024–0.085)	0.427 (0.286–0.642)	0.146 (0.115–0.180)	0.281 (0.055–0.502)	0.428 (0.262–0.574)	0.083 (0.001–0.254)	0.394 (0.252–0.519)	0.226 (0.001–0.751)
T2 & CTC2	12	0.995 (0.984–1.000)	0.092 (0.001–0.322)	0.949 (0.912–0.982)	0.127 (0.001–0.295)	0.952 (0.653–1.000)	0.542 (0.111–0.856)	0.938 (0.857–0.996)	0.378 (0.092–0.704)
T1 & CTC2	26	0.574 (0.379–0.742)	1.000 (0.997–1.000)	*0.005 (0.002–0.009)	0.889 (0.736–0.989)	*0.001 (0.001–0.002)	0.451 (0.053–0.826)	*0.001 (0.001–0.003)	0.464 (0.068–0.903)
T2 & CTC1	12	*0.001 (0.001–0.001)	*0.001 (0.001–0.001)	0.654 (0.553–0.753)	0.053 (0.001–0.175)	*0.004 (0.001–0.027)	0.060 (0.001–0.734)	*0.001 (0.001–0.001)	*0.005 (0.001–0.003)

* significant support

but a large number of smaller clades. These clades are distributed along the central spine of the unrooted maximum likelihood tree and have little internal structure. The T2, CTC1, and CTC2 samples form relatively compact clades, while the more represented T1 and T3 clades are generally intermixed. The MNT and MNTD test confirm this (Supplementary Table 2), with T2, CTC1, and CTC2 showing significant clustering signal. When the data is filtered to 50% and 90% data density (Supplementary Figure 1), the previously significant relationship disappears, likely due to the reduction of T2, CTC1 and CTC2 samples into a small number of cells.

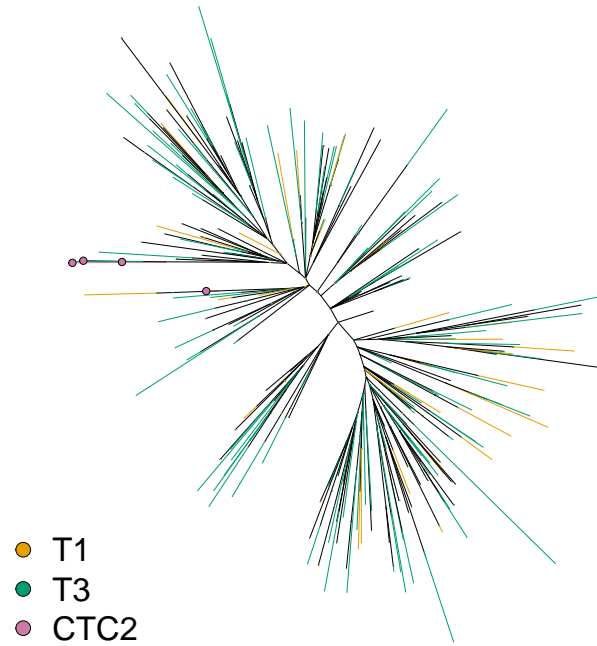
In the trees reconstructed from the SNV data, the CTC2 cells do cluster (Supplementary Figure 2), but this relationship is not significant (Supplementary Table 2), likely due to the small number of cells remaining compared to a large number of cells from the T1 and T3 datasets. Only the T1 and the putative relationship between T1 and CTC2 cells was supported, but this support disappeared when the data was further filtered. Due to the small number of SNVs for the dataset filtered to 90% data density, many cells were identical, with small or collapsed branches (Supplementary Figure 2c).

The data filtered with the stepwise filtering algorithm failed to show a convincing clustering signal for each sample. This is likely caused by the variable quality of our sample and thus the results should be interpreted in this context. On a dataset that is not burdened by similar issues, the stepwise filtering algorithm might be the preferred method. The partial clustering signal for the lowest density filtering shows that the phylogenetic methods can handle well large amount of missing data. This means that data reduction should be performed carefully to limit the required amount of computational burned and n to reduce the amount of missing data.

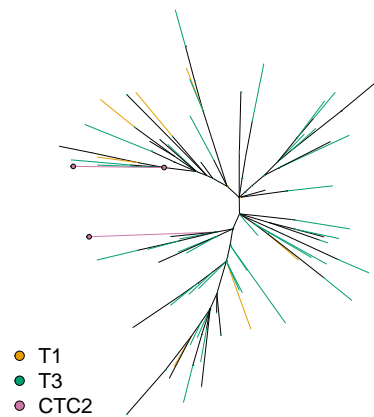


SUPPLEMENTARY FIGURE 1. Maximum Likelihood trees from the expression data for 20% (Supplementary Figure 1a), 50% (Supplementary Figure 1b), and 90% (Supplementary Figure 1c) data density. Terminal branches are colored according to cell's sample of origin (T1, T2, T3, CTC1, CTC2). The T2, CTC1 and CTC2 samples are marked with colored circles.

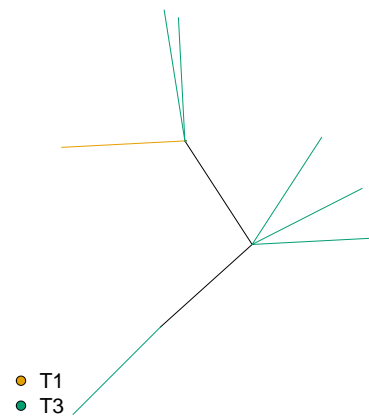
iv



(a) 20% data density



(b) 50% data density



(c) 90% data density

SUPPLEMENTARY FIGURE 2. Maximum Likelihood trees from the SNV data for 20% (Supplementary Figure 2a), 50% (Supplementary Figure 2b), and 90% (Supplementary Figure 2c) data density. Terminal branches are colored according to cell's sample of origin (T1, T2, T3, CTC1, CTC2). The T2, CTC1 and CTC2 samples are marked with colored circles.

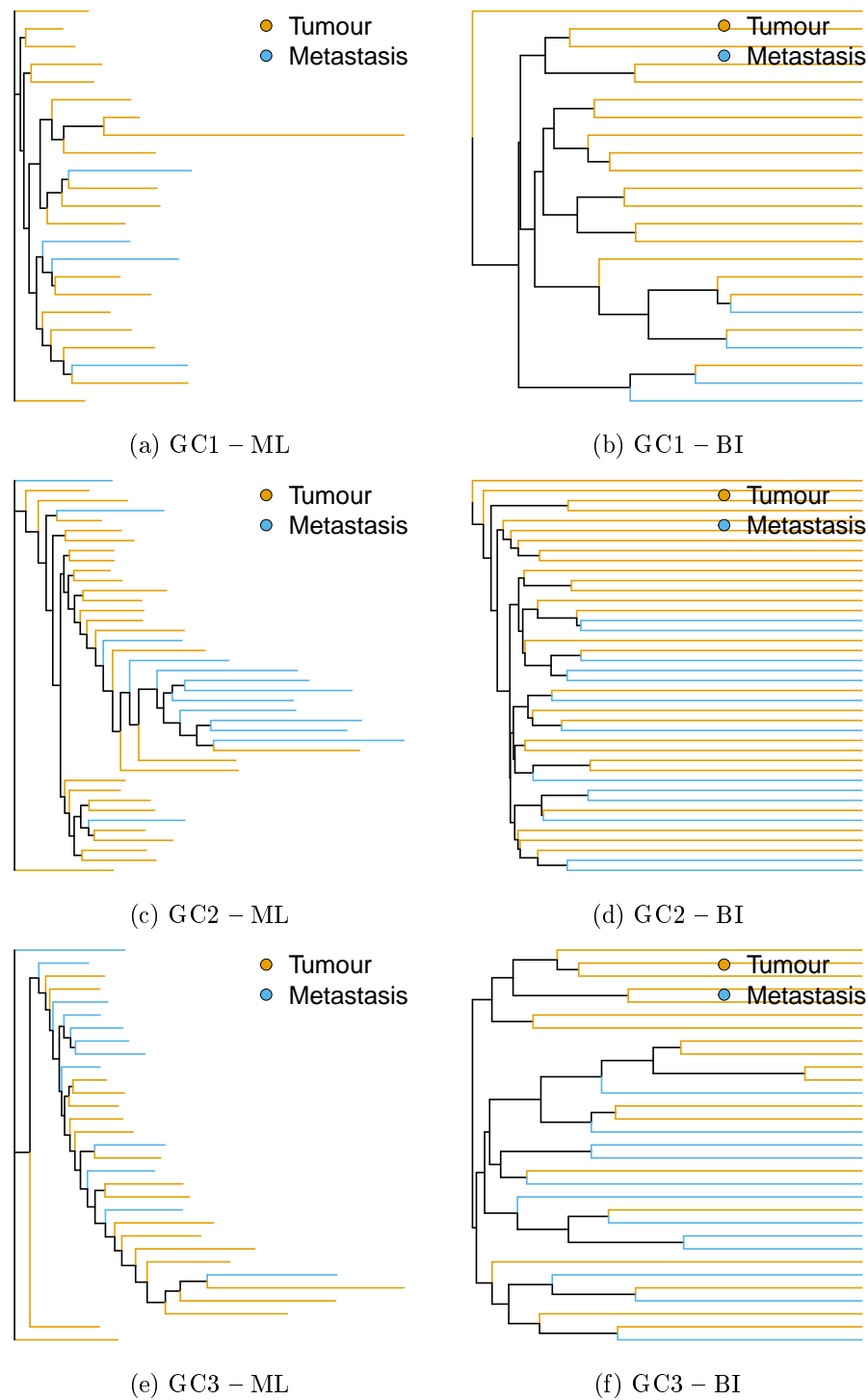
SUPPLEMENTARY TABLE 2. Test of phylogenetic clustering for the datasets filtered to 20%, 50% and 90% data density. Mean Pairwise Distance (MPD) and Mean Nearest Taxon Distance (MNTD) calculated for the Maximum Likelihood (ML) and Bayesian (BI) trees from the expression and SNV data. P-values for MPD and MNTD were calculated for each sample (T1, T2, T3, CTC1, CTC2) and expected clustering for cells isolated from a single individual (T1 with CTC1, and T2 with CTC2) and to test a possible mislabeling between CTC1 and CTC2 samples (T1 with CTC2, and T2 with CTC1). Significant p-values at $\alpha = 0.05$ after correcting for multiple comparisons using the False Discovery Rate method (Benjamini et al. 1995) are marked with an asterisk.

Data	Groups	20% data density			50% data density			90% data density		
		Cells	MPD	MNTD	Cells	MPD	MNTD	Cells	MPD	MNTD
Expression	T1	701	1.000	1.000	688	1.000	0.996	329	0.910	0.089
	T2	11	*0.001	*0.001	0	–	–	0	–	–
	T3	806	0.126	1.000	758	*0.001	*0.001	262	0.015	*0.002
	CTC1	58	*0.001	*0.001	3	0.156	0.099	0	–	–
	CTC2	51	*0.001	*0.001	5	1.000	1.000	2	1.000	1.000
	T1 & CTC1	759	0.992	0.999	691	1.000	0.998	329	0.910	0.089
	T2 & CTC2	62	*0.001	*0.001	5	1.000	1.000	2	1.000	1.000
	T1 & CTC2	752	1.000	1.000	693	1.000	1.000	331	0.972	0.148
	T2 & CTC1	69	*0.001	*0.001	3	0.156	0.099	0	–	–
SNV	T1	352	*0.001	*0.001	55	0.066	0.276	13	0.221	0.200
	T2	0	–	–	0	–	–	0	–	–
	T3	514	1.000	0.999	196	0.795	0.571	47	0.801	0.736
	CTC1	0	–	–	0	–	–	0	–	–
	CTC2	4	0.520	0.557	3	0.892	0.695	0	–	–
	T1 & CTC1	352	*0.001	*0.001	55	0.066	0.276	13	0.221	0.200
	T2 & CTC2	4	0.520	0.557	3	0.892	0.695	0	–	–
	T1 & CTC2	356	*0.001	*0.001	58	0.152	0.558	13	0.221	0.200
	T2 & CTC1	0	–	–	0	–	–	0	–	–

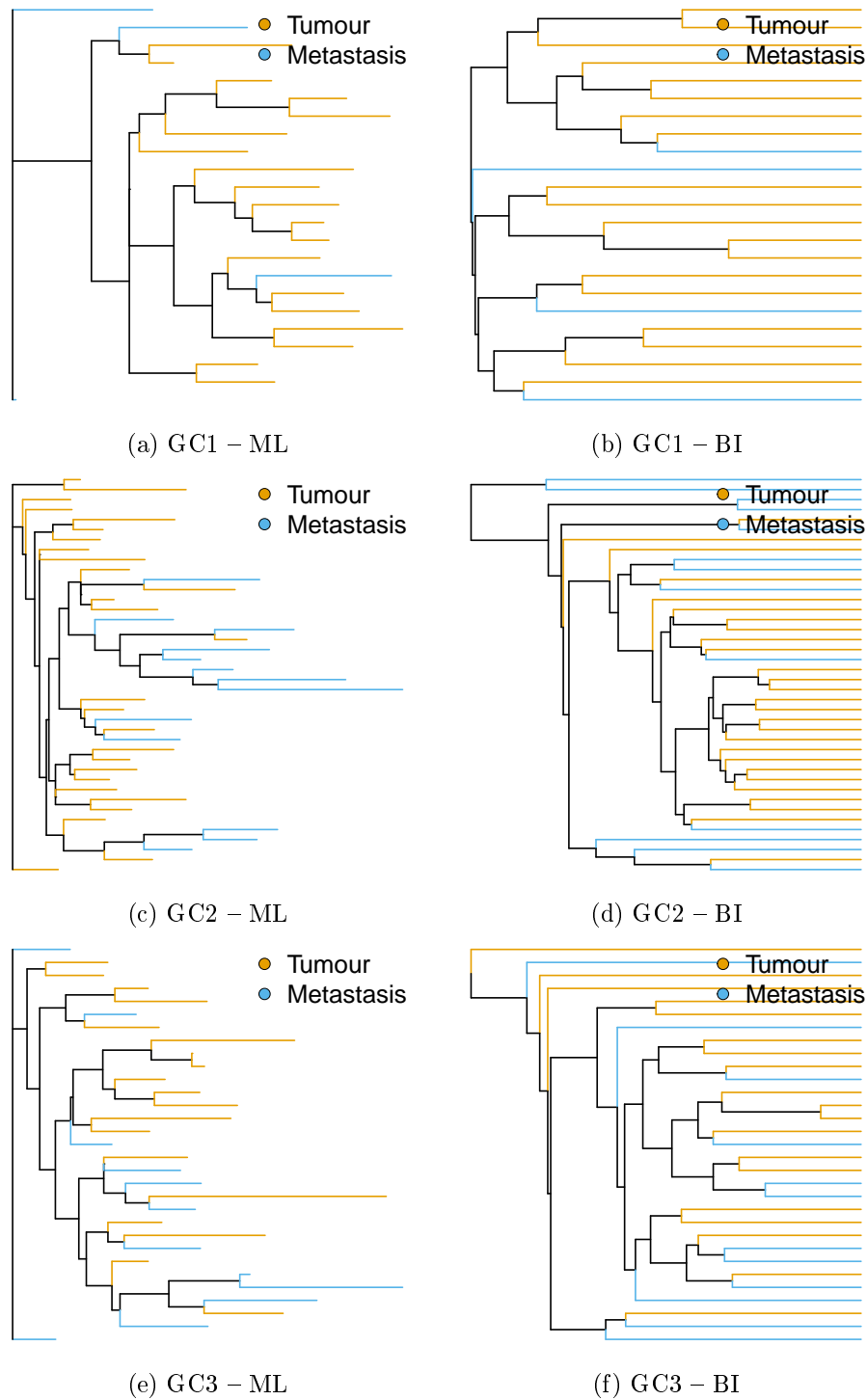
* significant support

vi

Wang et al. (2021).



SUPPLEMENTARY FIGURE 3. Maximum Likelihood Bayesian trees constructed from the expression data published by Wang et al. (2021). Terminal branches are colored according to cell's sample of origin.



SUPPLEMENTARY FIGURE 4. Maximum Likelihood Bayesian trees constructed from the SNV data published by Wang et al. (2021). Terminal branches are colored according to cell's sample of origin.