

Defining ancestry, heritability and plasticity of cellular phenotypes in somatic evolution

Joshua S. Schiffman^{1,2,†}, Andrew R. D’Avino^{1,2,3,†}, Tamara Prieto^{1,2,†}, Yakun Pang^{4,5}, Yilin Fan^{6,7}, Srinivas Rajagopalan^{1,2}, Catherine Potenski^{1,2}, Toshiro Hara^{6,7}, Mario L. Suvà^{6,7}, Charles Gawad^{4,5,8} and Dan A. Landau^{1,2,✉}

Summary

The broad application of single-cell RNA sequencing has revealed transcriptional cell state heterogeneity across diverse healthy and malignant somatic tissues. Recent advances in lineage tracing technologies have further enabled the simultaneous capture of cell transcriptional state along with cellular ancestry thus enabling the study of somatic evolution at an unprecedented resolution; however, new analytical approaches are needed to fully harness these data. Here we introduce PATH (Phylogenetic Analysis of Transcriptional Heritability), an analytical framework, which draws upon classic approaches in species evolution, to quantify heritability and plasticity of somatic phenotypes, including transcriptional states. The PATH framework further allows for the inference of cell state transition dynamics by linking a model of cellular evolutionary dynamics with our measure of heritability versus plasticity. We evaluate the robustness of this approach by testing a range of biological and technical features in simulations of somatic evolution. We then apply PATH to characterize previously published and newly generated single-cell phylogenies, reconstructed from either native or artificial lineage markers, with matching cellular state profiling. PATH recovered developmental relationships in mouse embryogenesis, and revealed how anatomic proximity influences neural relatedness in the developing zebrafish brain. In cancer, PATH dissected the heritability of the epithelial-to-mesenchymal transition in a mouse model of pancreatic cancer, and the heritability versus plasticity of transcriptionally-defined cell states in human glioblastoma. Finally, PATH revealed phenotypic heritability patterns in a phylogeny reconstructed from single-cell whole genome sequencing of a B-cell acute lymphoblastic leukemia patient sample. Altogether, by bringing together

perspectives from evolutionary biology and emerging single-cell technologies, PATH formally connects the analysis of cell state diversity and somatic evolution, providing quantification of critical aspects of these processes and replacing *qualitative* conceptions of “plasticity” with *quantitative* measures of cell state transitions and heritability.

Introduction

THE application of single-cell RNA sequencing (scRNAseq) across biology has revealed vast phenotypic diversity within healthy [Hammond et al., 2019, Papalexli and Satija, 2018, Plasschaert et al., 2018] and diseased [Nefitel et al., 2019, Wu et al., 2021] tissues. As genetic variation is limited within the soma, much of the heritable diversity of somatic phenotypes is attributed to non-genetic sources, such as epigenetic modifications. Indeed, the stable propagation of somatic phenotypes (*e.g.*, cell type [Zeng, 2022]) through mitotic divisions, sometimes called *epigenetic memory* [Fennell et al., 2022, Halley-Stott and Gurdon, 2013, Larsen et al., 2021, Shaffer et al., 2020], often relies on the heritable transmission of epigenetic marks, such as DNA methylation, histone modification, or the propagation of key transcription factors [Adam and Fuchs, 2016, Whyte et al., 2013]. Somatic cells, however, may also accumulate genetic variation over time [Li et al., 2020, Martincorena et al., 2015, 2018], for example enabling more proliferative phenotypes that can lead to cancer [Hanahan, 2022, Vogelstein et al., 2013]. In addition to cell-intrinsic sources of heritable phenotypic diversity, cell-extrinsic sources, such as the microenvironment [Gola and Fuchs, 2021, Hara et al., 2021] or morphogen gradients [Houchmandzadeh et al., 2002], may contribute to heritable cellular phenotypic diversity, as progeny often share the same microenvironment as parent cells. Crucially, not all cellular phenotypic variation is stable, and cells can also plastically toggle between phenotypes in somatic evolution. For instance, healthy skin cells can dedifferentiate to repair injuries [Donati et al., 2017, Gola

¹New York Genome Center, New York, NY, USA. ²Weill Cornell Medicine, New York, NY, USA. ³Tri-Institutional MD-PhD Program, Weill Cornell Medicine, Rockefeller University, Memorial Sloan Kettering Cancer Center, New York, NY, USA. ⁴Stanford University, Stanford, CA, USA. ⁵St. Jude Children’s Research Hospital, Memphis, TN, USA. ⁶Department of Pathology and Center for Cancer Research, Massachusetts General Hospital and Harvard Medical School, Boston, MA, USA. ⁷Broad Institute of Harvard and MIT, Cambridge, MA, USA. ⁸Chan Zuckerberg Biohub, San Francisco, CA, USA.

† These authors contributed equally. ✉ Corresponding author; dlandau@nygenome.org

81 and Fuchs, 2021] and cancer cells have been shown to toggle
82 between proliferative and invasive phenotypes [Karras et al.,
83 2022, Oren et al., 2021], or to morph and evade treatment
84 [Chan et al., 2022].

85 To approach these key aspects, it can be useful to con-
86 sider cellular phenotypic diversity from an evolutionary per-
87 spective. Somatic cells descend from a common ances-
88 tor, and following successive divisions, accumulate heritable
89 variation in the form of genetic, epigenetic or cell-extrinsic
90 changes. Throughout this process of *somatic evolution*, the
91 heritable variation within a population can be sculpted by
92 selection, which has important implications for organismal
93 health. Outcomes of somatic evolution, for instance, include
94 the initiation, relapse, and treatment resistance of cancers
95 [Fennell et al., 2022, Jan et al., 2012, Shaffer et al., 2017].
96 However, it is not yet clear to what degree epigenetic [Mazor
97 et al., 2016] or genetic [Househam et al., 2022, Turajlic et al.,
98 2019] variation contributes to the evolution and persistence
99 of malignant phenotypes [Nam et al., 2021]. To confront
100 the challenge of studying somatic evolution, we require an
101 integrative model of somatic evolution that considers cel-
102 lular phenotypic diversity and ancestry [Nam et al., 2021],
103 informed by technologies that deliver phenotypically anno-
104 tated single-cell phylogenetic trees [Biddy et al., 2018]. By
105 tracing cellular ancestries, we can begin to elucidate the
106 shared developmental origins of cell states and map differ-
107 entiation trajectories [Chan et al., 2019, Raj et al., 2018].
108 Furthermore, this framework can enable us to dissect the
109 heritability versus plasticity of somatic cellular phenotypes,
110 to define how evolution shapes somatic cellular populations.

111 Recently, an array of techniques for lineage tracing has been
112 advanced that can provide ancestry information at a single-
113 cell level [Baron and van Oudenaarden, 2019, Sankaran
114 et al., 2022]. In model organisms, cellular lineages or phylo-
115 genies can be reconstructed from *artificial* lineage markers
116 [Pei et al., 2020, Raj et al., 2018, Rodriguez-Fraticelli et al.,
117 2020, Spanjaard et al., 2018] that can be experimentally in-
118 sserted and edited. In contrast, retracing lineage histories
119 in human samples leverages *native* lineage markers, such as
120 patterns of genetic (copy number [Salehi et al., 2022, Wang
121 et al., 2021] or single nucleotide [Lodato et al., 2015, Lud-
122 wig et al., 2019]) or epigenetic (stochastic methylation [Gaiti
123 et al., 2019]) variation. Both artificial and native lineage
124 tracing approaches can be combined with other single-cell
125 modalities, like scRNAseq, to deliver phylogenetic trees with
126 phenotypically annotated leaves (terminal nodes).

127 Such phenotypically annotated cellular lineages emerge as a
128 formidable tool to study critical questions in biology, such
129 as mapping the ontogenetic relations between cells in de-
130 velopment [Bandler et al., 2021], and clinically important
131 features of cancer evolution, such as the stability of differ-
132 entiation hierarchies [Chaligne et al., 2021], and metastatic
133 dynamics [Quinn et al., 2021]. These experimental advances
134 need to be complemented by a broadly applicable analytical

framework, grounded in evolutionary biology, that could be
applied to examine how cellular state (as for example pro-
filed by scRNAseq) depends on ancestry (delivered by line-
age tracing). Such a framework would enable us to distin-
guish between mitotically stable and ephemeral phenotypic
states, and to make inferences about unobserved evolution-
ary dynamics. Tools for the analysis of multimodal single-
cell lineages, such as *Hotspot* [Detomaso and Yosef, 2021]
and *The Lorax* [Minkina et al., 2022], and others [Chaligne
et al., 2021, Fang et al., 2022, Jones et al., 2022, Wang et al.,
2022, Yang et al., 2022], are being developed to measure
heritability. Nonetheless, additional conceptual and ana-
lytic advances are needed to fully harness these datasets for
the study of somatic evolution. These advances will allow
us to account for technical and biological variables affect-
ing heritability measurements, and enable the integration of
heritability assessments with phenotypic transition proba-
bility measurements, within a comprehensive and easy-to-
implement analytical framework.

To address this challenge, we introduce **PATH**
(**Phylogenetic Analysis of Transcriptional Heritability**), an
analytical framework that draws upon classic approaches
in species evolution, to quantify heritability and plasticity
of somatic cellular phenotypes, such as transcriptional cell
states. PATH measures *phylogenetic correlations*, which
quantify the degree by which cellular phenotypes, broadly
defined (*e.g.*, transcriptional program, cell state or location),
depend on ancestry, as provided by single-cell phylogenies,
and *thus defines a measure of somatic heritability versus*
plasticity. PATH builds upon auto-correlative [Cheverud
and Dow, 1985, Gittleman and Kot, 1990] methods clas-
sically used to measure *phylogenetic signal* [Blomberg and
Garland, 2002], the phylogenetic clustering of species phe-
notypes. Furthermore, PATH generalizes this approach to
measure phylogenetic correlations *between* phenotypes (and
from across modalities), providing a measure of how distinct
phenotypes co-cluster on phylogenies, and thus defining a
pairwise measure of phylogenetic signal. Additionally, for
categorical phenotypes, such as cell type, PATH can trans-
form phylogenetic correlations, our measurement of heri-
tability versus plasticity, into inferences of transition rates
between cell types or states. Importantly, this transforma-
tion provides a concrete interpretation of what phylogenetic
signal measures, as the *pattern* of phylogenetic signal is di-
rectly linked with the *process* of cell type or state toggling.
Further, PATH represents a comprehensive, versatile quan-
titative framework that can handle sparsely sampled and
lowly resolved phylogenies, reconstructed under a range of
biological and technical variables.

We first demonstrate PATH’s capabilities through simula-
tions reflecting plausible biological and technical param-
eters of single-cell data, including cell sampling rate, phy-
logenetic reconstruction fidelity, cellular division and death
rate, and show that PATH reproducibly and accurately mea-

189 sures heritability versus plasticity across different contexts.
190 We show how the detection of heritability depends on sam-
191 pling and phylogenetic reconstruction fidelity, and how these
192 results can guide future lineage tracing experimental de-
193 sign and methods development. PATH can infer cell type
194 transition dynamics with high accuracy, comparable to a
195 classic maximum likelihood approach from species evolu-
196 tion [Lewis, 2001, Louca and Pennell, 2019, Pagel, 1994],
197 but with higher computational efficiency, a critical feature
198 considering the massive potential scale of phenotypically
199 annotated phylogenies in high throughput single-cell data.
200 We then apply PATH to published single-cell multi-omic
201 datasets, which use either native or artificial lineage trac-
202 ing (for human and model organism data, respectively), to
203 explore two broad themes, development and cancer. Specif-
204 ically, we examine mouse embryogenesis [Chan et al., 2019]
205 and zebrafish neural development [Raj et al., 2018], a model
206 of pancreatic cancer [Simeonov et al., 2021] and human
207 glioblastoma [Chaligne et al., 2021]. PATH quantitatively
208 maps cell fate trajectories during development, character-
209 izes the variable plasticity of transcriptional states along the
210 epithelial-to-mesenchymal transition in cancer and quanti-
211 fies the heritability and stability of cell states of the cor-
212 rupted neurodevelopmental hierarchy in glioblastoma. Fi-
213 nally, we apply PATH to newly generated single-cell whole
214 genome sequencing data from a patient B-cell acute lym-
215 phoblastic leukemia (B-ALL) sample with a phylogeny con-
216 structed from somatic mutations with accompanying protein
217 marker expression data. PATH reveals heritability of cellu-
218 lar phenotypes, and quantifies plasticity of immunotherapy-
219 targeted B-cell surface markers and calculates transition
220 rates between CD19 low, medium and high cell states. We
221 make PATH available to the community as a comprehen-
222 sive package, including software, analyses, and tutorials at
223 <https://github.com/landau-lab/PATH>.

224 Results

225 Heritability, plasticity and cell state transi- 226 tion dynamics

227 Evolutionary biology offers a collection of metrics for char-
228 acterizing heritable patterns of phenotypic variation, which
229 can be adapted to interrogate single-cell ancestries. The
230 degree to which phenotypic and ancestral similarity align
231 is quantified by *heritability* statistics (h^2 and H^2) [Gille-
232 spie, 2004], which are weighted measures of the phenotypic
233 correlation between relatives. These statistics have found
234 application in agriculture, as part of the breeder’s equa-
235 tion, enabling the prediction of a phenotypic response to an
236 artificial selection pressure [Gillespie, 2004]. Analogously,
237 through leveraging phylogenetic trees, the degree to which
238 related species phenotypically resemble each other, termed
239 *phylogenetic signal* [Blomberg and Garland, 2002], can be
240 quantified with various metrics (*e.g.*, Pagel’s λ [Househam

et al., 2022, Pagel, 1999], Blomberg’s K [Blomberg et al.,
2003], Moran’s I [Gittleman and Kot, 1990]), and is used to
make inferences about inheritance patterns and the evolu-
tionary lability of phenotypes. These metrics are sometimes
categorized as either statistic- or model-based [Münkemüller
et al., 2012], but nonetheless show strong agreement [Diniz-
Filho et al., 2012]. Signal statistics, such as Moran’s I , quan-
tify the phylogenetic dependency of a phenotype, whereas
model-based metrics, such as Pagel’s λ , assess the diver-
gence between a phenotype’s phylogenetic distribution with
a distribution expected by a model of random genetic drift.
PATH builds upon these approaches to characterize the her-
itability or plasticity of cellular states in somatic evolution.

Specifically, PATH adapts Moran’s I (**Methods: Phy-**
logenetic correlations), a measure of *phylogenetic auto-*
correlation and phylogenetic signal (but originally conce-
ived as a spatial auto-correlation metric [Moran, 1950]),
to quantify the heritability or plasticity of single-cell pheno-
types. Like classic heritability statistics, phylogenetic auto-
correlation is a measure of phenotypic similarity, weighted
by relatedness. Phylogenetic auto-correlation quantifies the
phylogenetic dependency of a single-cell measurement or
phenotype (broadly defined), such as cellular state, tran-
scriptional profile, or spatial location. Fundamentally, phy-
logenetic auto-correlation measures how much phenotypic
resemblance close relatives have to one another compared to
randomly chosen cells. If cells resemble close relatives much
more than randomly chosen cells, the phenotype will appear
highly heritable and phylogenetically auto-correlated. Such
a pattern might be observed for a genetically encoded phe-
notype, as for example a phenotype affected by chromosomal
copy number change. Alternatively, if closely related cells
resemble each other to the same degree as any other cells,
regardless of ancestry, the phenotype will appear plastic,
not heritable and not auto-correlated. Such a pattern could
reflect temporally transient states such as cell-cycle phase.
Generally, phylogenetic auto-correlation captures the tem-
poral stability or transience of a cell state, whether state is
defined by intrinsic (*e.g.*, mutation) or by extrinsic factors
(*e.g.*, interactions with the microenvironment). For exam-
ple, if there is rapid toggling between states within a single
generation, these states likely will not be auto-correlated in
phylogenetic space, in contrast to more stable cell states that
persist without transitioning for time scales longer than one
cell division. Furthermore, we can assess statistical signif-
icance by computing phylogenetic correlation z scores, ei-
ther analytically [Czaplewski and Reich, 1993] or by using a
leaf-permutation test (**Methods: Phylogenetic correla-**
tions). By measuring phylogenetic auto-correlations, PATH
provides a powerful framework for quantifying the temporal
stability and thus heritability versus plasticity of somatic
cell states (or phenotypes) using multi-omic platforms that
jointly capture the lineage history and the cell state of single
cells.

295 In addition to quantifying the lineage dependency of single
 296 cell states to define heritability versus plasticity, to under-
 297 stand the evolutionary relationships *between* cell states we
 298 measure *phylogenetic cross-correlations* (**Methods: Phy-**
 299 **logenetic correlations**). Phylogenetic cross-correlation
 300 quantifies the dependency of one cell state's distribution on
 301 the lineage patterning of another state. For example, again
 302 consider the phylogenetic distribution of a phenotype that
 303 depends on chromosomal copy number. If a chromosomal
 304 duplication occurs, cells with the extra chromosome, and
 305 affected phenotype, will be in close phylogenetic proximity
 306 to each other, and farther from cells without the chromo-

somal duplication. As such, each of the phenotypes, one
 307 affected and one unaffected by the duplication, will be auto-
 308 correlated, but because these phenotypes will be phyloge-
 309 netically segregated from each other they will be negatively
 310 cross-correlated. On the other hand, if distinct measure-
 311 ments co-cluster phylogenetically, such as the transcription
 312 levels of two genes located on a chromosomal copy vari-
 313 ant, such measurements will be positively cross-correlated.
 314 The phylogenetic cross-correlation of a cell state with it-
 315 self is also its auto-correlation, so to simplify terminology
 316 when possible, we refer to both phylogenetic auto- and cross-
 317 correlations as *phylogenetic correlations*.
 318

Figure 1

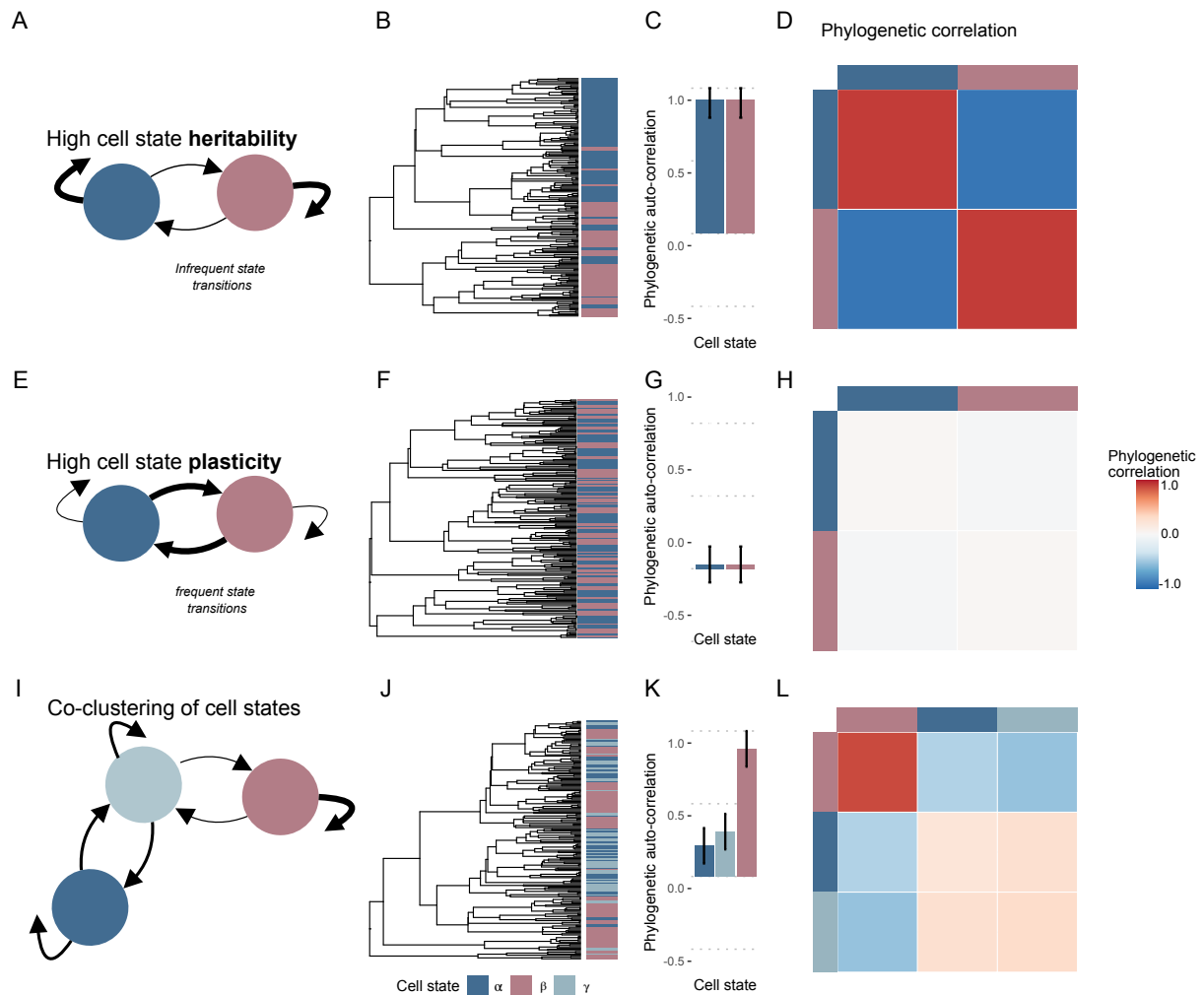


Figure 1: Phylogenetic correlations quantify the heritability versus plasticity of single-cell phenotypes

- A)** Diagram of highly heritable (categorical) cell state transition dynamics (**Methods: Markov model of cell state transitions**). Markov transition probabilities between states were simulated as $P_{\alpha\alpha} = P_{\beta\beta} = 0.9$, and $P_{\alpha\beta} = P_{\beta\alpha} = 0.1$ (meaning that cells had a 10% probability of switching states over each time point).
- B)** Phylogenetic tree containing 200 cells, simulated as a somatic evolutionary process (**Methods: Simulating phylogenies**), from simulated transition dynamics depicted in **A**, with birth rate = 1 and death rate = 0.
- C)** Phylogenetic auto-correlations (**Methods: Phylogenetic correlations**) for cell states depicted in **B**.
- D)** Phylogenetic cross-correlation (**Methods: Phylogenetic correlations**) heat map for cell states depicted in **B**. Diagonals are equivalent to bars shown in **C**.
- E)** Diagram of highly plastic (categorical) cell state transition dynamics (**Methods: Markov model of cell state transitions**). Markov transition probabilities between states were all the same ($P_{\alpha\alpha} = P_{\beta\beta} = P_{\alpha\beta} = P_{\beta\alpha} = 0.5$; meaning that cells had a 50% probability of switching states at any time).
- F)** Phylogenetic tree containing 200 cells, simulated as a somatic evolutionary process (**Methods: Simulating phylogenies**), from simulated transition dynamics depicted in **E**, with birth rate = 1 and death rate = 0.
- G)** Phylogenetic auto-correlations (**Methods: Phylogenetic correlations**) for cell states depicted in **E**.
- H)** Phylogenetic cross-correlation (**Methods: Phylogenetic correlations**) heat map for cell states depicted in **F**.
- I)** Diagram of a three-state system (**Methods: Markov model of cell state transitions**) in which states α and γ transition to each other at a rate higher than either transitions to state β . Markov transition probabilities between the three states were $P_{\alpha\alpha} = P_{\alpha\gamma} = P_{\gamma\gamma} = 0.5$, $P_{\alpha\beta} = P_{\beta\alpha} = 0$, $P_{\gamma\alpha} = 0.45$, $P_{\beta\gamma} = 0.1$, $P_{\gamma\beta} = 0.05$, and $P_{\beta\beta} = 0.9$.
- J)** Phylogenetic tree containing 200 cells, simulated as a somatic evolutionary process (**Methods: Simulating phylogenies**), from simulated transition dynamics depicted in **I**, with birth rate = 1 and death rate = 0.
- K)** Phylogenetic auto-correlations for cell states depicted in **J**.
- L)** Phylogenetic cross-correlation (**Methods: Phylogenetic correlations**) heat map for cell states depicted in **J**.
- Error bars in **C**, **G**, and **K** represent the analytical phylogenetic auto-correlation standard deviations calculated with the method from Czaplewski and Reich [1993].

To illustrate PATH, **Figure 1** depicts phylogenies that are the result of simulations of somatic evolution (**Methods: Simulating phylogenies**), where cells can transition between states. When cell states are heritable, meaning that state transitions occur infrequently (**Fig. 1A**), cells appear to phylogenetically group by state (*e.g.*, **Fig. 1B**), and thus states are positively auto-correlated and negatively cross-correlated (**Fig. 1C,D**). In contrast, for highly plastic dynamics where state transitions occur frequently (**Fig. 1E**), cells do not appear to phylogenetically group by state (*e.g.*, **Fig. 1F**), and states are lowly phylogenetically auto- and cross-correlated (**Fig. 1G,H**). The phylogenetic correlations between states can reflect evolutionary relationships; phylogenetic correlations increase or decrease with between-state transitions rates. For example, since transitions between state α and γ occur more frequently than transitions to β (**Fig. 1I**), α and γ co-cluster on the phylogeny (**Fig. 1J**) and are more phylogenetically correlated with each other than with β (**Fig. 1K,L**). Note that despite focusing on categorical cell states in **Figure 1**, phylogenetic correlations can also be computed for quantitative phenotypes (*e.g.*, gene expression level).

We hypothesized that as cell state phylogenetic patterning can be related to the rate of state transitions (as in **Fig-**

ure 1), the rates of these state transitions might be inferred from such patterns. To test this, we simulated categorical state transition dynamics on idealized phylogenies (*i.e.*, completely sampled and balanced, where every node has the same number of progeny; **Methods: Simulating phylogenies, Fig. S1A**). First, we confirmed a strong association between simulated transition rates and phylogenetic correlations (**Fig. S1B**, Spearman's $\rho = 0.89$). Next, we explicitly connected phylogenetic correlations with a mathematical model of state transition rates (**Methods: Phylogenetic correlations and cell state transitions, Box S1**). For categorical cell states, phylogenetic correlations characterize the frequencies at which states are found within cell pairs that share recent ancestry, and these frequencies can be anticipated given a model of state transitions. For example, the states found within a pair of sister cells will depend on the state of the sisters' shared parent and the rates at which transitions to other states can occur. For a highly heritable cell state in which transitions to other states occur infrequently, we will observe more sister cell pairs in the same such state than what we would expect given the state's frequency. Using this mathematical relationship we can transform phylogenetic correlations into transition rate estimates with high accuracy (**Methods: Inferring cell state transitions from phylogenetic correlations, Fig-**

397 **S1C, Box S1).**

398 **Measuring heritability, plasticity, and cell** 399 **state transition dynamics in somatic evolu-** 400 **tion**

401 The study of somatic evolution requires addressing an array
402 of complicating biological and technical features not repre-
403 sented by idealized phylogenies (*e.g.*, **Fig. S1A**). For in-
404 stance, when cell division is not synchronized within a pop-
405 ulation [Brody et al., 2018], meaning that different cell gen-
406 erations coexist, the resultant phylogenies will be more ad-
407 equately modeled in continuous-time. Additionally, not all
408 cells will leave the same number of progeny, resulting in less
409 balanced phylogenies. Moreover, in experimental contexts,
410 not all cells are successfully assayed, leading to incomplete
411 sampling. Other technical factors, such as sequencing depth
412 or barcode length, can limit the detection or accumulation of
413 heritable markers necessary to resolve close phylogenetic re-
414 lationships. As such, to test the robustness of PATH across
415 a wide range of biological and technical factors, we applied
416 PATH to phylogenies simulated with a more sophisticated
417 model of somatic evolution [Louca, 2020, Nee et al., 1994]
418 (**Methods: Simulating phylogenies**). In this model, cell
419 division and death occur, each with some probability, until
420 the population reaches a chosen size. Then only a fraction
421 of surviving cells is sampled and lineage relationships recov-
422 ered. Cell states are simulated along the sampled phyloge-
423 nies using a Markov model (**Methods: Markov model**
424 **of cell state transitions**). Cell division, death, sampling,
425 and state transition rates can be specified, thus providing a
426 more accurate representation of somatic evolution to assess
427 PATH’s applicability to complex somatic evolution datasets.

428 Consistent with our observations on idealized phylogenies
429 (**Figure S1**), in phylogenies produced by this sampled so-
430 matic evolutionary process, phylogenetic correlations remain
431 strongly related to cell state transitions. For instance, auto-
432 correlation, our measure of heritability, declines as state
433 transitions become more frequent. However, in addition to
434 declining with plasticity, phylogenetic auto-correlations also
435 decrease as sampling becomes sparser (**Fig. 2A**), under-
436 estimating heritability. Here, heritability is underestimated
437 because incomplete sampling leads to an overestimation of
438 lineage proximity in terms of node distance (**Fig. 2B**). In
439 other words, cells that may appear to be close relatives on
440 the tree (*e.g.*, separated by one node) may in fact be more
441 distant relatives due to the loss of unsampled intermediates
442 (due to cell death, incomplete sampling or incomplete phy-
443 logenetic reconstruction). As such, when sampling is low,
444 as might be the case when only hundreds or thousands of
445 cells from a tumor are collected, even the closest related
446 sampled cells from such lineages will usually represent fairly
447 distant relationships, thus affecting heritability estimates.
448 In these cases, only highly heritable phenotypes, reliably
449 propagated over the number of cell divisions separating the

450 closest related sampled cells will be detectable. These data
451 reveal that under sufficiently sparse sampling, heritable phe-
452 nototypes may appear plastic.

453 Next, we used PATH to infer state transition dynamics on
454 phylogenies simulated by the sampled somatic evolutionary
455 process. Since our inference approach transforms heritabil-
456 ity measurements – which are underestimated when sam-
457 pling is low – into transition rate estimates, transition in-
458 ference accuracy was highest when state heritabilities were
459 detectable (state auto-correlation z scores > 2 , **Fig. 2C,D**,
460 insets depict inferences for simulations in which heritabil-
461 ity was not detectable [z score ≤ 2]). Notably, transi-
462 tion inference accuracy (**Methods: Assessing cell state**
463 **transition inference accuracy**) with PATH is comparable
464 to state-of-the-art Maximum Likelihood Estimation (MLE)
465 methods (as implemented in Louca and Doebeli [2018]) tra-
466 ditionally used in evolutionary biology to infer character
467 transitions (**Fig. 2E, Fig. S2A,B**), but with signifi-
468 cantly faster compute times when analyzing a large number
469 of states (**Fig. 2F, Fig. S2C**) and/or cells (**Fig. 2G,**
470 **Fig. S2C**). PATH’s relative speed derives from the fact
471 that PATH transforms a statistic (phylogenetic correlation)
472 into a transition probability, whereas MLE uses an optimiza-
473 tion algorithm to search for the most likely state transition
474 probabilities and often requires many more calculations.

475 Another important confounder in harnessing phylogenetic
476 trees to measure heritability is the fidelity of phylogenetic
477 reconstruction. Intuitively, this can be understood in the
478 context of artificial lineage tracing techniques that stochas-
479 tically scar or cut genetic barcodes (*e.g.*, Molecular recorder
480 [Chan et al., 2019] and scGESTALT [Raj et al., 2018]), where
481 a limited number of cut sites can result in phylogenetic re-
482 construction errors. To understand this, beyond simulating
483 phylogenies as a sampled somatic evolutionary process, we
484 also simulated the reconstruction of these phylogenies by
485 employing a model of CRISPR/Cas9 scarring (**Methods:**
486 **Phylogenetic reconstruction**). To do this, each cell in a
487 simulated evolving population contains a *barcode*, or a set
488 of mutable and heritable sites that can be modified (*i.e.*,
489 scarred) stochastically. In contrast to our previous approach
490 in which true phylogenies were recovered, here phylogenies
491 were reconstructed from the differences between barcodes re-
492 trieved from cells in the terminal population, much as they
493 would be for lineage tracing experiments. Comparing re-
494 constructed with true phylogenies, we observe that as the
495 number of mutable sites or barcode length increases, phy-
496 logenetic reconstruction accuracy improves (**Fig. S2D**).
497 Concordant with reconstruction accuracy, state transition
498 inferences using PATH also improve (**Fig. 2H**).

499 Since the accuracy of state transition inferences using PATH
500 is affected by reconstructed branch lengths, which scale phy-
501 logenetic distances by time, inference will be impeded when
502 branch lengths are inaccurate, and not possible when branch
503 lengths are absent (which is common for single-cell phyloge-

504 nies using artificial scarring methods). PATH can compen- 523
 505 sate for this by imputing terminal branch lengths, indepen- 524
 506 dent of phylogenies, if cell population sizes can be approxi- 525
 507 mated (**Methods: Inferring cell state transitions from** 526
 508 **phylogenetic correlations, Imputing branch lengths**). 527
 509 PATH achieves this because under the model of sampled 528
 510 somatic evolution, the degree by which sampling leads to 529
 511 an overestimate of phylogenetic proximity can be calculated 530
 512 (**Fig. 2B, Fig. S2E,F**) and accommodated. In other 531
 513 words, under incomplete sampling, in which close phylo- 532
 514 genetic relationships are overestimated due to the loss of 533
 515 unsampled intermediate cells, from the sampling rate (and 534
 516 independent of the reconstructed phylogeny), we can esti- 535
 517 mate how many intermediates are unsampled, and rescale 536
 518 branch lengths accordingly. Replacing measured branch 537
 519 lengths with model-imputed lengths significantly improves 538
 520 the accuracy of state transition inferences using PATH, par- 539
 521 ticularly for low fidelity phylogenetic reconstructions where 540
 522 branch lengths are often less accurate (**Fig. 2H**). Thus, us-

ing PATH, state transitions can be accurately inferred for 523
 low fidelity phylogenies and when branch lengths are absent 524
 (in contrast to MLE), making PATH a powerful tool for 525
 the analysis of phylogenies produced by molecular scarring 526
 technologies. 527

In conclusion, these simulated datasets demonstrate that 528
 PATH, through the measurement of phylogenetic correla- 529
 tions, provides a comprehensive framework to analyze cell 530
 state heritability and plasticity in somatic cell populations, 531
 and can transform these measurements into inferences of 532
 state transition dynamics. PATH can accommodate a wide 533
 range of biological and technical features associated with 534
 somatic evolution. Thus, observable patterns of heritability 535
 and plasticity are robustly linked to the (often unobservable) 536
 processes that produce them, providing insights into cell lin- 537
 eage histories and somatic evolutionary dynamics. Having 538
 explored PATH's capabilities on simulated datasets, we next 539
 sought to apply PATH to published single-cell lineage trac- 540
 ing datasets in two broad contexts, development and cancer. 541

Figure 2

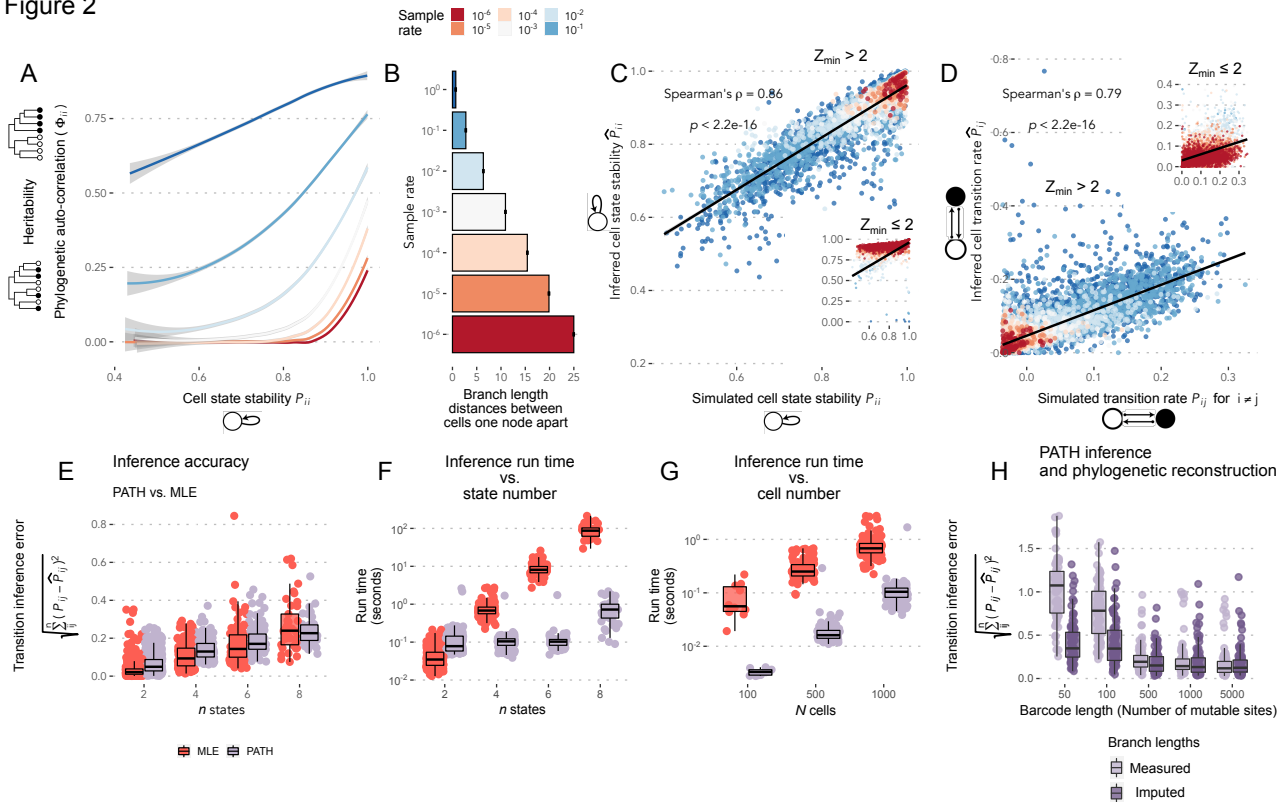


Figure 2: Measuring heritability, plasticity, and cell state transition dynamics in somatic evolution

A) Simulated cell state stability (Markov self-transition probability, **Methods: Markov model of cell state transitions**) for state 1 versus measured phylogenetic auto-correlation under different sampling rates (**Methods: Phylogenetic correlations**). Phylogenies contain 1,000 cells and Markov cell state transition dynamics were randomly generated for three-state systems. Phylogenies simulated as a sampled somatic evolutionary process (**Methods: Simulating phylogenies**) with birth rate 1 and death

rate 0. Lines colored by sampling rate depict LOESS regression lines with 95% confidence intervals (light gray). 550

B) Mean branch length (in units of time) distance between cell pairs only one-node apart on phylogenies versus cell sampling rate for phylogeny simulations. 551 552

C) Simulated cell state stability (Markov self-transition probability) for state 1 versus PATH-inferred state stability for systems with phylogenetic auto-correlation z scores > 2 . Colors represent sampling rates. Inset shows systems with at least one phylogenetic auto-correlation z score ≤ 2 , and uses the same regression line. 553 554 555

D) Simulated versus PATH-inferred cell state transition probability from state 1 to state 2 for three-state systems with phylogenetic auto-correlation z scores > 2 . Colors represent sampling rates. Inset shows systems with at least one phylogenetic auto-correlation z score ≤ 2 , and uses the same regression line. 556 557 558

E) Comparing the state transition dynamic inference accuracy of PATH (light purple) with Maximum Likelihood Estimation (MLE; orange). Inference error is calculated as the Euclidean distance between inferred and simulated transition probability matrices (equation shown on y-axis label), and the number of possible states in a simulated system is shown on the x-axis (**Methods: Assessing cell state transition inference accuracy**). Panel depicts simulations for 1,000 cell phylogenies, sampled at a rate of 10^{-2} , excluding simulations in which either inference method failed (which were usually due to the complete absence of some cell states). 559 560 561 562 563 564

F) Same as **E** but measuring compute time. 565

G) Comparing PATH and MLE compute times while varying phylogenetic tree size (number of cells; x-axis) fixing systems to four cell states, and sampled at 10^{-2} . All inferences filtered to simulations surpassing the minimum phylogenetic auto-correlation z score threshold of 2. 566 567 568

H) Comparing state transition inference of PATH using two different node depth estimation methods: (light purple) using measured branch length distances, and (dark purple) using imputed branch lengths (**Methods: Imputing branch lengths**) from estimated cell sampling rates. Simulations are for three-state systems simulated on 1,000 cell sampled somatic evolutionary phylogenies (**Methods: Simulating phylogenies**). Phylogenies were reconstructed by using the UPGMA algorithm on the cell pairwise Hamming distances between simulated lineage barcodes that were stochastically scarred at rate $s = 0.01$ (**Methods: Phylogenetic reconstruction**). 569 570 571 572 573 574

576 PATH quantifies ancestry and divergence of 577 germ layers and cell types during mouse em- 578 bryogenesis

579 Embryogenesis and organogenesis require the organization 580 of the progeny of progenitor cells, which are restricted in 581 number, location and levels of potency, into complex tissues. 582 Single-cell lineage tracing methods provide sufficient resolu- 583 tion to map the cellular trajectories and interactions that 584 underlie this exquisitely regulated organization. We reason- 585 ed that the application of PATH to such datasets would 586 enable quantification of cell differentiation patterns through 587 calculation of (i) phylogenetic auto-correlations that can be 588 interpreted in this developmental context as cell state com- 589 mitment strength and (ii) phylogenetic cross-correlations to 590 determine relationships between tissue layers and cell types, 591 and to understand gene expression across development.

592 We first asked whether PATH is able to reconstruct known 593 cell fate relationships and dynamics in the well-characterized 594 context of murine gastrulation (**Fig. 3A**). To accomplish 595 this, we applied PATH to published mouse embryogenesis 596 data [Chan et al., 2019], comprising single-cell phyloge- 597 nies with matching single-cell transcriptional data. The au- 598 thors leveraged a CRISPR/Cas9 lineage tracing construct to 599 study early murine development, isolating embryos at E8.5 600 and constructing phylogenies from the edited barcodes (**Fig. 601 3B, Fig. S3A**). We applied PATH to these data to measure

576 phylogenetic correlations for cellular phenotypes at multiple 577 levels of resolution, and gained insight into the commitment 578 and divergence patterns of cellular phenotypes from their 579 origin layers in the blastocyst through gastrulation, and ul- 580 timately to their differentiated tissue in the E8.5 embryo. 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601

602 As expected, all blastocyst layers with sufficient representa- 603 tion had high auto-correlation in both replicates, indicating 604 that a cell from a particular blastocyst layer is more likely 605 to produce progeny that are also found in the same layer, re- 606 inforcing what is known about the rigidity of developmental 607 programs [Thowfeequ and Srinivas, 2022]. Germ layers de- 608 rived from outside of the epiblast had high auto-correlation 609 in all replicates that had sufficient cell recovery, while tissues 610 that shared a common origin in the epiblast had lower auto- 611 correlations (**Fig. S3B**). Thus, the non-epiblast-derived 612 layers show evidence of earlier fate commitment, while the 613 more plastic phenotype of the epiblast is consistent with 614 its later divergence [Thowfeequ and Srinivas, 2022]. PATH 615 also accurately reconstructed the patterns of shared ancestry 616 between blastocyst layers and germ layers (**Fig. 3C**). No- 617 tably, phylogenetic correlations recovered the dual contribu- 618 tion of both embryonic- and extraembryonic-derived tissues 619 to the endoderm [Kwon et al., 2008, Nowotschin et al., 2019, 620 Pijuan-Sala et al., 2019] (**Fig. 3C**). This highlights PATH's 621 ability, by leveraging phylogenies, to identify phenotypically 622 similar but ancestrally distinct cells. 623 624 625 626 627

Figure 3

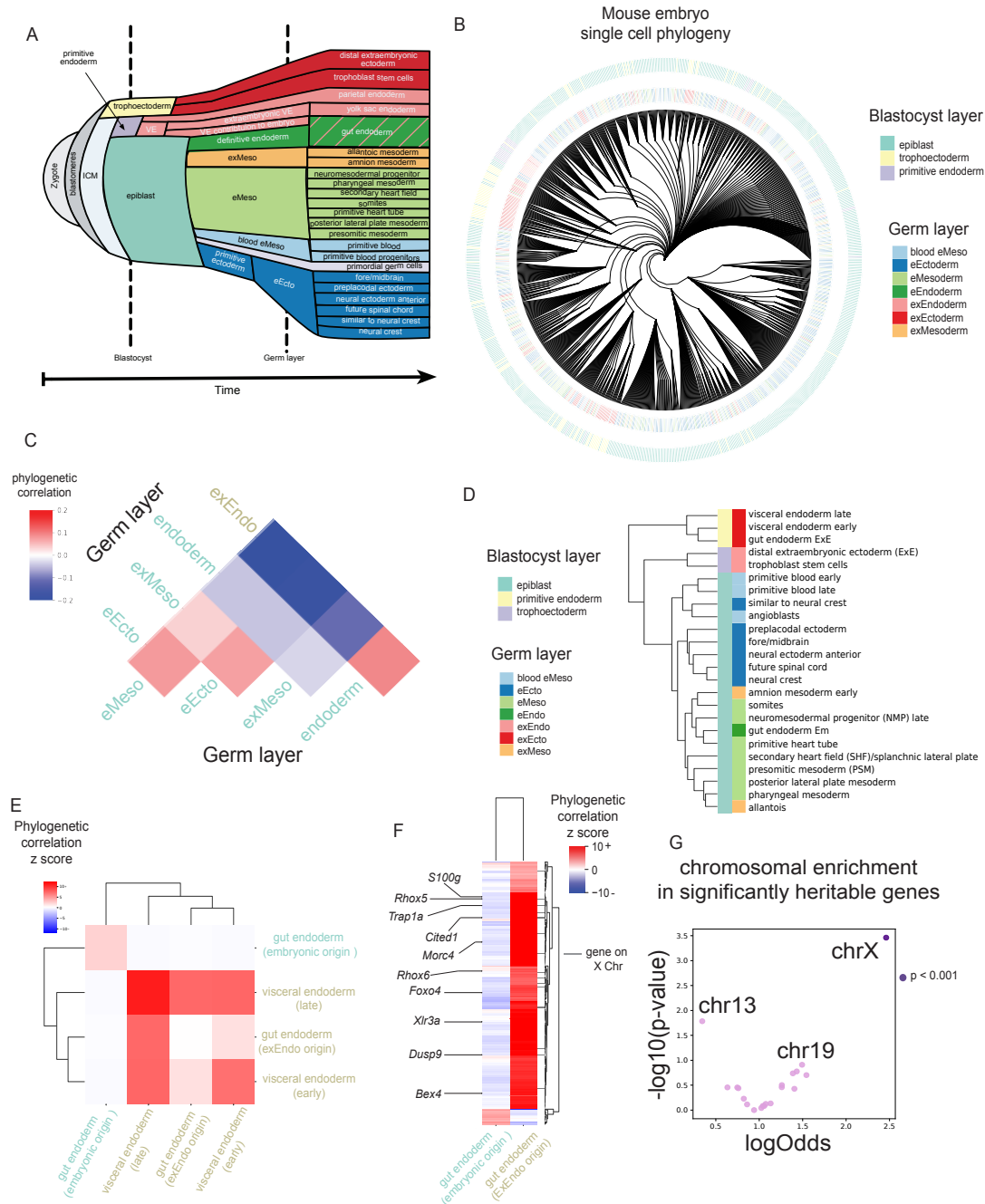


Figure 3: PATH quantifies ancestry and divergence of germ layers and cell types during mouse embryogenesis

A) Schematic of mouse embryogenesis adapted from Thowfeeq and Srinivas [2022]. VE, visceral endoderm; ICM, inner cell mass; e prefix, embryonic; ex prefix, extraembryonic.

B) Single-cell phylogeny from mouse embryo 6 from Chan et al. [2019], containing 700 randomly chosen of 1,722 cells for visualization. Each leaf represents a single cell. Leaves are colored by blastocyst or germ layer of origin. e prefix, embryonic; ex prefix, extraembryonic.

C) Germ layer phylogenetic correlations for embryo 2. Labels colored by cell type blastocyst origin: visceral endoderm, gold;

epiblast, green.

D) Hierarchical clustering of tissue types by phylogenetic correlation using Ward's method. Only tissues with more than 30 cells were used. Tissues colored by germ and blastocyst layer of origin. Phylogenetic correlations can be found in **Fig. S3C**. ExE, extraembryonic; EM, embryonic.

E) Phylogenetic correlation z score of gut endoderm cells annotated by their source tissue in the blastocyst and visceral endoderm (early and late). Labels colored by cell type blastocyst origin: visceral endoderm, gold; epiblast, green.

F) Phylogenetic correlation z scores between genes and tissue assignment. Genes on the X chromosome are denoted with a gray bar (right) with select X-chromosome genes labeled (left). Cell state labels colored by cell type blastocyst origin: visceral endoderm, gold; epiblast, green. The complete set of phylogenetic correlations are in **Table S1**.

G) Enrichment of highly heritable genes at the whole chromosome level (with chromosome 13, 19 and X labeled). Log odds ratio and p-value ($p < 10^{-3}$, Fisher's exact test) of number of highly heritable genes (z score > 3) on each chromosome compared to all other chromosomes Only expressed genes were considered for comparison (top 2,000 most variable genes across phylogeny, see **Methods: Mouse embryogenesis**).

After implementing PATH at the level of the blastocyst and germ layers, we sought to quantify the degree of shared origin of higher resolution, transcriptionally defined cell types derived from each germ layer (**Fig. 3D**). Cell types that share ancestry will likely be highly phylogenetically correlated. Indeed, PATH analysis correctly identified important developmental relationships between primitive blood cells (early and late); and neural crest and future spinal cord. Interestingly, PATH also identified the shared origins of the embryonic splanchnic lateral plate and extraembryonic allantois cells in the nascent mesoderm [Thowfeequ and Srinivas, 2022], highlighting PATH's ability to identify shared ancestry from progeny that have diverged into different germ layers (**Fig. S3C,D**). Of note, we again observed high cross-correlation between the endoderm and extraembryonic endoderm-derived tissues in the gut endoderm (**Fig. 3C**), now at the level of cell type (**Fig. 3E**). This higher resolution analysis revealed that extraembryonic-derived endoderm tissue cross-correlates almost exclusively with cells from the late visceral endoderm (arising around E8.0 in the extraembryonic endoderm), as opposed to the early visceral endoderm (arising around E7.0 in the extraembryonic endoderm) [Grosswendt et al., 2020] or embryonic-derived gut endoderm. Given that the intercalation of extraembryonic endoderm into the gut endoderm occurs between E7.5 and E8.5 [Nowotschin et al., 2019], this analysis nominates a specific cell population from the extraembryonic visceral endoderm contributing to the definitive endoderm.

Having examined the phylogenetic correlations of embryonic germ layers and cell types, we then took advantage of the versatility of PATH to evaluate the heritability of gene expression programs in these populations of endoderm cells. We calculated phylogenetic correlations between each population of endoderm cells (originating in the epiblast or the primitive endoderm) and gene expression across the tree. We found distinct gene expression profiles phylogenetically correlated with each population of endodermal cells (**Fig. 3F**). In concordance with prior work, we found that

Rhox5 and *Trap1a*, two X-linked genes, had high phylogenetic correlation with endoderm cells with extraembryonic origin [Nowotschin et al., 2019, Pijuan-Sala et al., 2019]. Interestingly, we found that genes on the X chromosome beyond *Trap1a* and *Rhox5* were significantly enriched in this heritable expression program (**Fig. 3F,G**). This signal is grounded in the differential imprinting patterns between extraembryonic and embryonic cells: extraembryonic endoderm cells have paternally imprinted X-inactivation [Takagi and Sasaki, 1975] imbuing them with a unique expression pattern that has been shown to persist after intercalation into the visceral endoderm [Loda et al., 2022]. These results demonstrate PATH's ability to explore patterns and timing of coordinated gene expression during development, including epigenetically propagated signals.

PATH identifies cell fate-determining factors across anatomical, defined tissue and gene expression layers during neurogenesis in zebrafish

One notable aspect of PATH is its ability to quantify relationships between different types of phenotypic information, providing the opportunity to leverage not only transcriptional information from scRNAseq data, but also any available spatial, anatomical or temporal information. As such, we can perform multi-modal analysis to characterize relationships between these phenotypic annotation layers, and thus draw inferences about their interactions (for example, we can use the phylogenetic cross-correlations of individual genes with either cell or tissue type to nominate cell fate determination factors). To explore this capability, we applied PATH to prospectively lineage-traced developing zebrafish brains [Raj et al., 2018]. The data in Raj et al. [2018] comprise cells annotated not only by single-cell transcriptional profiling but also by the anatomic region from which they were dissected. These multi-layer annotations enabled us to investigate neuronal development dynamics within, between

725 and across anatomically distinct brain regions.
 726 We first used PATH to examine phylogenetic correlations of
 727 different brain regions. Neuronal tissue had been collected
 728 from two whole brains and anatomic regions were manually
 729 separated during dissection, resulting in three main regions
 730 (forebrain, midbrain, hindbrain; **Fig. 4A,B**). By projecting
 731 anatomic region on the reconstructed phylogeny and apply-

ing PATH, we found that each defined anatomic location
 had high phylogenetic auto-correlation, indicating that neu-
 732 ronal cells within a brain region share recent ancestry
 733 (Fig 4C). As expected, the cells with ambiguous annotations (la-
 734 beled “mix”) had much lower phylogenetic auto-correlations,
 735 most likely due to heterogeneous sampling that diluted the
 736 phylogenetic signal.
 737
 738

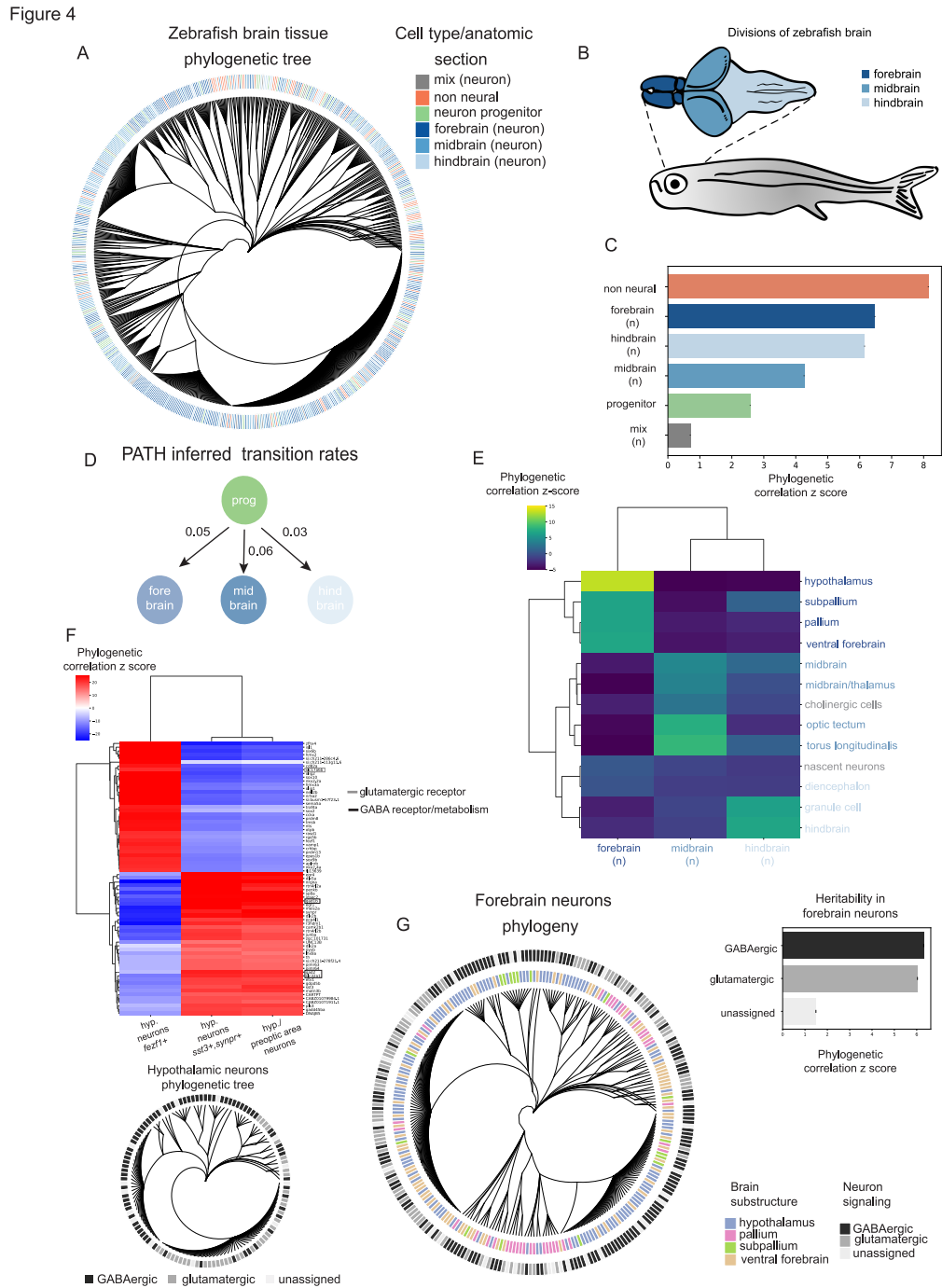


Figure 4: PATH identifies cell fate-determining factors across anatomical, defined tissue and gene expression layers

during neurogenesis in zebrafish

A) Single-cell phylogeny from zebrafish brain 3 (replicate 1) from Raj et al. [2018]. Each leaf represents a single cell ($N = 750$). All cell type and anatomic section annotations are as defined in Raj et al. [2018], by scRNAseq and manual dissection, respectively. Cells colored in orange are non-neurons, cells in green are neural progenitors. Neuronal cells (blue hues and gray) are colored by the anatomic location from which they were dissected. Non-neural and neuron progenitor cells lack anatomical annotation. Cells labeled “mix” were from dissections with ambiguous anatomical origin (see **Methods: Zebrafish brain development**).

B) Zebrafish brain schematic. Forebrain, midbrain and hindbrain have been labeled.

C) Cell type/anatomic-section phylogenetic auto-correlations. Mature neurons are labeled “n” and annotated by dissection site (blues, gray); neuronal progenitors are labeled in green and non-neural cells are in orange.

D) PATH inferred transition probabilities between neuron progenitor cells (prog) and neurons from each anatomic brain region. Branch lengths imputed by approximating the cell sampling rate to be 10^{-4} to infer transition probabilities. Values rounded to the nearest hundredth.

E) Phylogenetic correlation z scores between anatomic site and transcriptionally assigned brain substructure across all neurons. Substructures are colored by brain location from **A**.

F) Phylogenetic correlation z scores between (top 2,000 most variably) expressed genes and individual hypothalamus clusters (defined by Raj et al. [2018] from select marker genes). The 35 most auto-correlated genes per cluster are shown, and a complete set of phylogenetic correlations are in **Table S2**. Phylogenetic tree of hypothalamic neurons annotated by GABA/Glut signaling (**Fig. S4C**) (see **Methods: Zebrafish brain development**).

G) (Left) phylogeny of all forebrain neurons ($N = 270$), leaves annotated by brain substructure assignment and GABA and glutamatergic signaling. (Right) phylogenetic auto-correlation of GABA and glutamatergic signaling across all forebrain neurons.

To characterize potential developmental trajectories between neurons and neuronal progenitors, we next used PATH to infer transition dynamics between them, segregating neurons by their anatomic region. Notably, we found that the progenitor cell pool contributes at similar rates to the forebrain, midbrain and hindbrain (**Fig. 4D**), consistent with the findings of Raj et al. [2018] suggesting that progenitor cells were multipotent at the time of barcoding.

As the versatility of PATH allows not only for comparisons within the same category of data (*e.g.*, brain region), but also for integrated analysis across different layers of phenotypes, we next aimed to examine the phylogenetic correlation of anatomical brain regions with higher-resolution brain structure information derived from scRNAseq marker data. PATH analysis showed that these brain structures cross-correlate with their expected anatomical region (**Fig. 4E**), demonstrating the ability to correctly integrate transcriptionally and anatomically derived single-cell annotations across a phylogeny.

We next focused our analysis on the hypothalamus, a complex brain structure that is essential for the maintenance of homeostasis in an organism’s adaptive response to its environment. This structure is composed of a variety of anatomically and molecularly distinct neuron subtypes which respond to and release distinct sets of neuropeptides and hormones [Benevento et al., 2022]. Given this complexity, the transcriptional and phylogenetic dynamics underlying the functional organization of the hypothalamus were of interest for us to explore within the PATH framework. Using gene

clusters defined by Raj et al. [2018] using scRNAseq, we first assessed the phylogenetic correlations of transcriptionally distinct clusters (**Fig. S4A**) of hypothalamic neurons. This analysis showed that *tac1+*, *nrgna+*, neurons were highly cross-correlated with neurons from the preoptic area (POA), indicating a shared cellular ancestry. The expression of both of these genes was negatively cross-correlated with *fezf1+* neurons, indicating distinct histories (**Fig. S4A**). To explore the molecular underpinnings of these differences in developmental origins we cross-correlated gene expression with hypothalamic neuron subtype (**Fig. S4A**) across the phylogeny of forebrain neurons to determine which genes were most strongly cross-correlated with these cell types (**Fig. 4F**). Interestingly, we found that genes required for glutamatergic signaling (*slc17a6b*) were highly cross-correlated with *fezf1+* neurons, while those genes required for GABAergic signaling (*gad1b*, *gad2*, *slc32a1*) were highly cross-correlated with POA and *tac1+*, *nrgna+*, neurons, indicating that use of GABAergic or glutamatergic signaling is a heritable trait in cells of the differentiating hypothalamus (**Fig. 4F**). Indeed, we found that glutamatergic and GABAergic signaling were heritable in the forebrain (**Fig. 4G, Fig. S4B,C**), consistent with lineage tracing studies that found high heritability of GABAergic signaling in the murine forebrain [Bandler et al., 2021]. Thus, PATH is able to connect gene expression profiles to cell state through lineage information in an unbiased, quantitative manner, and uncovers the contribution of biologically meaningful cell populations underlying the observed patterns of heritability.

Quantifying cell state transitions during metastasis

Malignant populations harbor significant cell state diversity and the characterization of their relative heritability and plasticity is currently a major goal of the cancer field [Bell et al., 2019, Fennell et al., 2022, Oren et al., 2021, Shaffer et al., 2020]. Tumor single-cell phylogenies provide a unique opportunity to distinguish between cancer cell state heritability versus plasticity. Cancer cell state diversity has been associated with critical disease aspects such as tumor growth [Nefitel et al., 2019], treatment response [Fennell et al., 2022], and metastatic spread [Karras et al., 2022], emphasizing the need to define the heritability versus plasticity of cancer cell states. Notably, in comparison to primary tumors, in most contexts there is a lack of established, recurrent genetic drivers of metastasis [Rogiers et al., 2022]. Thus, other non-genetic factors likely play a major role in metastasis. We therefore applied PATH to correlate lineage dynamics with key non-genetic features, including location and cell state, of metastatic tumors. We re-analyzed data from a murine model of metastatic pancreatic cancer with inducible CRISPR/Cas9 based lineage recording and scRNAseq [Simeonov et al., 2021]. Metastatic tumors are thought to arise by the dissemination of a single or a small number of clones from the primary tumor [El-Kebir et al., 2018, Gudem et al., 2015, Hu et al., 2019, Navin et al., 2011, Turajlic et al., 2018]. By leveraging PATH's ability to integrate data of different modalities, we tested this assumption by assessing the shared ancestry of metastatic tumor cells harvested from distinct anatomical sites: primary tumor (pancreas), lung metastatic tumor, liver metastatic tumor, peritoneal metastatic tumor, tumors forming at the site of the surgical lesion and circulating tumor cells (CTCs). Cellular tissues of origin were highly phylogenetically auto-correlated (Fig. 5A,B), consistent with the established model in which a small number of founder cells seed metastases, creating site-specific clonal bottlenecks. Importantly, the quantification provided by PATH allowed for direct comparison of harvest site-specific lineages, revealing patterns of clonal seeding in metastasis. For instance, surgical lesions (which formed on the peritoneal surgical incision site) and peritoneal metastases had negative phylogenetic correlation, (Fig. S5A) suggesting that they had distinct origins despite their physical proximity. As expected, CTCs, which may have many distinct clonal origins, had lower phylogenetic auto-correlation than solid tissues (Fig. 5B).

The epithelial-to-mesenchymal transition (EMT) plays a crucial role in metastasis [Dongre and Weinberg, 2019, Lambert et al., 2017, Thiery, 2002], and thus Simeonov et al. [2021] calculated an EMT score for each tumor cell, reflective of that cell's position along a transcriptional continuum from highly epithelial to mesenchymal cells. Low scores correspond to more epithelial characteristics and high scores correspond to more mesenchymal characteristics. Of note, there is an ongoing discussion in the field regarding

whether EMT is best modeled as a series of functionally discrete, transcriptionally and epigenetically distinct intermediate states or a continuum of transcriptional hybrid states [McFaline-Figueroa et al., 2019, Pastushenko and Blanpain, 2019, van Dijk et al., 2018]. Because we can simultaneously observe both cellular position within the EMT continuum and on the phylogeny, this dataset offers a unique opportunity to investigate this question (Fig. 5C).

First, phylogenetic auto-correlation revealed the high heritability of cellular position on the EMT transcriptional continuum (Fig 5D). This finding can be contrasted with phylogenetic auto-correlation measurements of cellular position within the cell cycle, which can serve as a negative control, as position within the cell cycle is not usually expected to depend on ancestry [Chaligne et al., 2021] (Fig 5C,D).

Next, we asked how heritability and plasticity varied across the EMT continuum. Cells had been assigned EMT scores ranging from 0, denoting a completely epithelial cell to > 30 denoting a completely mesenchymal cell [Simeonov et al., 2021]. We partitioned cells along the continuum using units of 1 (bin #1 includes cells with EMT scores from 0 to 1, bin #2 includes cells from 1-2, etc.), merging bins at the extremes (all cells with a score of 7 or less were assigned to a single bin, as were cells that scored higher than 30) because these bins had low cellular representation. We calculated phylogenetic correlations for each individual bin, revealing four distinct groups of cross-correlated states along the EMT continuum defined by varying degrees of heritability (Fig. 5E; Fig. S5B,C, Table S3). Specifically, one group of phylogenetically correlated states corresponds to the epithelial and early transition states (T1), indicating that cells in this part of the EMT continuum tended to remain in the T1 state and were less likely to transition to other states. Likewise, mesenchymal (M) cells were also highly phylogenetically auto-correlated, indicating temporal stability of the mesenchymal state. However, cells in bins in the middle part of the continuum (later transition states; T2, T3) appeared less heritable, suggesting that these states were more plastic (Fig. 5E, Fig. S5B). These results were robust to different bin sizes (Fig. S5D), suggesting that these results are not an artifact of the binning procedure. Intriguingly, these results imply that despite tumor cells occupying a continuum of EMT transcriptional states, the states at the extremes of the continuum show a higher degree of heritability, whereas intermediate cells states show a higher degree of plasticity. As our analysis above showed a high degree of phylogenetic similarity within the same metastatic location, we further ruled out that EMT heritability is driven by variability in the representation of EMT states across metastatic sites (Fig. 5F). Furthermore, these results were replicated within each metastatic location, and consistently showed the T1 state to be the most heritable within each tissue, and the T2/T3 states to be more plastic, suggesting that patterns of cell state heritability were not driven by tumor location.

Figure 5

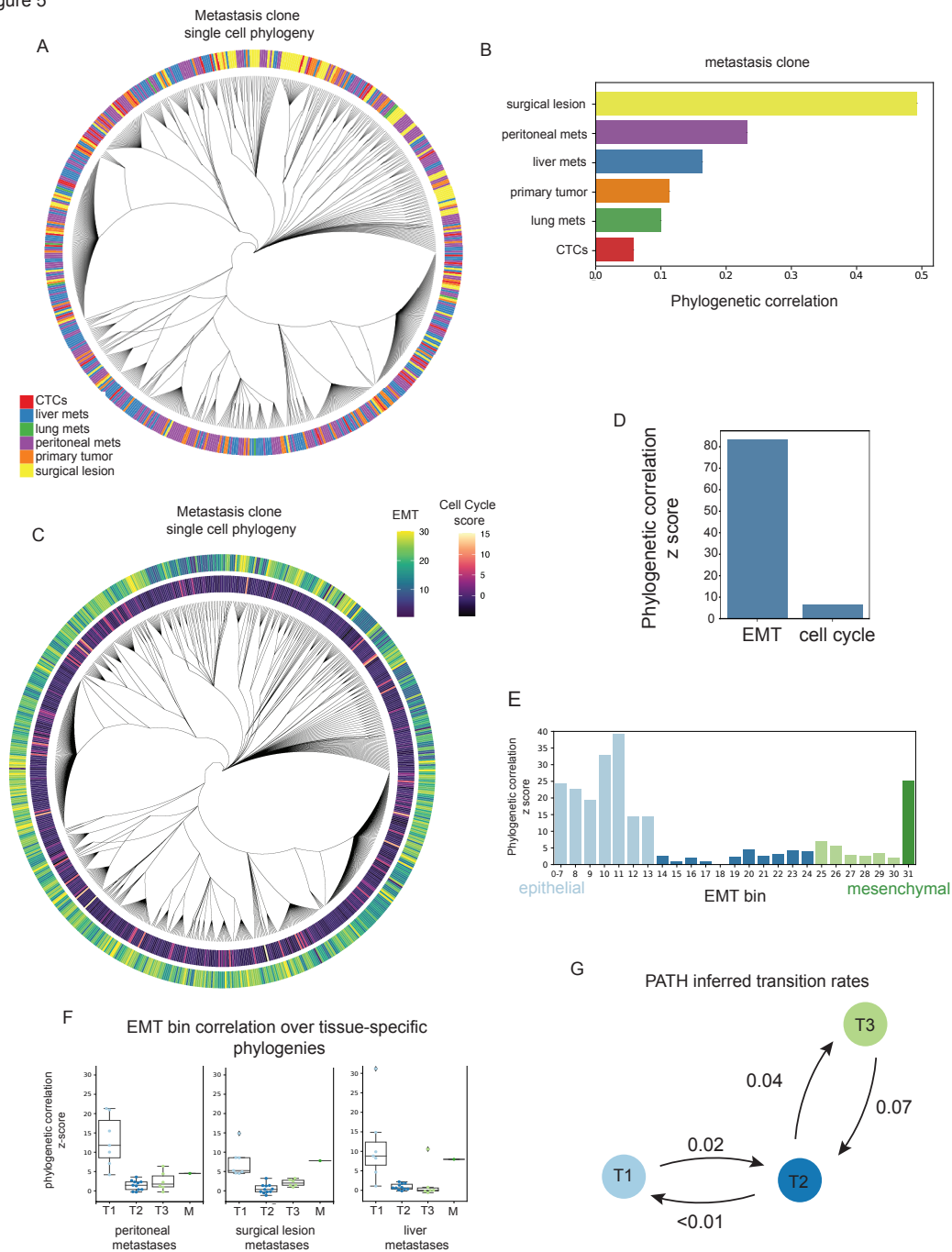


Figure 5: Quantifying cell state transitions during metastasis

A) Single-cell phylogeny from Mouse 1, Clone 1 from Simeonov et al. [2021], containing 700 randomly chosen of 7,968 cells for visualization. Each leaf represents a single cell. Leaves are colored by their harvest site. CTCs denote circulating tumor cells. Mets, metastases.

B) Phylogenetic auto-correlation of tumor cells annotated by harvest site. Bars colored by harvest site, as in **A**.

C) Single-cell phylogeny from **A**, with cells colored by EMT and cell cycle score (G2M score).

D) EMT and cell cycle phylogenetic auto-correlations across all tumor cells (N = 7,958).

- E)** EMT bin phylogenetic auto-correlations (z scores) using all cells. Bins are colored by transition states derived from **Fig. S5B**. 940
- F)** Box and whisker plot of EMT bin phylogenetic correlations (z scores) across phylogenies that contain cells from only one harvest site. Dots correspond to EMT bins. Bins are grouped and colored by transition state membership. Boxes represent the interquartile range (IQR); the center line represents the median; minima and maxima shown represent 1.5-IQR. 941 942 943
- G)** PATH inferred transition probabilities between states (T1, T2, T3) using all cells (N = 7,968). Values rounded to the nearest hundredth. Transition probability inferences use imputed branch lengths by approximating a sampling rate of 10^{-6} (see **Methods: Mouse model of pancreatic cancer**). 944 945 946

948 Finally, to quantify cell state transitions from the initial epithelial state to the more plastic later states, we used PATH 949 to infer transition dynamics between early (T1), middle (T2) 950 and late (T3) EMT states. We observed that transitions 951 out of the early epithelial state (T1) into more plastic states 952 along the continuum (T2) occurred with some frequency, 953 but transitions in the reverse direction going from a later 954 plastic state back to an early epithelial state were rare. In 955 contrast, we found marked plasticity between later intermediate 956 states (T2 and T3) (**Fig. 5G**). These results suggest 957 that EMT represents neither a smooth continuum of 958 hybrid states nor an equally discretized cell state trajectory, 959 but instead comprises punctuated states with different 960 transition probabilities. These analyses indicate an integration 961 of the two proposed models of EMT: cells undergoing 962 EMT are transcriptionally continuous (as reported by 963 [McFaline-Figueroa et al., 2019, Pastushenko and Blanpain, 964 2019, Simeonov et al., 2021, van Dijk et al., 2018]), but their 965 lineage dynamics reveal functionally and heritably distinct 966 states in EMT (as reported from functional transplantation 967 assays in mice by Pastushenko et al. [2018]). These findings 968 highlight the power of combining single-cell multi-omics 969 data with phylogenetic information to draw conclusions that 970 would not be possible through analyzing either data type 971 alone. 972

973 Elucidating heritable transcriptional modules and cell state transition dynamics in human glioblastoma 974 975

976 While artificial lineage tracing is a powerful approach in 977 model organisms, it cannot be applied to reconstruct phylogenetic 978 relationships in human data. Recent advances in multi-modal 979 single-cell sequencing enable joint lineage reconstruction and cell 980 phenotyping in primary human samples [Sankaran et al., 2022]. 981 To examine this exciting frontier, we applied PATH to phenotypically 982 annotated retrospective phylogenies reconstructed from human 983 single-cell data leveraging stochastic DNA methylation changes 984 as native lineage barcodes (**Methods: Human patient glioblastoma**) 985 [Chalighe et al., 2021, Gaiti et al., 2019]. 986

987 Having observed the high heritability of harvest site location 988 across multiple tumors in metastasis (**Fig. 5A,B**), 989 we set out to test whether a cell's spatial location within 990 a single tumor was stable. We applied PATH to MGH105,

an IDH-wildtype (WT) glioblastoma (GBM) patient specimen 991 in which cells were sampled from four distinct tumor 992 locations (**Fig. 6A**) [Chalighe et al., 2021, Neftel 993 et al., 2019]. We found that each of the locations (inset, 994 **Fig. 6A**) were highly phylogenetically auto-correlated (leaf- 995 permutation test, **Fig. 6B**), indicating that spatially proximal 996 tumor cells were also more proximal in terms of ancestry, 997 consistent with our expectations for a solid tumor malignancy. 998 999

GBM harbors significant cell state diversity, which can be 1000 classified according to the expression four major gene modules, 1001 defined as neural progenitor-like (NPC-like), oligodendrocyte 1002 progenitor-like (OPC-like), astrocyte-like (AC-like), and 1003 mesenchymal-like (MES-like) [Neftel et al., 2019]. By 1004 measuring transcriptional signatures for these modules in 1005 each cell, GBM cells can be classified into four distinct 1006 transcriptionally-defined cell states. These cell states can 1007 be further grouped by function; for instance, we define the 1008 stem-like cells as cells that highly express one of the progenitor 1009 (NPC- or OPC-like) gene modules. The stem-like and AC-like 1010 states each resemble a known neurodevelopmental program, 1011 and thus can be collectively considered as neurodevelopmental-like. 1012 In contrast, the MES-like state 1013 does not reflect a developmental brain expression program 1014 and its emergence has been associated with both genetic and 1015 non-genetic factors, including interaction with immune cells 1016 and hypoxia [Hara et al., 2021]. 1017

The cell state heterogeneity in GBM has been a challenge 1018 for successful implementation of targeted therapies [Nicholson 1019 and Fine, 2021], so understanding the mechanisms and dynamics 1020 of cell state plasticity could provide insight into more effective 1021 treatment regimens. To examine the potential heritability or 1022 plasticity of these cell states, we re-analyzed MGH115, a human 1023 patient-derived GBM sample with annotated phylogeny with (i) 1024 continuous gene transcriptional module scores (generated from 1025 module-specific gene expression using matched scRNAseq) and 1026 (ii) categorical cellular state annotation based on the per cell 1027 maximum transcriptional module score (**Fig. 6C**). The stem-like 1028 (NPC-/OPC-like) and MES-like transcriptional modules displayed 1029 high phylogenetic auto-correlations, suggesting that in this 1030 specimen, the expression of these genes is in part heritable. 1031 The AC-like module, however, was not significantly phylogenetically 1032 auto-correlated, suggesting that the transcriptional 1033 1034

state was more plastic in this patient sample (**Fig. 6D**).

As the MES-like state does not recapitulate any neurodevelopmental expression program and has been reported to be influenced by non-genetic factors [Hara et al., 2021, Neftel et al., 2019], it is distinct from the other GBM cell states. Interestingly, recent work has demonstrated that the MES-like state is driven by interactions between the tumor cells and immune cells, and has suggested that the targeted induction of the MES-like cell state together with immunotherapy may represent a novel opportunity for therapeutic intervention [Hara et al., 2021]. The neurodevelopmental-like transcriptional modules (NPC-/OPC-/AC-like) were more phylogenetically correlated with each other than any individual module was with the MES-like module (**Fig. 6E**). However, among the neurodevelopmental transcriptional modules, the AC-like module was the most phylogenetically correlated with the MES-like module, suggesting that transit between neurodevelopmental-like (NPC-/OPC-/AC-like) and MES-like states is driven by the AC-like state. To explore these relationships between GBM states further, we next used the phylogenetic correlations of GBM cell states, as determined by the per cell maximum transcriptional module scores, to infer cell state transition probabilities. This analysis revealed that stem-like cells primarily differentiated into AC-like cells, which could either dedifferentiate back into a stem-like state [Chaligne et al., 2021] or progress to the MES-like state (**Fig. 6F**). Notably, this inference suggests that, in this patient, the MES-like state derives from transitioning AC-like cells. This observation is consistent with recent findings that show that many MES-like cells have AC-like properties [Chanoch-Myers et al., 2022] and that the receptors (*e.g.*, OSMR, EGFR, PDGFRB, and AXL) for ligands that drive transition into the MES-like state are expressed in AC-like cells but not stem-like cells [Hara et al., 2021]. PATH transition inferences from another human patient-derived GBM sample MGH122, from Chaligne et al. [2021], agreed with inferences from MGH115, revealing that of the neurodevelopmental-like cell states, AC-like cells appear to transition to the MES-like state at the highest rate (**Fig. S6A**).

To experimentally corroborate these cell state transition inferences obtained from primary human samples, we leveraged the artificial Molecular recorder approach [Chan et al., 2019] to trace gliomasphere phylogenies, using MGG23 [Wakimoto et al., 2011], a human patient-derived gliomasphere model (**Methods: Gliomasphere phylogenies, Fig. 6G**). Gliomaspheres are spheroid GBM cultures capable of recapitulating parental tumor cellular diversity [Laks et al., 2016], and thus represent an appropriate setting to measure cell state heritability versus plasticity. Two gliomasphere MGG23 replicates were grown *in vitro* for 4 weeks, at which point phylogenies were reconstructed using recovered barcodes, and cells were annotated according to their scRNAseq profiles. Consistent with the human patient data

(**Fig. 6E**), PATH measurements in the gliomasphere model also showed higher phylogenetic correlations between the neurodevelopmental-like modules, than between any of the neurodevelopmental-like and MES-like modules (**Fig. 6G**). Furthermore, among the neurodevelopmental-like modules, the AC-like module was, as in patient sample MGH115, the most correlated with the MES-like module. Thus, using both native and artificial approaches for phylogenetic tracing in primary human samples and an *in vitro* model, respectively, we observed a strong phylogenetic relationship between the AC- and MES-like transcriptional programs; consistent with a model in which the MES-like cell state primarily derives from the AC-like state.

Finally, after analyzing the heritability of predefined glioblastoma gene transcriptional modules, using gene set enrichment analysis (GSEA) [Subramanian et al., 2005] we next profiled the heritability of the 3,000 most variably expressed genes in MGH115 (**Table S4**), ranked by their autocorrelation z scores, to discover heritable modules in an unbiased fashion. Consistent with **Fig. 6D**, this revealed an overrepresentation of five (NPC1/OPC/AC/MES1/MES2) GBM gene modules. This analysis further revealed that targets of the Polycomb repressive complex 2 (PRC2) constituents (*i.e.*, targets of EED, SUZ12, EZH2), as well as sets of genes with promoters characterized by high CpG density and the repressive histone mark H3K27me3, in multiple stem cell contexts, were also enriched among heritably expressed genes in glioblastoma (**Fig. 6H, Table S5**). Similarly, brain tissue genes with bivalent promoters that are dually marked by both H3K27me3 and the activating mark H3K4me3, were also enriched among heritably expressed genes (**Fig. 6H**). This promoter methylation pattern represents a poised functional state that generally resolves to repressed (H3K27me3-only) or active (H3K4me3-only) states as cells differentiate. Promoter H3K27me3 levels are maintained primarily by targeting of the chromatin modifying PRC2, preventing differentiation by repressing lineage-specific gene expression [Boyer et al., 2006]. Notably, activity at PRC2-targeted sites is a key switch in the differentiation and maintenance of glioma stem cells [Natsume et al., 2013, Suvà et al., 2009].

To understand the relationships between these highly heritable gene modules, we next analyzed the enrichment of gene sets within distinct heritable gene modules defined by cross-correlations, with Over-Representation Analysis (ORA) [Korotkevich et al., 2021]. Hierarchical clustering of the phylogenetic correlations between the top 100 most auto-correlated genes revealed two heritable gene modules in MGH115 (**Fig. S6B, Table S6**). The first heritable module was enriched for gene sets associated with the neurodevelopmental-like glioma cell states (NPC1/OPC/AC), EED (a PRC2 subunit) target genes, and genes with high CpG density promoters with H3K27me3. This result is consistent with our previous

1143 observation that PRC2-target genes are preferentially hy- 1155
 1144 pomethylated, accessible and activated in the stem-like cell 1156
 1145 states [Chaligne et al., 2021]. The second heritable module 1157
 1146 was enriched for genes associated with the MES-like state 1158
 1147 and gene signatures associated with hypoxia. These results 1159
 1148 suggest that in patient MGH115, glioblastoma cells could 1160
 1149 occupy one of two heritable transcriptional states, either 1161
 1150 neurodevelopmental-like or mesenchymal-like. Cells could 1162
 1151 transit between these two states, primarily when occupying 1163
 1152 the more astrocyte-like end of the neurodevelopmental-like 1164
 1153 spectrum. Further, the neurodevelopmental-like module, in 1165
 1154 particular the stem-cell like states, is likely heritably main-

tained by PRC2 activity. These findings further highlight
 PATH's ability to extract epigenetically grounded and bi-
 ologically relevant expression profiles from single cell tran-
 scriptional and phylogenetic data in an unbiased manner.

In summary, the application of PATH to primary human
 glioblastoma samples identified the expected phylogenetic
 similarity by spatial location, nominated AC-like cells as the
 candidate precursor for MES-like cells, and highlighted the
 role of PRC2 in stable propagation of stem-like cell states.
 Thus, PATH can provide critical insight as to the biology
 underlying transcriptional cell state diversity in cancer.

Figure 6

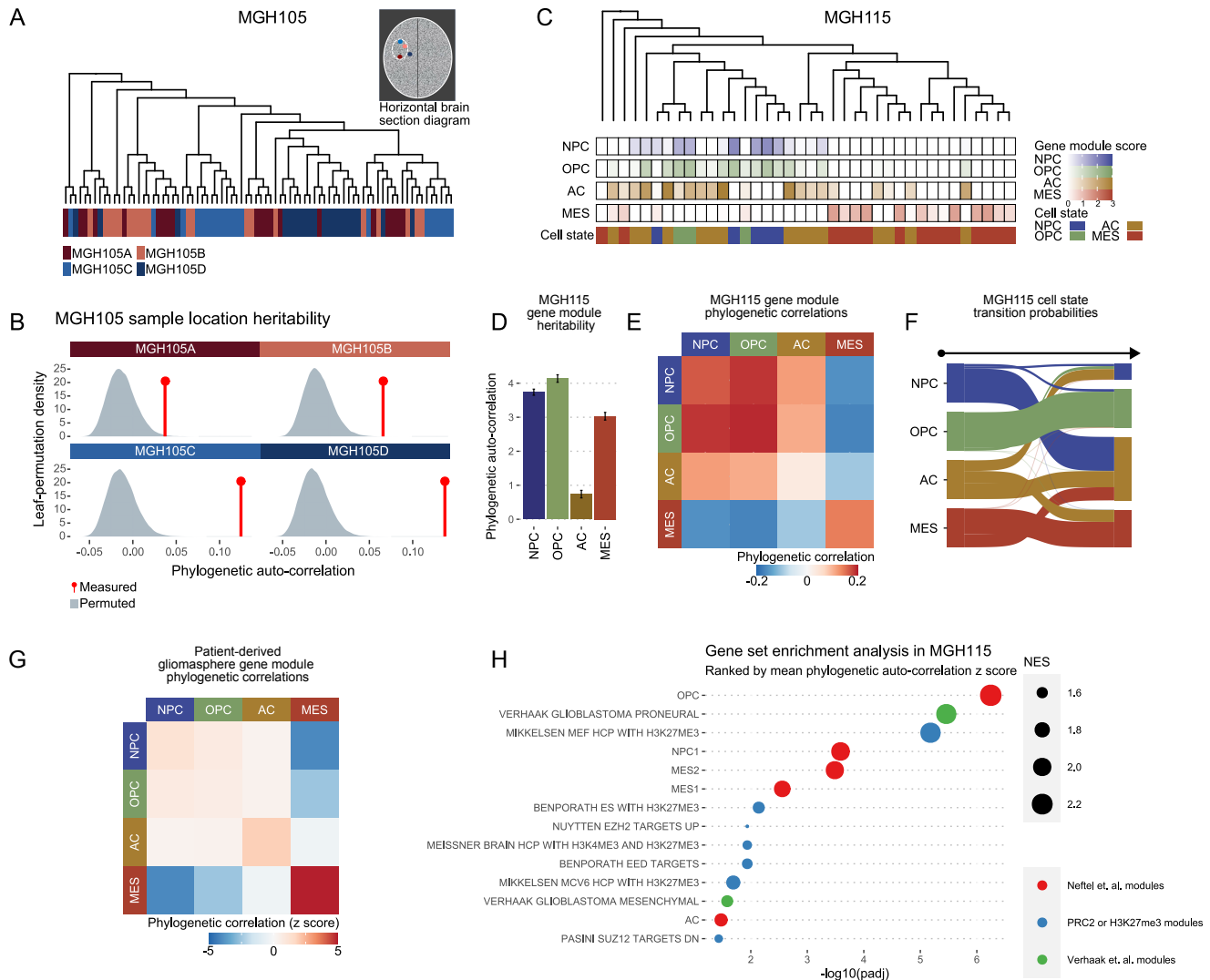


Figure 6: Heritable transcriptional modules and cell state transition dynamics in human glioblastoma

A) Human GBM sample (MGH105) single-cell consensus phylogeny containing 80 cells (20 from each tumor location) with tumor sample location projected onto leaves. Inset is a schematic of the four MGH105 patient tumor sample locations.

- B)** Leaf-permutation test (10^6 permutations) of tumor sample location phylogenetic auto-correlation. Density plot depicts leaf-permutation auto-correlations and red lines show measured (non-permuted) phylogenetic auto-correlations. 1171
1172
- C)** Human GBM patient sample (MGH115) single-cell phylogeny (replicate 6) containing 38 cells with GBM gene module scores and categorical cell states projected onto leaves. 1173
1174
- D)** Replicate mean (across 9 MGH115 phylogeny replicates) phylogenetic auto-correlation z scores for GBM gene module scores for patient sample MGH115. 1175
1176
- E)** Replicate mean phylogenetic correlation heat map for patient sample MGH115 GBM gene modules. 1177
- F)** Sankey plot of replicate mean Markov transition probabilities inferred from categorical state phylogenetic correlations in patient sample MGH115 phylogeny replicates. Probabilities shown are shown for $\hat{P}(\tau)$ (**Methods: Inferring cell state transitions from phylogenetic correlations**). 1178
1179
1180
- G)** Replicate mean phylogenetic correlation z score heat map for gliomasphere GBM gene modules, using one-node weighting. 1181
- H)** Dot plot of enriched pathways from GSEA of chemical and gene perturbation curated gene sets (C2:CGP) and six GBM gene modules (NPC1-/NPC2-/OPC-/AC-/MES1-/MES2-like) [Nefitel et al., 2019] for patient sample MGH115, with genes ranked by their phylogeny-replicate mean phylogenetic auto-correlation z scores (**Methods: Phylogenetic correlations, Methods: Human patient glioblastoma**). Only select gene sets are depicted; other significantly enriched gene sets can be found in **Table S5**. Dot sizes are proportional to GSEA normalized enrichment scores (NES). 1182
1183
1184
1185
1186
- GBM gene modules (NPC-/OPC-/AC-/MES-like) were shortened to (NPC/OPC/AC/MES). 1187

Quantifying cell state heterogeneity in B-cell acute lymphocytic leukemia (B-ALL) using single-cell whole genome sequencing

An exciting next frontier in the analysis of somatic evolution in humans is using somatic mutations as native lineage barcodes for lineage tree reconstruction from single-cell whole genome sequencing (scWGS). Current approaches often rely on costly and low-throughput single-cell cloning followed by WGS [Lee-Six et al., 2018], as somatic mutation rates are low and many scWGS methods suffer from high error and dropout rates, impacting the ability to call somatic variants with high confidence from single cells. To circumvent these challenges, and to explore PATH application to newly generated single-cell phylogenies constructed from the whole genome sequencing of single cells, we harnessed primary template-directed amplification [Gonzalez-Pena et al., 2021], a scWGS method based on a quasi-linear amplification that allows for high reproducibility and low allelic dropout. We aimed to construct a high-resolution lineage tree from scWGS of a B-ALL patient sample (**Fig. 7A**) with accompanying flow cytometry data for cell surface markers, and then apply PATH to determine the heritability versus plasticity of therapeutically relevant phenotypes in tumor cells.

To leverage somatic mutations as native lineage barcodes, we generated whole genome sequences for 86 cells (~8x coverage) sampled from a patient with B-ALL (**Methods: B-ALL analysis**) and quantified levels of cell surface markers that represent both more immature B cell states (CD34, CD10 and CD38) and more mature B cell states (CD19,

CD20 and CD45) [Welner et al., 2008]. We used 55,251 single nucleotide variants (SNVs) to construct a high-resolution phylogeny (**Methods: B-ALL analysis**), annotated with genetic (copy number deletion, exonic SNVs excluded from tree reconstruction) and phenotypic (cell surface marker expression) information, with sorting time as a control for a random, non-heritable trait (**Fig. 7A, Table S7**). To determine the heritability of each trait, we applied PATH to these data to calculate phylogenetic correlations. As expected, genetic variation was highly heritable and sorting time, a random control, was not heritable (**Fig. 7B**). However, the phenotypic information was more variable; the majority of markers had intermediate phylogenetic scores that were between those of the genetic and random traits, with CD34 and CD20 displaying the highest heritability (**Fig. 7B**). These results showed that PATH can be used to analyze single-cell phylogenies generated from scWGS data and to measure the heritability of cell-surface protein expression markers in tumor cells.

To more deeply explore the biology of these tumor phenotypic traits, we next calculated the phylogenetic cross-correlation between the significantly heritable cell surface markers (**Fig. 7C**). PATH showed that a marker associated with more immature B cells (CD34) negatively cross-correlated with markers associated with more mature B cells (CD19, CD20 and CD45), which in turn were strongly cross-correlated with one another. These results indicated that this B-ALL sample comprised tumor cells with heritable earlier and later B cell states, suggesting that some structure of the normal B cell differentiation trajectory is retained in this sample.

Figure 7

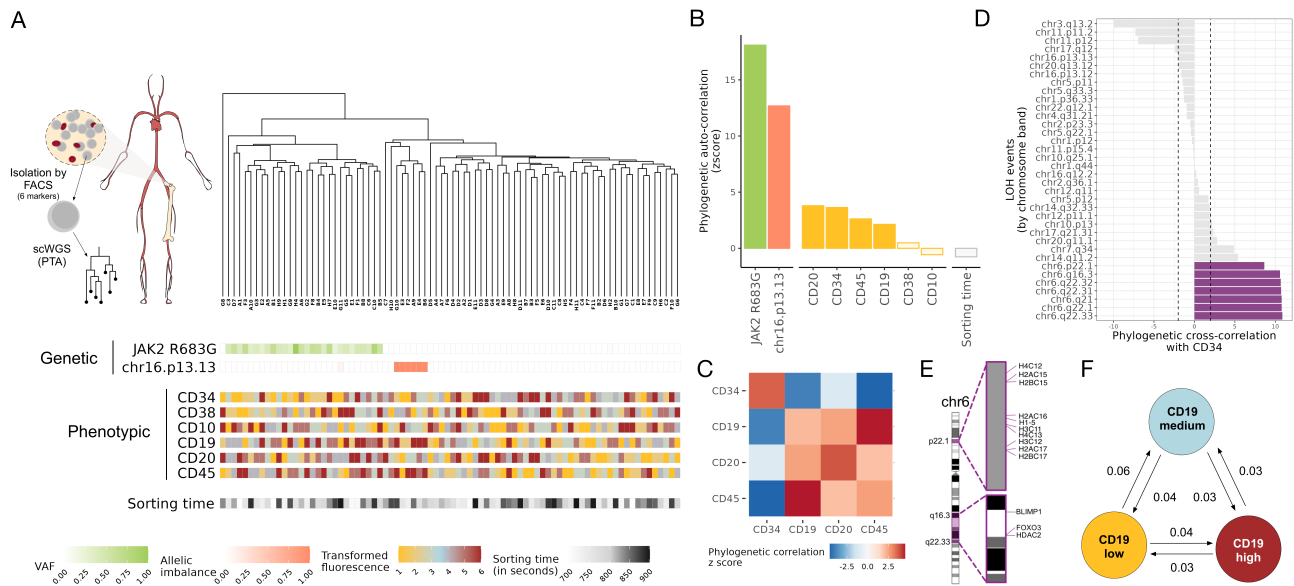


Figure 7: Quantifying cell state heterogeneity in B-ALL using single-cell whole genome sequencing

A) Top left- schematic of single-cell whole genome sequencing (scWGS) by primary template-directed amplification (PTA) of bone marrow isolated B cells sorted using six cell surface markers from a B-cell acute lymphocytic leukemia patient. Single-cell whole genome sequences were used to construct a single-cell phylogeny.

Top right- Lineage tree constructed from single-cell whole genome sequences from a B-ALL patient sample (N=82 cells; ~8x coverage).

Bottom- Genetic [allelic imbalance of germline heterozygous SNPs indicating a copy-number deletion at chr16; variant allele frequency (VAF) of single-nucleotide variant (SNV) of *JAK2*], Phenotypic (fluorescence of cell surface markers) and Random (sorting time) traits projected onto leaves. Cell surface markers used for cell sorting: CD34, CD10 and CD38 represent more immature B cell states, CD19, CD20 and CD45 represent more mature B-cell states.

B) Phylogenetic auto-correlation z scores for genetic (copy-number deletion and SNV as in **B**), phenotypic (cell surface protein markers) and random (sorting time) factors.

C) Phylogenetic correlation z score heat map for heritable cell surface protein markers.

D) Phylogenetic cross-correlation z scores for CD34 and copy number deletions. Phylogeny annotated with genome-wide copy number deletion map can be found in **Fig. S7**.

E) Chromosomal regions of deletions in clones with high CD34 expression.

F) PATH inferred transition probabilities between states (CD19 low, medium and high) using all cells. Values rounded to the nearest hundredth.

1270 Taking advantage of the multimodality of the single-cell lineage data, we next sought to identify genetic features that
1271 correlated with CD34 expression, a marker that displayed
1272 high heritability and that reflects a more immature B cell
1273 state. To associate genetic and phenotypic features, we
1274 calculated phylogenetic correlations between copy number
1275 deletions and CD34 expression. PATH identified high phy-
1276 logenetic correlations between CD34 expression and chro-
1277 mosome 6p22.1 and 6q16-q22 region deletions (**Fig. 7D,**
1278

Fig. S7), indicating that tumor clones that harbored these
1279 specific deletions also had higher CD34 expression. To iden-
1280 tify potential genetic contributors that are associated with
1281 CD34 expression in these tumor clones, we more closely ana-
1282 lyzed the deleted chromosomal regions and their impacted
1283 genes. Interestingly, these regions harbor genes that encode
1284 important B cell differentiation factors including PRDM1,
1285 FOXO3 and HDAC2 on 6q, as well as a histone gene clus-
1286 ter on 6p (**Fig. 7E**). Notably, it has been shown in B-cell
1287

lymphoma that deleterious mutations in histones H1B/H1-5 can cause remodeling of the chromatin state [Yusufova et al., 2021], leading to expression of stem cell genes, which is consistent with the earlier B cell state phenotype that we observed in cells harboring these deleted regions in this B-ALL sample. Therefore, it is possible that copy number loss of these regions and deletion of these genes could potentially contribute to the emergence of an earlier, more stem cell-like state (CD34 high). Indeed, 6p22.1 is known to be relatively frequently deleted in B-ALL and 6q16-q22 in DLBCL [Brady et al., 2022, Chapuy et al., 2018], further supporting the link between these deletions and a more stem-like state in this sample. Thus, PATH showed that quantifying the heritability of phenotypes and analyzing cross-correlation with genotypic features nominates candidate genotype-to-phenotype associations.

Finally, we sought to harness the ability of PATH to quantify transition dynamics between cell states to interrogate the plasticity of B-ALL targets of immunotherapy. In contrast to acute myeloid leukemia, where tumor cells develop from a more restricted window of cells from across the hematopoietic developmental trajectory [Miles et al., 2020, Zeng et al., 2022], B-ALL is considered more functionally plastic based on transplantation assays [Rehe et al., 2013] and cell-of-origin studies [Johnsen et al., 2014]. However, there is limited direct evidence of lineage-informed cell state plasticity and transitions directly in human samples at the single-cell level. Importantly, B-cell markers including CD19 have been used as targets for chimeric antigen receptor T (CAR-T) cell therapy [Davila and Brentjens, 2016, Maude et al., 2014], and while this approach has had success, there remain limitations in efficacy and sustained response [Schroeder et al., 2022]. B-ALL relapse after treatment with CD19-targeted CAR-T cells can be driven by genetic loss of CD19 [Xu et al., 2019], but other mechanisms, including the intrinsic plasticity of cell states associated with CAR-T target expression, could affect treatment implementation and success. We note that while PATH showed that CD19 expression had positive phylogenetic auto-correlation, (Fig. 7B), this marker had lower heritability compared to other analyzed markers and was substantially lower than the heritability of genetic traits, suggesting that CD19 expression was at least partially plastic. Indeed, PATH quantification of the transitions between high, medium and low CD19 expression states (Methods: B-ALL analysis) showed that while CD19 expression states were largely stable, we detected transitions between all three states. In particular, the low CD19 expression state was more likely to transition to the medium state, while the high CD19 expression state was about equally likely to transition to medium or low states (Fig. 7F). Thus, these results showed that there is a low level of fluid transitions between high, medium and low CD19 states, suggesting that in this B-ALL sample, while CD19 expression was a heritable trait with a positive phylogenetic correlation, it also exhibited a degree of plasticity

between these expression level states. Altogether, these results and analyses highlighted the power of single-cell whole genome sequencing for phylogenetic analysis of human tumor cells, as well as the ability of PATH to quantify the heritability of therapy-relevant traits in a lineage-informed manner in order to gain insights into the plasticity of tumor cell states across subclones of a phylogeny.

Discussion

The cells that comprise a multicellular organism derive from a single ancestral cell, thus remaining nearly genetically identical. Despite this genetic similarity, somatic cells within a multicellular organism encompass vast functional and phenotypic diversity. This phenotypic diversity can be maintained across mitotic divisions through the heritable transmission of both cell-intrinsic factors, such as epigenetic marks [Bintu et al., 2016, Halley-Stott and Gurdon, 2013] (*e.g.*, DNA methylation and histone modifications) and cell-extrinsic factors (*e.g.*, microenvironment). Each somatic cellular division, however, presents an opportunity to introduce changes to these heritable factors, for example in the form of heritable genetic or epigenetic changes. The phenotypic effect of these changes, however, is highly context dependent. In the case of cancer, mutations in putative cancer driver genes do not always lead to tumorigenesis and depend on cellular identity. For example, the malignant competence of BRAF mutations is dependent on the transcriptional background [Baggiolini et al., 2021], and some somatic mutations that confer a proliferative advantage are masked when found in progenitor cells [Nam et al., 2019]. As the presence of phenotypic variation provides a substrate for natural selection, an understanding of how these phenotypes are differentially encoded and inherited will help us dissect how cells in the soma evolve throughout the lifespan. To achieve this, however, we need an integrative model of somatic evolution informed by phenotypically annotated phylogenies. As such, scRNAseq is not sufficient and must be coupled with technologies that can also deliver information on cell ancestry.

To address this gap, PATH delivers an analytic framework needed for analyzing novel multi-omic lineage tracing single-cell datasets. PATH achieves this by building upon approaches from quantitative genetics and evolutionary biology used to measure heritability and phylogenetic signal [Blomberg and Garland, 2002] and adapts these to a somatic context. Specifically, PATH offers a bivariate generalization of phylogenetic signal in the form of phylogenetic correlation. Using phylogenetic correlations, PATH measures the ancestral dependency of single-cell phenotypes to infer their heritability versus plasticity. Additionally, for categorical phenotypes, such as a cell state or identity, PATH can transform phylogenetic correlations into state transition probabilities and thus allows for the inference of unobserved cellular dynamics. Importantly, this transformation also makes the

1395 classic interpretation of phylogenetic signal more concrete,
1396 as phenotypic transition dynamics are directly linked with
1397 the measurement of phylogenetic signal.

1398 In step with the rapid advancement of lineage tracing tech-
1399 nologies, other frameworks, such as *Hotspot* [Detomaso and
1400 Yosef, 2021] and *The Lorax* [Minkina et al., 2022], have been
1401 developed to study the lineage dependency of phenotypes
1402 in the single-cell context. Unlike other approaches, how-
1403 ever, PATH can connect such measurements with a model
1404 of evolutionary dynamics and infer (categorical) phenotypic
1405 transition probabilities. Leveraging this connection, PATH
1406 allowed us to study how technical (*e.g.*, sampling and recon-
1407 struction fidelity) and biological variables affect heritability
1408 measurements. This can inform our interpretations, for ex-
1409 ample, as PATH makes it clear that when sampling is suffi-
1410 ciently sparse, heritable phenotypes will likely appear plas-
1411 tic.

1412 Other methods have also been advanced to estimate state
1413 transitions from phylogenies. For instance, if representing
1414 phenotypic (*e.g.*, cell type) transitions as a Markov model,
1415 transition probabilities can be fit using Maximum Likelihood
1416 Estimation (MLE) [Louca and Pennell, 2019] or inferred
1417 with kin correlation analysis (KCA) [Hormoz et al., 2015,
1418 2016]. PATH’s inference approach is more akin to KCA, as
1419 it transforms correlations into transitions; however, PATH
1420 can additionally be applied to subsampled phylogenies and
1421 when branch length measurements are absent. MLE, on the
1422 other hand, is commonly used in evolutionary biology to
1423 infer phenotypic transitions from species phylogenies. This
1424 approach takes the structure of the entire phylogeny into
1425 account (as opposed to just phylogenetic correlations) and
1426 searches for optimal transition rates. PATH’s accuracy is
1427 comparable to MLE, but computationally faster, particu-
1428 larly for larger trees with many phenotypes. This ability
1429 to accurately handle large trees with speed renders PATH
1430 suitable for analyzing single-cell phylogenies, which often
1431 contain many states, and an ever growing number of cells.

1432 Using PATH, we studied previously published developmen-
1433 tal lineage tracing datasets in early stages of embryologi-
1434 cal development [Chan et al., 2019] and brain organogen-
1435 esis [Raj et al., 2018]. In murine development, we were
1436 able to analyze phylogenetic correlations between the blasto-
1437 cyst, the germ layers and specialized tissues, reconstructing
1438 known developmental trajectories and importantly, captur-
1439 ing the dual origin of the gut endoderm from both the epi-
1440 blast and primitive endoderm [Kwon et al., 2008, Rothová
1441 et al., 2022, Saykali et al., 2019], which would not be achiev-
1442 able with scRNAseq alone. This highlights the ability of
1443 PATH to distinguish between phenotypic and ancestral simi-
1444 larity. We further showed that, consistent with a model of
1445 epigenetic inheritance and our understanding of imprinting
1446 throughout development [Loda et al., 2022], a unique X-
1447 chromosome expression profile is inherited by gut cells with
1448 extraembryonic origins. In zebrafish brain development, we

used PATH to show how anatomic proximity influences re-
latedness of neurons in the developing brain and further
highlighted PATH’s ability to coordinate transcriptional and
anatomic data to show a shared lineage between substruc-
tures in the fore, mid and hind brain. As multi-modal single-
cell technologies improve, PATH could be applied to coordi-
nate transcriptional data with other modalities, beyond
anatomic location, to interrogate fundamental questions in
development. We also observed a striking pattern of stable
lineage commitment for both excitatory (glutamatergic) and
inhibitory (GABAergic) neurons in the forebrain. As lineage
tracing techniques improve, using PATH we may eventually
be able to more finely map the transitions undergirding cell
state differentiation hierarchies in these functionally com-
plex organs and reveal the factors responsible for maintain-
ing and modifying lineage commitments.

Many scRNAseq analyses have revealed cell state diversity in
cancer, but representing only a snapshot, have been unable
to determine how temporally stable or transient such cell
states are. Using PATH on lineage traced scRNAseq data,
we can bypass this constraint, to quantify cell state tem-
poral dynamics. To demonstrate this potential, we applied
PATH to two previously published single-cell cancer datasets
[Chaligne et al., 2021, Simeonov et al., 2021]. First, we ob-
served that spatial location was highly stable: metastatic
tissue location in a mouse model of pancreatic cancer, and
tumor region in a human glioblastoma. Second, we used
PATH to study transcriptional stability. It is not yet clear
whether cancer cell state diversity predominantly reflects
transient transcriptional fluctuations akin to entering and
exiting the cell cycle, or more stable transcriptional changes
analogous to cell fate commitment in development. In both
cancer datasets, we observed the heritability of transcrip-
tionally defined cell states in two of the largest drivers of
cancer cell state diversity – position along the EMT con-
tinuum in pancreatic cancer, and in the stem cell hierarchy
in glioblastoma. Interestingly, in both of these cancers, cell
states were not uniformly plastic/heritable. Future appli-
cation of PATH to other cancers could guide future treat-
ments, such as the strategic targeting of specific transcrip-
tional states, or the therapeutic modulation of state transi-
tion rates, in order to drive tumors to extinction.

Underscoring this potential, our analysis of newly gener-
ated data from a B-ALL patient demonstrated that using
a powerful new single-cell whole genome sequencing ap-
proach (PTA) enabled construction of a high-resolution tu-
mor cell phylogeny, and that application of PATH to this an-
notated tree yielded a detailed cancer profile encompassing
genetic, phenotypic and ancestral dimensions. This PATH
profile provided quantitative measurements of the heritabil-
ity and plasticity of cell surface marker expression, reveal-
ing heritability of early vs. late B cell differentiation states,
and linking these state biases with potential underlying ge-
netic aberrations. Moreover, PATH analyses also quantified

1503 the plasticity of the therapeutically-relevant B-ALL marker
1504 CD19, which has been successfully used as a target of CAR-
1505 T immunotherapy [Schroeder et al., 2022]. As cell state
1506 plasticity in the expression level of a therapeutic target can
1507 serve as a potential evolutionary therapeutic escape mech-
1508 anism, we propose that such information could potentially
1509 serve to prioritize therapeutic targets for clinical develop-
1510 ment.

1511 We speculate that as sequencing costs continue to fall, clinical
1512 single-cell whole genome sequencing for phylogeny re-
1513 construction and analysis of tumor samples could become
1514 more accessible, rendering such approaches feasible.

1515 In conclusion, somatic evolution represents an exciting frontier
1516 in evolutionary biology, where asexually reproducing
1517 somatic cells evolve over the multicellular organism’s life
1518 span. Studying this frontier requires analytical advances in
1519 step with technological advances that provide multi-modal
1520 single-cell annotation with high resolution phylogenetic in-
1521 formation. We envision that PATH can thus help transform
1522 qualitative key concepts in multicellular somatic biology
1523 such as fate-commitment, heritability and plasticity
1524 into precise measurements, with broad impact on our under-
1525 standing of organismal health and disease. As future tech-
1526 nology evolves to capture phylogenetic information with epi-
1527 genetic and spatial information, we further envision that the
1528 adaptability of the PATH framework will enable the linkage
1529 of cell state heritability and the mode of inheritance propa-
1530 gation (*e.g.*, genetic, epigenetic, cell-extrinsic) to define the
1531 fundamental principles of somatic evolution.

1532 **Limitations** Mathematical models represent an idealized
1533 situation, and in practice, can be robust to small violations
1534 to their assumptions. As outlined in the results and methods
1535 sections, several assumptions are made in PATH’s cell state
1536 transition inference model (*e.g.*, transitions are Markovian,
1537 cell states are near their equilibrium proportions). These
1538 assumptions should be (nearly) met if transition rates only
1539 depend on a cell’s current and not prior states, and when
1540 sampling is not biased. Other assumptions, such that cell
1541 birth or death rates do not differ as a function of cell state,
1542 could be violated and impact inferences. Specifically, if some
1543 cell states have much higher proliferation rates than others,
1544 inferred transition rates could be biased. Such a scenario
1545 represents an opportunity for future model development.
1546 However, such a model would likely rely on accurate
1547 branch length measurements and higher resolution single-
1548 cell phylogenies than are typically available now. Transition
1549 inference accuracy is also most reliable when heritability is
1550 significantly detected, as demonstrated in **Fig. 2C,D**, and
1551 inferences from phylogenies with insignificant phylogenetic
1552 correlations should be interpreted cautiously.

1553 Additionally, the robustness of PATH measurements is dependent
1554 on the quality and resolution of the lineage data, and analysis
1555 of sparsely sampled trees can lead to underes-

1556 timation of heritability, as shown by our simulations. Re-
1557 latedly, PATH is subject to the standard problems affect-
1558 ing single-cell analyses, including data dropout, accuracy of
1559 cell state assignment algorithms, completeness of gene set
1560 modules and batch effects. These limitations may constrain
1561 the analysis of currently available datasets; however, we an-
1562 ticipate that with advances in lineage tracing and single-
1563 cell multiomics technologies, PATH’s utility will expand as
1564 single-cell lineage tree data continue to improve.

1565 Most single-cell phylogenies do not include branch length
1566 estimates, which can further confound inferences. PATH,
1567 however, was designed to accommodate some of these lim-
1568 itations, by imputing branch lengths, and by focusing on
1569 closer (one-node apart) phylogenetic relationships.

1570 As more multi-omic single-cell lineage tracing experiments
1571 are conducted, and lineage tracing and other technologies
1572 further mature, allowing for even higher resolutions of phy-
1573 logenetic relationships and phenotypic states, more subtle
1574 evolutionary dynamics could be teased apart with PATH.
1575 If multiple layers of information, in addition to transcrip-
1576 tional phenotype and ancestry, such as location or microen-
1577 vironment, are gathered for each cell, measured phyloge-
1578 netic correlations across these layers could help dissect the
1579 encoding of heritable phenotypes. That is, phylogenetic cor-
1580 relations between phenotypes and microenvironments could
1581 help determine whether a heritable phenotype is encoded
1582 intrinsically (*e.g.*, via genetic or epigenetic mechanisms) or
1583 extrinsically (*e.g.*, via shared microenvironment stimuli).

1584 **Conclusion** In summary, throughout a multicellular or-
1585 ganism’s lifetime, its constituent somatic cells continuously
1586 evolve, accumulating heritable phenotypic variation. When
1587 positively selected, heritable phenotypic variation deleteri-
1588 ous to the organism as a whole may also lead to disease
1589 states or malignancy, which itself represents a “runaway”
1590 evolutionary process. PATH formally connects the analysis
1591 of cell state diversity and somatic evolution, and quantifies
1592 critical aspects, replacing *qualitative* conceptions of “plas-
1593 ticity” with *quantitative* measures of cell state transition
1594 and heritability. The application of PATH thus powerfully
1595 brings together approaches from evolutionary biology and
1596 single-cell technology, to study complex dynamics governing
1597 somatic evolution – an exciting novel frontier in multicellu-
1598 lar biology.

1599 Acknowledgments

1600 We thank members of the Landau laboratory and Norbert
1601 Fehér for thoughtful discussions throughout the develop-
1602 ment of this work. We thank Nir Yosef for critical comments
1603 on the manuscript. We thank Aaron McKenna and Bushra
1604 Raj for sharing data and code related to the scGESTALT
1605 phylogenies. We thank Alexander Meissner’s group and the
1606 authors of Chan *et al.* 2019 for sharing their cell type assign-

ment data and code. CG is supported by a Burroughs Wellcome Fund Career Award for Medical Scientists, National Institutes of Health Director’s New Innovator Award (DP2-CA239145), and Chan Zuckerberg Investigator Award. DAL is supported by the Burroughs Wellcome Fund Career Award for Medical Scientists, the Valle Scholar Award, the William Rhodes and Louise Tilzer-Rhodes Center for Glioblastoma at NewYork-Presbyterian Hospital (NYPH 203205-01), the Sontag Foundation (Distinguished Scientist Award, SFI 203261-01), the National Institutes of Health Director’s New Innovator Award (DP2-CA239065), Leukemia Lymphoma Scholar Award and the Mark Foundation Emerging Leader Award. This work was supported by the National Heart Lung and Blood Institute (R01HL157387-01A1), National Cancer Institute (R01 CA242020, R01 CA251138, and P50 254838), a Tri-Institutional Stem Cell Initiative award and the National Human Genome Research Institute, Center of Excellence in Genomic Science (RM1HG011014). DAL and MLS are jointly supported by NCI R01CA258763 and a grant from the STARR Cancer Consortium.

Competing interests

MLS is equity holder, scientific co-founder, and advisory board member of Immunitas Therapeutics. CG is a co-founder, equity holder, and board member of BioSkrby Genomics. DAL has served as a consultant for Abbvie, AstraZeneca and Illumina, and is on the Scientific Advisory Board of Mission Bio, Pangea, Alethiomics, and C2i Genomics; DAL has received prior research funding from BMS, 10x Genomics, Ultima Genomics, and Illumina unrelated to the current manuscript.

Author contributions

JSS, ARD, TP and DAL conceived the project and designed the study. JSS developed PATH and performed simulations. ARD, JSS, TP and SR performed analyses. YF and TH generated the gliomasphere data. YP and CG generated the single-cell PTA data. JSS, ARD, TP, MLS, CG and DAL helped interpret the results. MLS, YF, CG and TH provided critical comments on the manuscript. JSS, ARD, TP, CP and DAL wrote the manuscript. All authors reviewed and approved the manuscript.

Code availability

The code used to measure phylogenetic correlations and to infer cell state transitions is available as part of our PATH R software package at <https://github.com/landau-lab/PATH>. Code used for data processing and analysis will be made available upon publication.

Methods

Phylogenetic correlations

To quantify the distribution of a single-cell measurement, such as transcriptional state, across a phylogeny, we use Moran’s I [Moran, 1950], a classic measure of spatial auto-correlation. We also import its bivariate generalization, a measure of spatial cross-correlation [Chen, 2015, Wartenberg, 1985] to quantify pairwise phylogenetic cross-correlations [Chaligne et al., 2021]. For this study, we refer to both phylogenetic auto- and cross-correlations as *phylogenetic correlations*.

To compute the phylogenetic auto-correlation of a single variable (Moran’s I), we need a measurement of pairwise distances between cells, provided by the phylogeny, and a standardized observation per cell (with mean subtracted and normalized by population standard deviation).

For example, the expression of a particular gene in N cells could be represented by the N -dimensional vector x , where each element represents an expression score per cell. This vector is then standardized, producing the vector $z_x = (x - \mu_x) / \sigma_x$, where μ_x and σ_x are the mean and population standard deviation of x , respectively.

Pairwise phylogenetic distances (*e.g.*, node or branch length distances), represented by the elements of the square N -dimensional matrix L , are transformed into a phylogenetic weight matrix W , with a chosen weighting function f_w , such that $W = f_w(L)$. This function first weights each off-diagonal element of L , and then sets diagonal elements of L to 0. An example of a weighting function is the inverse of phylogenetic distance (*i.e.*, for $i \neq j$, $W_{ij} = 1/L_{ij}$, otherwise $W_{ij} = 0$). Another example of a weighting function that we use throughout this study is to select only a specific phylogenetic distance (*e.g.*, for $L_{ij} = d$ and $i \neq j$, $W_{ij} = L_{ij}$, otherwise $W_{ij} = 0$), where d is either a chosen branch or node distance. These weights are then normalized such that they sum to 1, resulting in a normalized weight matrix, \bar{W} . The phylogenetic auto-correlation of x is then defined as,

$$\phi_x = z_x^T \bar{W} z_x,$$

where superscript T signifies the matrix transpose.

The phylogenetic cross-correlation between two different single-cell measurements (bivariate Moran’s I), is calculated similarly, where both z_x and z_y are standardized single-cell measurements or observations corresponding to the vectors x and y ,

$$\phi_{yx} = z_x^T \bar{W} z_y.$$

All pairwise phylogenetic (auto- and cross-) correlations can be computed simultaneously if single-cell measurements are

in matrix form. Single-cell measurements are represented by the $N \times n$ dimensional matrix X , in which its N rows represent individual cells and its n columns represent distinct measurements (such as the expression of n distinct genes). When measuring phylogenetic correlations for a categorical states, in which a cell can occupy only one of a set of possible states at any given time (*e.g.*, cell type), each column of X denotes a distinct cell state, and the state of each cell is indicated by a 1 in the appropriate column, and 0s in the remaining columns. For example, if the i th cell is in the second of two possible cell states, then $X_{i,1} = 0$, and $X_{i,2} = 1$. For all measurement types, the columns of the single-cell measurement matrix X are standardized, as above, to produce the $N \times n$ dimensional matrix Z , which is then used to compute the square n -dimensional phylogenetic correlation matrix,

$$\Phi = Z^T \bar{W} Z.$$

Note that the diagonal elements of Φ correspond to phylogenetic auto-correlations. Furthermore, phylogenetic correlation z scores can be calculated by performing a leaf-permutation test or analytically with moments from Czaplewski and Reich [1993]. Phylogenetic correlations and analytical z scores can be computed with the function `xcor()` in our R software package. Additionally, normalized phylogenetic weight matrices can be computed using either `one_node.tree.dist()`, `inv.tree.dist()`, or `exp.tree.dist()` from our *PATH* R package.

Note that phylogenetic correlations depend on the structure of the matrix \bar{W} , thus weighting functions should be chosen carefully. For the purposes of this study, we predominantly use a weighting function that only includes cells that are each other's nearest phylogenetic neighbor, specifically cells that are separated by a node distance of one.

Simulating phylogenies

In this study we use two approaches to simulate single-cell phylogenies. We simulate *idealized phylogenies*, which are completely sampled, discrete-time, bifurcating, ultrametric, and balanced phylogenies that contain $N = 2^g$ cells, where g is the number of generations that have occurred since the root. Additionally, each branch length, which corresponds to one generation, has a length of one. To generate an idealized phylogeny we use the function `pmtree(b = 1, d = 0, n = N, type = "discrete")` from the R software package *phytools* [Revell, 2012].

We also simulate phylogenies using what we refer to as a *sampled somatic evolutionary process*, which is a sampled and continuous-time birth-death process, using the function `generate_tree_hbd_reverse()` from the R software package *castor* [Louca, 2020, Louca and Doebeli, 2018]. In contrast to idealized phylogenies, these phylogenies can be

imbalanced, and contain any number of cells that represent a fraction of the total somatic population. For these simulations, parameters for cell division (or birth), and cell death, the sampling rate, and the total number of sampled cells can be specified. Here, phylogenetic branch lengths correspond to time in continuous units, and not to generations, as in idealized phylogenies.

Cell state transition dynamics are represented as a discrete- or continuous-time Markov model (**Methods: Markov model of cell state transitions**) on idealized, and sampled somatic evolutionary phylogenies, respectively. Markov cell state transitions are simulated on both types of phylogenies using the *castor* function, `simulate_mk_model()`.

Markov model of cell state transitions

We model cell state transition dynamics as a Markov chain [Grimmett and Stirzaker, 2020], in both discrete- and continuous-time.

For a discrete-time Markov chain comprising n possible cell states, the transition probabilities (corresponding to one unit of time) are stored in a square n -dimensional *transition matrix*, P . Individual elements of the transition matrix are referred to by their subscript coordinates, such that P_{ij} refers to the transition probability located in row i and column j and represents the probability of switching from state i to state j . The probability that a cell in state i transitions to state j after t discrete time-steps is given by P_{ij}^t (note: superscript t reflects matrix, not element-wise, powers). As elements represent probabilities, each row of P must sum to 1.

Discrete-time chains might be more intuitive when recording times in non-overlapping generations, and continuous-time might be more appropriate when generation times vary and/or overlap. A continuous-time Markov chain has a *transition rate matrix*, Q . Each element, Q_{ij} records the infinitesimal transition rate between states indexed by their row and column. The transition probability matrix can be recovered by matrix exponentiating the rate matrix, that is $P = \exp(Q)$, and the transition probability of switching from state i to state j after a (continuous) t amount of time is given by $P(t) = \exp(Qt)$. Lastly, each row of Q must sum to 0.

The *stationary distribution* of a Markov chain, if also a *limiting distribution*, represents the expected frequencies of each cell state at equilibrium, and is represented by the n -dimensional vector π . For large t , the transition matrix P^t , if it has a limiting distribution, converge to the matrix Π , where each row of Π is equivalent to the vector π . This means that after a sufficiently long amount of time, the probability of transitioning from any state to state j is equal to state j 's equilibrium frequency, π_j . For chains with symmetric transitions, where transitions to and from a state

occur with equal probability (*i.e.*, $P_{ij} = P_{ji}$), the equilibrium frequency for each state is $1/n$, where, recall n is the number of possible cell states.

Finally, Markov chains are *reversible* if the products of the transition probabilities between two states and their stationary frequencies of origin are the same, *i.e.* $\pi_i P_{ij} = \pi_j P_{ji}$. Note that the reversibility of a Markov chain does not imply that transitions are symmetric, and that asymmetric Markov chains can also be reversible.

We connect Markov cell state transition dynamics with phylogenetic correlations in **Phylogenetic correlations and cell state transitions**, and use this connection to infer cell state transition dynamics from phylogenetic correlations in **Inferring cell state transitions from phylogenetic correlations**.

Phylogenetic correlations and cell state transitions

Phylogenetic auto-correlations measure the phenotypic similarity of closely versus randomly related cells (with respect to ancestry). More generally, the phylogenetic cross-correlation of two phenotypes, is a measure of the relationship between those phenotypes in closely related, as compared to, randomly chosen cells (**Methods: Phylogenetic correlations**). When measuring categorical states on phylogenies, if we use a phylogenetic weighting function that retains only specified phylogenetic distances and omits all others, phylogenetic correlations measure the difference between *state-pair* frequencies in closely (as specified by the retained distances) versus randomly related cell pairs. Here, *state-pair* refers to the states represented in a pair of chosen cells.

For example, on idealized phylogenies (**Methods: Simulating phylogenies**), if we apply a phylogenetic weighting function that preserves all branch lengths equal to two, and sets all other phylogenetic distances to zero, the phylogenetic correlation between two states will be a measure of the difference between the frequencies at which pairs of states are found within sisters versus random cell pairs. On idealized phylogenies, sister cells are separated by a branch length of two, because the branches that connect each of them to their parent, represent one generation, and thus have a branch length of one. Similarly, if a weighting function that retained only branch lengths equal to four is used, the resultant phylogenetic correlations, for an idealized phylogeny, would measure the difference between state-pair frequencies in first-cousins versus random cell pairs. In general, if we use a weighting function on an idealized phylogeny that only retains phylogenetic branch lengths equal to $2t$, phylogenetic correlations would measure the difference between the frequencies at which specific state-pairs are found within pairs of cells that share a most recent common ancestor (MRCA) t generations ago, versus randomly chosen cell pairs (with

replacement).

To illustrate, consider an idealized N -cell phylogeny and n possible cell states, in which the pairwise phylogenetic branch lengths between cells, represented by the square N -dimensional matrix L , and each cell's categorical state, recorded in the $N \times n$ dimensional matrix X (as in **Methods: Phylogenetic correlations**), are known. First, a weighting function that only retains phylogenetic branch lengths equal to $2t$ is applied, such that $W(t) = f_w(L, t)$, and the sum of the weights in $W(t)$ are normalized to equate to 1, resulting in the normalized phylogenetic weight matrix $\bar{W}(t)$. The frequency in which cells phylogenetically separated by a branch length distance of $2t$ are in states i and j is given by the ij th element of the square n -dimensional frequency matrix,

$$F(t) = X^T \bar{W}(t) X.$$

Note, that on a phylogeny, because the order of the cells within a pair is arbitrary, for $i \neq j$, the frequency of observing either the state-pair ij or state-pair ji , is given by the sum of the frequencies $F(t)_{ij} + F(t)_{ji}$. Additionally note that in the specific context of idealized phylogenies, state-pair frequencies as in $F(t)$ are equivalent to *kin correlations* [Hormoz et al., 2016].

These state-pair frequencies can be transformed into phylogenetic correlations, $\Phi(t)$, by first subtracting the random (with replacement) state-pair frequencies, and then normalizing by the cell state population covariances, where μ and σ are the respective n -dimensional state frequency and population standard deviation vectors (and division is element-wise),

$$\Phi(t) = (X^T \bar{W}(t) X - \mu\mu^T) / \sigma\sigma^T.$$

If cell state does not depend on ancestry, then we would not expect state-pair frequencies to substantially differ in closely and randomly related cells, resulting in low (near zero) phylogenetic correlations. However, if cell states can be inherited, but also sometimes stochastically transition, we would expect phylogenetic correlations to be generally non-zero. This is due to the fact that, if heritable, the states for cells that share a MRCA t generations ago will each depend on the state of the same ancestral cell. As such, state-pair frequencies and therefore phylogenetic correlations as measured above, will depend on how heritable each cell state is, and how often each state transition to another state occurs. In other words, the difference between state-pair frequencies in closely related versus random cells, might be attributable to underlying cell state transition and inheritance dynamics. To make this more concrete, below we link a Markov model of cell state transition dynamics with cell state phylogenetic correlations.

For cell state transition dynamics that can be represented as a Markov chain (**Methods: Markov model of cell state transitions**), we can predict state-pair frequencies for a given pairwise phylogenetic distance, from the transition probabilities P (a square n -dimensional matrix, where n is the number of cell states) and the limiting distribution π (an n -dimensional vector). For an intuitive example, consider the situation where a pair of sister cells (that share a parent) are in the same specific state. One way sister cells can end up in the same state is by both inheriting the same parental state, and subsequently not transitioning to another cell state. Alternatively, if the sister cells did not inherit their current state, they could have each independently transitioned from the parent's state to the same new state. The probability of observing sister cells in the same specific state is then determined by summing the probabilities for each different scenario that could lead to such an outcome. The probability of each scenario is computed by taking the probability that the unobserved ancestral cell (here the parent) was in a particular state, given by π , and multiplying by the relevant transition probabilities, provided by P . For the situation in which there are only two possible cell states, the probability of observing the state-pair ij (where one cell is in state i and its sister is in state j) is,

$$\pi_1 P_{1i} P_{1j} + \pi_2 P_{2i} P_{2j}.$$

More generally, for n possible cell states, the probability of observing each state-pair (where one cell is in state i and the other is in state j , and i and j can range from 1 through n), in two cells that share a MRCA t generations ago, where $D = \text{diag}(\pi)$ and superscript T is the matrix transpose, is

$$\left(P^{tT} D P^t \right)_{ij}.$$

If the cell state transitions are reversible, then $P^T D = (D P)^T = D P$, and the probability of observing each state-pair in cells separated by a phylogenetic distance of $2t$ can be simplified to be,

$$(D P^{2t})_{ij}.$$

These equations show that, for Markov transition dynamics at equilibrium, the probabilities of observing each possible state-pair are determined by the probability that the shared ancestor was in a particular state, multiplied by the probability that such a state transitioned to the two descendant cell states observed t generations later, and then summed for each possible ancestral state. For reversible chains, this is also equivalent to the probability of starting in one of the descendant states, followed by a transition to the other descendant state after the $2t$ time-steps that separates them.

Using these equations, we can compute expected phylogenetic correlations for cell state transitions. This is achieved by subtracting the probability of observing randomly chosen cells (with replacement) from the state-pair probabilities and normalizing by the cell state covariances,

$$\left(P^{tT} D P^t - D \Pi \right) / \Sigma.$$

For reversible transitions, this simplifies to,

$$D (P^{2t} - \Pi) / \Sigma.$$

An illustration for these calculations for two cell states is depicted in **Box S1**. Notice that as t increases, $P^{2t} \rightarrow \Pi$, and all phylogenetic correlations thus approach 0. This means that as cell pairs become more distantly related, their state-pair frequencies should approach those as if the two cells comprising the pair were drawn at random from the population. Also note that the closer transition probabilities are to cell state equilibrium frequencies, the less heritable cell states will appear. Furthermore, in this context, a high cell state phylogenetic auto-correlation would imply that the probability of transitioning to any other state is relatively low, and thus that the cell state is highly heritable.

In the context of species evolution, the auto-correlative method of measuring phylogenetic signal was not based on an evolutionary model, in contrast to signal metrics like Pagel's λ , and thus considered more difficult to interpret biologically [Münkemüller et al., 2012]. Here, not only do we define a bivariate measure phylogenetic signal using phylogenetic correlations, but we illuminate a connection between the measurement of phylogenetic auto- and cross-correlations with a model of evolutionary dynamics. This relationship with (categorical) phenotypic transitions thus clarifies the interpretation of what phylogenetic correlations measure. Finally, although we only make the connection explicit for categorical phenotypic states, phenotypic “covariance structures” (which will affect phylogenetic correlations) can be linked with a variety of evolutionary processes, including models for the evolution of continuous phenotypic states [Hansen and Martins, 1996].

The relationship between phylogenetic correlations and reversible cell state transition dynamics, can be used to infer unknown transition probabilities from phylogenetic correlations, as demonstrated in **Inferring cell state transitions from phylogenetic correlations**.

Inferring cell state transitions from phylogenetic correlations

Idealized phylogenies

For reversible Markov chains with a limiting distribution (**Methods: Markov model of cell state transitions**)

operating on idealized phylogenies (**Methods: Simulating phylogenies**, and **Phylogenetic correlations and cell state transitions**), transition probabilities can be inferred by converting phylogenetic correlations back into state-pair frequencies (not centered or normalized) and then dividing each row i by \hat{D}_{ii} , the corresponding cell state frequencies at a branch length distance of $2t$ (where \hat{D} is an estimate of D),

$$\hat{P}^{2t} = \hat{D}^{-1}F(t).$$

To arrive at the transition probabilities for a specific length of time, appropriate matrix powers or roots can be taken. For instance,

$$\hat{P} = \sqrt[2t]{\hat{D}^{-1}F(t)}.$$

In this setting, using idealized phylogenies, this formulation is equivalent to inferring transition probabilities using *kin correlation analysis* (KCA) [Hormoz et al., 2016], and conceptually similar to an approach for approximating nucleotide substitution rates [Yang and Kumar, 1996].

Finally, note that in this context, if the Markov chain does not have a limiting distribution, for instance, if it is periodic, we might not be able to infer the correct transition probabilities. For example, in the situation where there are two possible cell states, and the transition probabilities to and from each state are $P_{12} = P_{21} = 1$, and the self-transition probabilities are $P_{11} = P_{22} = 0$, then the states of every observed cell (in the terminal generation) will be the same, but different from the states in the cells from the previous generation. For this case, we would correctly infer that the self-transition probability of the state observed in the terminal generation is 1 after $2t$ time-steps, however, our estimates for an odd number of time-steps would be incorrect.

Phylogenies from a sampled somatic evolutionary process

Phylogenies resulting from a sampled somatic evolutionary process (**Methods: Simulating phylogenies**) contain only a sampling of the somatic population under study and continuous and non-uniform branch lengths. These factors must be taken into account in order to successfully infer transition probabilities. To accomplish this, we take the state-pair frequency matrix (used to compute phylogenetic correlations) at a *node-depth* of d , $F(d)$, by applying a weighting function that omits all phylogenetic distances that do not correspond to a node-depth equal to d , and the mean of the corresponding branch length distances τ . For each node-depth, we can approximate the transition matrix as,

$$\hat{P}(\tau) = \hat{D}^{-1}F(d).$$

This is an estimate of the transition probability matrix for a time proportional to the mean branch length distance between cells d nodes apart. For a completely sampled idealized phylogeny, $\tau = 2$.

More generally, we estimate $P(t)$ (for time t), to be

$$\hat{P}(t) = f_r(e^{\frac{\hat{Q}(\tau)}{\tau}t}),$$

where $\hat{Q}(\tau) = \log \hat{P}(\tau)$, and $f_r(\cdot)$ normalizes rows so that each sums to 1.

For circumstances in which branch lengths are unknown or inaccurate, for a node-depth of one, τ can be imputed if the cell sampling can be approximated and a model of somatic evolution is assumed. This can be accomplished by using branch lengths from simulated phylogenies from our somatic evolutionary process (**Methods: Simulating phylogenies**), or approximated analytically (**Methods: Imputing branch lengths**). Cell state transition dynamics can be inferred with the function `PATH.inference()` in our R software package.

All inferred transition rates for the analyzed datasets were determined in this manner, using either $\hat{P}(\tau)$ (as in **Figs. 6F, S6A**) or $\hat{P}(t = 1)$ (as in **Figs. 4D, 5G, 7F**).

Phylogenetic reconstruction

To simulate evolution, phylogenetic reconstruction, analysis and inference, we first simulate trees as a sampled somatic evolutionary process, a continuous birth-death process, (**Methods: Simulating phylogenies**) under various parameter schemes, in which the sampled tree size, and the birth, death, and sampling rates can vary. Once phylogenies are simulated, two distinct Markov processes are run: (1) a process simulating cell state transition dynamics, and (2) a process simulating the mutation/scarring of heritable cellular barcodes. The first Markov model is as described in the section **Markov model of cell state transitions**, and the second Markov model simulates barcode scarring and is a simple two-state, continuous-time, and symmetric model, with one rate parameter s , that runs independently for each mutable site contained within a cell's heritable barcode. The elements of the 2-dimensional square barcode scarring transition rate matrix are given by $Q_{11} = Q_{22} = -s$, and $Q_{12} = Q_{21} = s$.

Once both cell state transition dynamics and barcode mutations are simulated, a phylogeny is reconstructed – ignoring the true simulated phylogeny – with the unweighted pair group method with arithmetic mean (UPGMA) algorithm on pairwise-barcode Hamming distances. Branch lengths (evolutionary distances) are estimated from the number of barcode differences, using $-0.5 \log(1 - 2(h/l))/s$, where h is the Hamming distance, l is barcode length, and s is the barcode cut rate.

Reconstructed phylogeny error is scored by computing the normalized Robinson-Foulds distance [Robinson and Foulds, 1981] and Mean Path Length distances [Steel and Penny, 1993] between the reconstructed and true trees. Phylogenetic correlations (using a node-depth of one weighting function) computed for the true and reconstructed tree are also compared by taking their mean differences. Lastly, transition inference is performed using two approaches (**Methods: Inferring cell state transitions from phylogenetic correlations**), by either using measured (determined by the Hamming distances) or imputed (**Methods: Imputing branch lengths**; determined using estimated parameters of a sampled somatic evolutionary process) branch lengths to derive $\hat{P}(1)$ from $\hat{P}(\tau)$. Accuracy for both methods is assessed by measuring the Euclidean distances between the inferred and true/simulated transition probabilities.

Imputing branch lengths

For phylogenies in which branch lengths are unknown or potentially inaccurate, we can impute the phylogenetic branch lengths used to infer transition rates (**Methods: Inferring cell state transitions from phylogenetic correlations**) by using the sampled somatic evolutionary process model (**Methods: Simulating phylogenies**), using two approaches. In both cases, branch lengths are imputed by using either measurements or estimates to parameterize our sampled somatic evolution model. For the first, more exact, approach, we directly measure branch lengths that correspond to a node depth of one in simulations that use the estimated parameters. For the second, more approximate approach, we use an analytical expression, given a somatic evolutionary model parameterization, for computing the expected lengths of phylogenetic *pendant edges*, which are proportional to the branch length distances that separate cells phylogenetically one node apart. For a sampled somatic evolutionary process, pendant edge lengths are expected to be [Stadler and Steel, 2012],

$$\varepsilon = \frac{\gamma \log(\gamma/\xi) - \gamma + \xi}{(\gamma - \xi)^2},$$

where ξ is the product of the cell birth and sampling rates, and γ is the net growth rate, given by the cell birth minus cell death rates. Using this expression, we can impute the approximate branch length distance between cells separated by one node, to be 2ε . For $\gamma = 1$ (where ξ is equal to the sampling rate, $N_{\text{sample}}/N_{\text{population}}$), as sampling becomes sparse, $\varepsilon \approx \log(N_{\text{population}}/N_{\text{sample}}) - 1$, and branch length distances at a node-depth of 1 are expected to be proportional to the logarithm of the (inverse) sampling rate.

To test the robustness of our cell state transition inference approach when using imputed branch lengths, we input a sampling rate estimate by randomly selecting a rate within

one order of magnitude above or below the true simulated rate. That is, if the simulated sampling rate was 10^{-6} , we randomly select a sampling rate estimate between 10^{-5} and 10^{-7} , for imputing branch lengths when inferring transition rates using PATH.

Assessing cell state transition inference accuracy

To assess the accuracy of our inferences using PATH, we simulated phylogenies across a range of parameters, varying the cell sampling, birth and transition rates, as well as the number of cells and possible cell states. To generate a random n -dimensional transition rate matrix, for each cell state, $(n - 1)$ numbers are drawn from a uniform random distribution, ranging between 0 and 0.1, and sequentially assigned to each off-diagonal matrix element per row. As rows must sum to 0, the remaining (diagonal) element in each row is set to the negative sum of these randomly drawn values. After parameters are chosen and a transition rate matrix is randomly generated, phylogenies are simulated (**Methods: Simulating phylogenies**) and phylogenetic correlations (**Methods: Phylogenetic correlations**) and inferences (**Methods: Inferring cell state transitions from phylogenetic correlations**) are computed.

We also compared cell state transition rate inference accuracy with MLE. To do this, we used the function `fit_mk()` from the R *castor* package [Louca, 2020, Louca and Doebeli, 2018] to estimate the transition rate matrix \hat{Q} from a simulated phylogeny (**Methods: Simulating phylogenies**). To assess the accuracy of inferences using either PATH or MLE, we compute the Euclidean distance between the inferred transition probability matrix \hat{P} , for $t = 1$, and the true transition probability matrix P . Inferences using both PATH and MLE were performed on the same simulated phylogenies, and accuracies compared.

Mouse embryogenesis

Normalized RNA matrices and phylogenies were downloaded from Gene Expression Omnibus (GEO) series GSE117542 and imported into R (v. 4.1.3). Cell type annotations were provided upon request by the corresponding authors of the original publication [Chan et al., 2019]. Blastocyst layer annotations were inferred from germ layer membership. Phylogenies were extended by connecting node identifiers with single-cell barcodes using a dictionary provided in pickle files. We analyzed phylogenies for embryos 2 and 6 from [Chan et al., 2019]. Originally, these phylogenies contained one cell per subclone; however, we added the remaining cells to the phylogeny as leaves descending from the same node. Phylogenetic correlations (**Methods: Phylogenetic correlations**) were calculated using the one-node depth weighting function. For categorical states (*e.g.*, cell type) phylogenetic correlations, weight matrices were first

2177 row-normalized before sum normalizing.

2178 To calculate enrichment of heritable genes on each chromo- 2228
2179 some, the top 2,000 most variably expressed genes (calcu- 2229
2180 lated using *Seurat* [Hao et al., 2021]) were segregated by 2230
2181 chromosome. Each set of variable genes (on each chromo- 2231
2182 some) was further divided into genes that were “heritable” 2232
2183 (z score ≥ 3) or “non-heritable” (z score < 3). For each 2233
2184 chromosome, a Fisher’s Exact test comparing the number of 2234
2185 “heritable” and “non-heritable” genes on that chromosome 2235
2186 to those on all other chromosomes combined was performed. 2236

2187 Zebrafish brain development

2188 Normalized RNA matrices and cell annotation tables were 2237
2189 downloaded from GEO series GSE105010 and imported into 2238
2190 R (v. 4.1.3). Zebrafish [Raj et al., 2018] phylogenies were 2239
2191 obtained by parsing json files using code provided by the 2240
2192 authors. We used zebrafish 3 (“rep 1”) and 5 (“rep 2”) 2241
2193 phylogenies from [Raj et al., 2018]. Phylogenetic correla- 2242
2194 tions (**Methods: Phylogenetic correlations**) were calcu- 2243
2195 lated using one-node weighting function, and for categor- 2244
2196 ical states, weight matrices were row-normalized before sum 2245
2197 normalizing. 2246

2198 Minor changes were made to the cell annotation provided in 2247
2199 the original study. In **Fig 4A** and **Fig 4C**, neuronal cells 2248
2200 originally annotated as “S1/S2” (forebrain/midbrain) and 2249
2201 “Mix” were both considered as “Mix”. All cell types that 2250
2202 were not neurons or neuronal progenitors were considered 2251
2203 non-neural. 2252

2204 To impute phylogenetic branch lengths (**Methods: Im-** 2253
2205 **puting branch lengths**) for PATH transition inferences 2254
2206 (**Methods: Inferring cell state transitions from phy-** 2255
2207 **logenetic correlations**), we estimated a cell sampling rate 2256
2208 of 10^{-4} , which assumes that there were approximately 10^6 2257
2209 cells per brain [Marhounová et al., 2019]. 2258

2210 To classify forebrain neurons as either GABA+, Gluta- 2260
2211 matergic (Glut+), or “unassigned”, GABA and Glut marker 2261
2212 gene sets were scored across forebrain neuron cells in 2262
2213 the rep1 fish ($N = 270$) using the *Scanpy* [Wolf et al., 2263
2214 2018] `score_genes()` function. Cells with a positive score 2264
2215 (greater than 0) for either GABA or Glut marker gene set 2265
2216 were classified accordingly (no cells had a positive score for 2266
2217 both categories). Cells with scores of 0 in both gene sets 2267
2218 were considered “unassigned”. 2268

2219 Mouse model of pancreatic cancer

2220 Phylogenies, RNA count matrices and phenotype tables 2270
2221 were downloaded from GEO series GSE173958 and imported 2271
2222 into R (v. 4.1.3). As the available RNA matrices for the 2272
2223 murine pancreatic cancer model [Simeonov et al., 2021] were 2273
2224 counts, we normalized them using *Seurat* (v. 4.2.0) [Hao 2274
2225 et al., 2021]. Also, given that each mouse had been injected 2275
2226 with different parental clones whose relationships cannot be 2276

2227 established, we could only study the annotated lineages of 2228
2228 each clone independently. We analyzed the phylogeny from 2229
2229 “Mouse 1 Clone 1” from [Simeonov et al., 2021], which was 2230
2230 chosen because it contained the most cells of any clone an- 2231
2231 notated with an EMT score. All cell annotations were used 2232
2232 as published in the original paper. Phylogenetic correla- 2233
2233 tions (**Methods: Phylogenetic correlations**) were computed 2234
2234 with the one-node depth weighting function, and for cate- 2235
2235 gorical states, weight matrices were row-normalized prior to 2236
2236 sum normalizing. 2237

2237 EMT bins were created to discretize the EMT score across 2238
2238 the EMT continuum according to the following: cells were 2239
2239 partitioned along the continuum using units of 1 (bin #1 2240
2240 includes cells with EMT scores from 0 to 1, bin #2 includes 2241
2241 cells from 1-2, etc.), merging bins at the extremes (all cells 2242
2242 with a score of 7 or lower were assigned to a single bin, as 2243
2243 were cells that scored higher than 30) because these bins had 2244
2244 low cellular representation. To check for robustness, we re- 2245
2245 peated the binning procedure using other intervals (0.5,2,3) 2246
2246 as shown in **Fig. S5D**. 2247

2247 To impute phylogenetic branch lengths (**Methods: Im-** 2248
2248 **puting branch lengths**) for PATH transition inferences 2249
2249 (**Methods: Inferring cell state transitions from phy-** 2250
2250 **logenetic correlations**), we estimated a cell sampling rate 2251
2251 of 10^{-6} , which assumes that there were approximately 10^9 2252
2252 cells per tumor [Del Monte, 2009]. 2253

2253 Human patient glioblastoma

2254 Glioblastoma (GBM) phylogenies and corresponding scRNA- 2255
2255 seq data (including gene module scores) were obtained 2256
2256 from Chaligne et al. [2021]. Patient sample MGH105 was 2257
2257 chosen because tumor location was annotated, and patient 2258
2258 samples MGH115 and MGH122 were chosen because each 2259
2259 exhibited significant gene module transcriptional heritabil- 2260
2260 ity in the original paper. The MGH105 phylogeny is a 2261
2261 maximum-likelihood (ML) consensus tree, containing 80 2262
2262 cells, 20 cells from each location (MGH105A, MGH105B, 2263
2263 MGH105C, and MGH105D). Analyses of patient sample 2264
2264 MGH115 used 9 ML phylogeny search replicates for the same 2265
2265 38 cells from the original paper. Analyses of MGH122 used 2266
2266 10 ML phylogeny search replicates and the same 45 cells 2267
2267 from the original paper. Phylogenetic correlations were com- 2268
2268 puted by using the inverse node-distance weighting function 2269
2269 (**Methods: Phylogenetic correlations**). 2270

2270 PATH inferred transition rates (**Fig. 6F, Methods: In-** 2271
2271 **fering cell state transitions from phylogenetic corre-** 2272
2272 **lations**) were computed using categorical cell states (NPC- 2273
2273 /OPC-/AC-/MES-like), with states defined by the corre- 2274
2274 sponding per cell maximum module score, as in Chaligne 2275
2275 et al. [2021]. Note that, in the original paper, the NPC-like 2276
2276 and MES-like modules combine the NPC1-/NPC2-like and 2277
2277 MES1-/MES2-like modules, respectively. PATH inferred 2278
2278 transitions $\hat{P}(t = \tau)$ correspond to a time scale proportional 2279

2279 to the mean branch length distance separating cells one node
2280 apart, τ .

2281 Gene set enrichment analysis (GSEA) and Over-
2282 Representation Analysis (ORA) were performed using the
2283 functions `fgsea()` and `fora()` from the R software pack-
2284 age *fgsea* [Korotkevich et al., 2021]. For both analyses, the
2285 3,000 most variably transcribed genes (selected using the
2286 `SCTransform()` function from the R software package *Seu-*
2287 *rat* [Hao et al., 2021] on scRNAseq data) in patient sample
2288 MGH115 were ranked by their phylogeny-replicate mean
2289 phylogenetic auto-correlation z scores (Table S4).

2290 In both analyses, we measured the enrichment of gene sets
2291 from the chemical and genetic perturbation (C2:CGP) col-
2292 lection from the molecular signatures database (MSigDB)
2293 [Subramanian et al., 2005], as well as the GBM gene modules
2294 (NPC1-/NPC2-/OPC-/AC-/MES1-/MES2-like) defined in
2295 Neftel et al. [2019], and filtered out sets with fewer than 20
2296 genes. For both analyses (GSEA and ORA), pathway en-
2297 richment p-values were adjusted “padj” with the Benjami-
2298 Hochberg procedure (BH), to account for multiple compar-
2299 isons. Enriched pathways (BH adjusted $p < 0.05$) using
2300 GSEA that are presented in Fig. 6H were chosen manually
2301 (due to putative relevance) from a list of enriched pathways
2302 (Table S5).

2303 ORA was performed on two gene clusters (“Cluster 1”
2304 and “Cluster 2” in Fig. S6B), which were determined
2305 by hierarchical clustering, using Ward’s method, of the
2306 replicate-mean cross-correlations between the top 100 most
2307 significantly auto-correlated genes (across the phylogeny-
2308 replicates, see Table S4) in patient sample MGH115. All
2309 3,000 of the most variable genes were used to define the “uni-
2310 verse” or “background” genes to test for over-representation.
2311 All enriched gene sets (BH adjusted $p < 0.05$) for Cluster
2312 1, and a manually chosen subset for Cluster 2, are shown in
2313 Fig. S6B. A complete list of ORA enriched gene sets found
2314 in Clusters 1 and 2 from Fig. S6B can be found in Table
2315 S6.

2316 Gliomasphere phylogenies

2317 Patient-derived human GBM cells (MGG23) [Wakimoto
2318 et al., 2011] were grown in Neurobasal Medium (Thermo
2319 Fisher Scientific) supplemented with 1/2 x N2 and 1 x
2320 B27 (Thermo Fisher Scientific), 1% Penicillin/Streptomycin
2321 (Thermo Fisher Scientific), 1.5 x Glutamax (Thermo Fisher
2322 Scientific), 20 ng/mL of EGF and 20 ng/mL of FGF2
2323 (Shenandoah Biotechnology). The Molecular Recorder cas-
2324 sette PCT62 [Chan et al., 2019] was introduced into MGG23
2325 cells using piggyBac-mediated transposition (Systems Bio-
2326 sciences). Lineage tracing was initiated by infecting cells
2327 with lentivirus expressing Cas9-EGFP, followed by FACS
2328 sorting for EGFP-positive cells. Cells were subsequently
2329 grown *in vitro* for 4 weeks and lineage traced with the Molec-
2330 ular Recorder approach for two replicates. scRNAseq li-

2331 braries were generated using the Chromium Next GEM Sin-
2332 gle Cell GEM, Library & Gel Bead Kit v3.1, Chromium
2333 Single Cell Feature Barcode Library Kit, Chromium Next
2334 GEM Chip G, and 10x Chromium Controller (10x Ge-
2335 nomics) according to manufacturer instructions. Single-cell
2336 gene expression libraries were sequenced with paired-end,
2337 28 and 91-base reads on a NextSeq 2000 sequencer (Illu-
2338 mina). The Cas9-edited Molecular Recorder barcodes were
2339 PCR amplified from single-cell cDNA libraries as previously
2340 described [Chan et al., 2019] and sequenced with paired-
2341 end, 28 and 272-base reads on a NextSeq 2000 sequencer
2342 (Illumina). Phylogenies were reconstructed using *Cassiopeia*
2343 [Jones et al., 2020] using the `VanillaGreedySolver()` with
2344 default parameters for each subclone per replicate. ScR-
2345 NAseq data for each replicate were processed independently
2346 using the R package *Seurat* [Hao et al., 2021], by normal-
2347 izing and scaling RNA count data after subsetting for cells
2348 with $< 25\%$ mitochondrial DNA and > 200 RNA features.
2349 GBM gene modules [Neftel et al., 2019] were assigned using
2350 the *Seurat AddModuleScore()* function. Within each repli-
2351 cate, subclone phylogenies (3 for the first replicate and 6 for
2352 the second replicate) were joined at their roots before com-
2353 puting phylogenetic correlations. Phylogenetic correlations
2354 were computed for GBM gene modules using the one-node
2355 only weighting function, and z scores were computed analyt-
2356 ically per replicate. Replicate mean phylogenetic correlation
2357 z scores are shown in Fig. 6G.

2358 B-ALL analysis

2359 A blood sample was extracted from a 16yo B-ALL patient
2360 after treatment for four weeks with prednisone, daunoru-
2361 bicin, vincristine, and pegaspargase (AALL1131). Rare
2362 single persistent blasts were sorted into a 96 well plate
2363 based on dim expression of CD45 and CD19 positivity. In
2364 addition, CD10, CD20, CD34, and CD38 expression were
2365 recorded for each cell. An unsorted remission bone mar-
2366 row sample was used as a germline control. In addition,
2367 a pre-treatment unsorted bulk sample was obtained from
2368 the patient at the time of diagnosis. Eighty-six cells with
2369 *a priori* tumorigenic phenotype were amplified using primary
2370 template-directed amplification (PTA) protocol [Gonzalez-
2371 Pena et al., 2021]. Libraries were constructed with the
2372 Illumina DNA Prep with Enrichment Kit. All libraries
2373 were subjected to whole-exome sequencing at the Chan
2374 Zuckerberg Biohub on an Illumina NovaSeq6000. The un-
2375 enriched libraries were whole-genome sequenced at the New
2376 York Genome Center on an Illumina NovaSeq6000 platform.
2377 WGS reads were mapped to hg38 using BWA mem and fur-
2378 ther processed following GATK best practices guidelines
2379 [Van der Auwera and O’Connor, 2020]. Somatic single nu-
2380 cleotide variants (SNVs) were detected using an in-house
2381 pipeline combining cell genotyping based on *GATK Haplo-*
2382 *typeCaller* [Poplin et al., 2017] and somatic detection based
2383 on *Mutect2* [Cibulskis et al., 2013]. Cell H3 was removed
2384 from the WGS analysis given that it was suspected of being

2385 a replicate of H4 because WGS and WES allele frequencies
 2386 at exonic mutations of H3 did not match. Phylogenetic
 2387 trees were built with *CellPhy* [Kozlov et al., 2022] using
 2388 the SNV mutations which were not overlapping with deletions.
 2389 We detected haplotypic deletions (genomic regions
 2390 containing only the maternal or only the paternal haplotypes)
 2391 based on phasing of germline heterozygous SNPs
 2392 [Delaneau et al., 2019]. Large chromosomal gains were not
 2393 detected by cytogenetics analyses so we assumed our sam-
 2394 ples were mainly diploid for the deletion detection analysis.
 2395 Mutations were mapped to the phylogeny using *treemut*
 2396 (<https://github.com/NickWilliamsSanger/treemut>).
 2397 The phylogeny was time-scaled using *rtreefit*
 2398 (<https://github.com/NickWilliamsSanger/rtreefit>). FACS
 2399 data were analyzed using the R package *flowCore*. Fluor-

2400 rescence values were compensated and logicle-transformed.
 2401 Three cells were identified as healthy based on their pheno-
 2402 type, their lower mutation burden and chromosomal dele-
 2403 tions, and they were removed from the tree in order to only
 2404 analyze the tumor population. Fluorescence values were
 2405 discretized based on frequency using the R package *arules*.
 2406 Phylogenetic correlations were computed analytically on
 2407 the discretized fluorescence values using the inverse-node-
 2408 distance weighting (**Methods: Phylogenetic correla-**
 2409 **tions**). We also classified cells into three states based on
 2410 the discretized CD19 fluorescence (low: 1-2, medium: 3-4,
 2411 high: 5-6) and calculated PATH transition rates among
 2412 those states (**Methods: Inferring cell state transitions**
 2413 **from phylogenetic correlations**).

Supplemental Figures

Figure S1

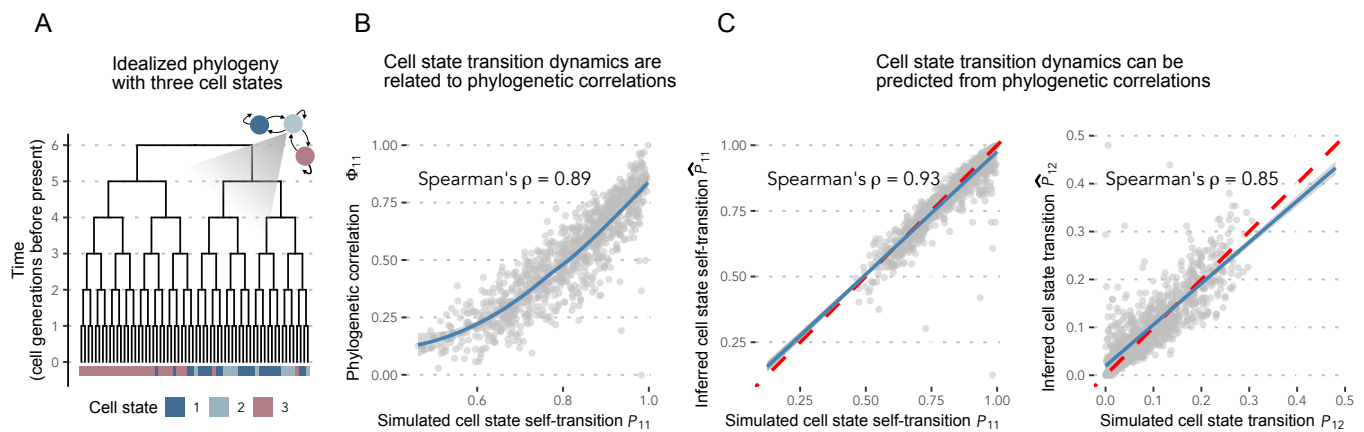


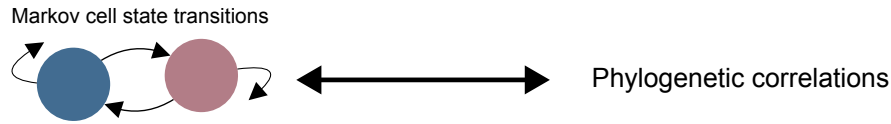
Figure S1: Cell state transition dynamics predict phylogenetic correlations

A) Simulated idealized phylogeny containing $2^6 = 64$ cells (**Methods: Simulating phylogenies**) in which cells can transition between three possible cell states. Cell state transitions are represented as a discrete-time Markov chain (**Methods: Markov model of cell state transitions**).

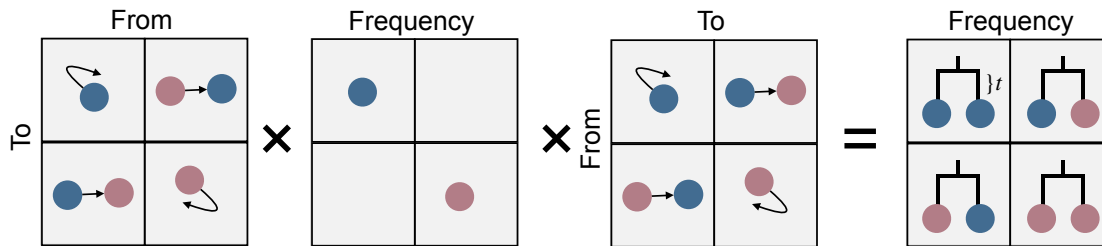
B) Simulated cell state transition dynamics (**Methods: Simulating phylogenies**) and measured phylogenetic auto-correlations (**Methods: Phylogenetic correlations**) for the first cell state for 1,000 independent simulations on idealized phylogenies, containing 64 cells as in A, in which state transition probabilities were randomly generated for each trial. Phylogenetic correlations were computed using a weighting function that included only sister cells (one-node only, as described in **Methods: Phylogenetic correlations and cell state transitions**). LOESS regression line (blue) with 95% confidence interval (light gray) is shown. Spearman's rank correlation coefficient = 0.89, $p < 2.2e - 16$.

C) (Left) Simulated versus PATH-inferred (**Methods: Inferring cell state transitions from phylogenetic correlations**), by transforming the phylogenetic auto-correlations measured in B, cell state self-transition (*i.e.*, stability) probabilities. Spearman's rank correlation coefficient 0.93, $p < 2.2e-16$. (Right) Simulated versus PATH-inferred (**Methods: Inferring cell state transitions from phylogenetic correlations**) cell state transition probabilities from state 1 to 2, on idealized phylogenies (**Methods: Simulating phylogenies**). Spearman's rank correlation coefficient 0.85, $p < 2.2e-16$. Dashed red lines both have slope 1 and pass through the origin. Linear regression lines (blue) with 95% confidence intervals (light gray) are shown for both plots.

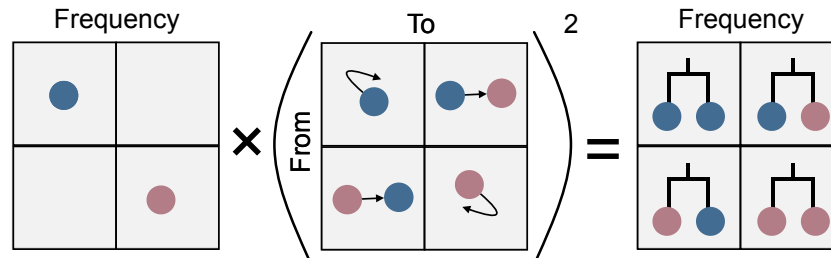
Box S1: Cell state transition dynamics and phylogenetic correlations



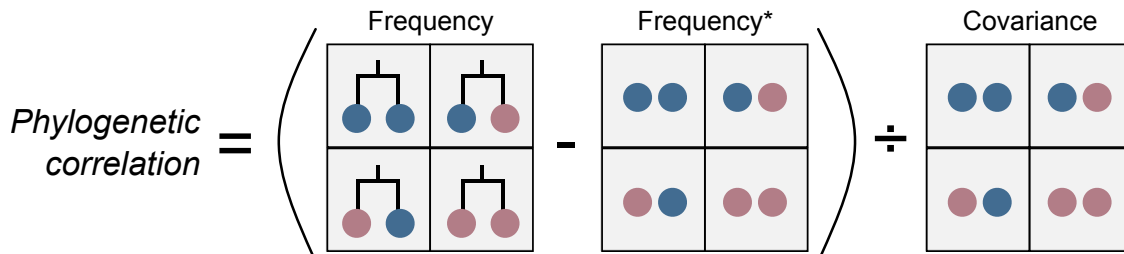
We can connect cell state transition dynamics (P^t) to phylogenetic cell state pair frequencies $F(t)$, for a given ancestral relationship t (e.g., sister cells [i.e., $t = 1$] or first-cousins [i.e., $t = 2$]) with, $(P^t)^T D P^t = F(t)$, where $D = \text{diag}(\mu)$, is the diagonal matrix of cell state frequencies, and T signifies the matrix transpose. This relation, for two cell states, is illustrated below.



For reversible Markov dynamics, this mathematical relation simplifies to, $DP^{2t} = F(t)$.



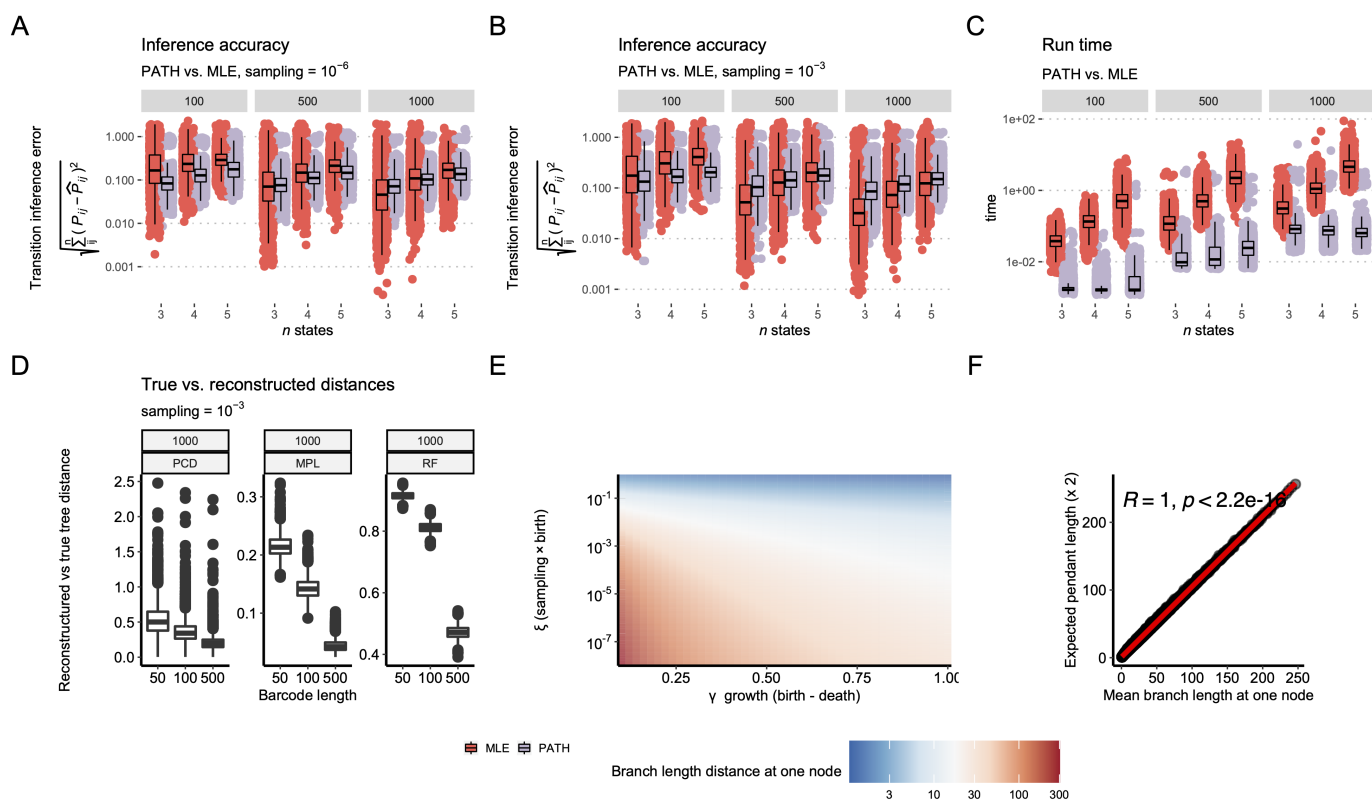
State pair frequencies can be transformed into phylogenetic correlations $\Phi(t)$, by standardizing: $\Phi(t) = (F(t) - \mu\mu^T)/(\sigma\sigma^T)$, with $\sigma^2 = \mu - \mu^2$.



*This matrix represents the sampling probabilities – with replacement – of observing illustrated cell state pairs. Similarly, the covariance matrix represents population covariances.

Finally, for reversible dynamics, state transitions can be directly inferred from state pair frequencies, $P^{2t} = D^{-1}F(t)$.

Figure S2



2437

2438

Figure S2: PATH inferences and simulations of somatic evolution

2439

A) Transition inference error (Euclidean distance between inferred and true transition probabilities) using PATH or MLE for 3, 4, or 5 cell states in a phylogeny composed of either 100 (left), 500 (middle), or 1,000 (right) cells, representing a sample of 10^{-6} of the total population. Each parameter combination was simulated 1,000 times and inferences are shown for all simulations in which neither PATH nor MLE inference failed.

2440

2441

2442

2443

B) Same as **A** but with a sampling rate of 10^{-3} .

2444

C) Run times corresponding to simulations depicted in **A**.

2445

D) Phylogenetic correlation difference (PCD, left), Mean Path Length distance (MPL) [Steel and Penny, 1993] (center), and Robinson-Foulds distance (RF) [Robinson and Foulds, 1981] (right) between simulated true and reconstructed phylogenies (Methods: Phylogenetic reconstruction). Phylogenies were simulated 1,000 times for each barcode length (x-axis).

2446

2447

2448

E) Expected pendant edge lengths for a sampled somatic evolutionary process, as a function of birth, death and sampling rates (Methods: Imputing branch lengths).

2449

2450

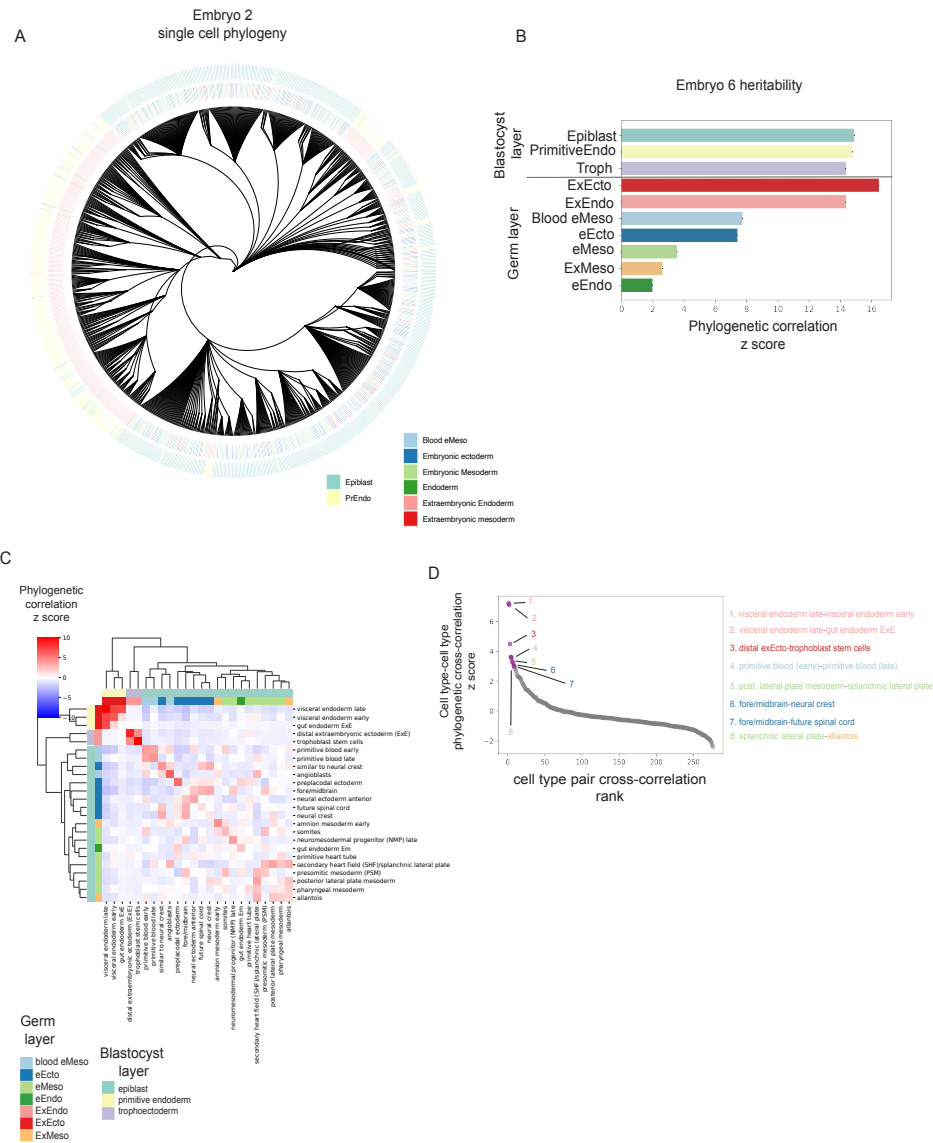
F) Correspondence between simulated branch lengths at a node depth of one and expected pendant lengths, while varying sampled somatic evolutionary process parameters.

2451

2452

2453

Figure S3



2454

2455

Figure S3: PATH quantifies ancestry and divergence of germ layers and cell types during mouse embryogenesis

2456

A) Single-cell phylogeny for mouse embryo 2 from Chan et al. [2022], containing 700 of 1,113 randomly chosen cells for visualization. Each leaf represents a single cell. Leaves are colored by their assignment to a blastocyst or germ layer of origin based on transcription profiles. e prefix, embryonic; ex prefix, extraembryonic. PrEndo, primitive endoderm.

2457

2458

2459

B) Blastocyst and germ layer phylogenetic auto-correlations for embryo 6 ($N = 1,722$ cells).

2460

C) Hierarchical clustering of tissue types in embryo 6 by phylogenetic correlation using Ward's method. Only tissues with more than 30 cells present in the sample were considered for analysis. Tissues colored by their germ layer and blastocyst layer of origin. ExE, extraembryonic; EM, embryonic.

2461

2462

2463

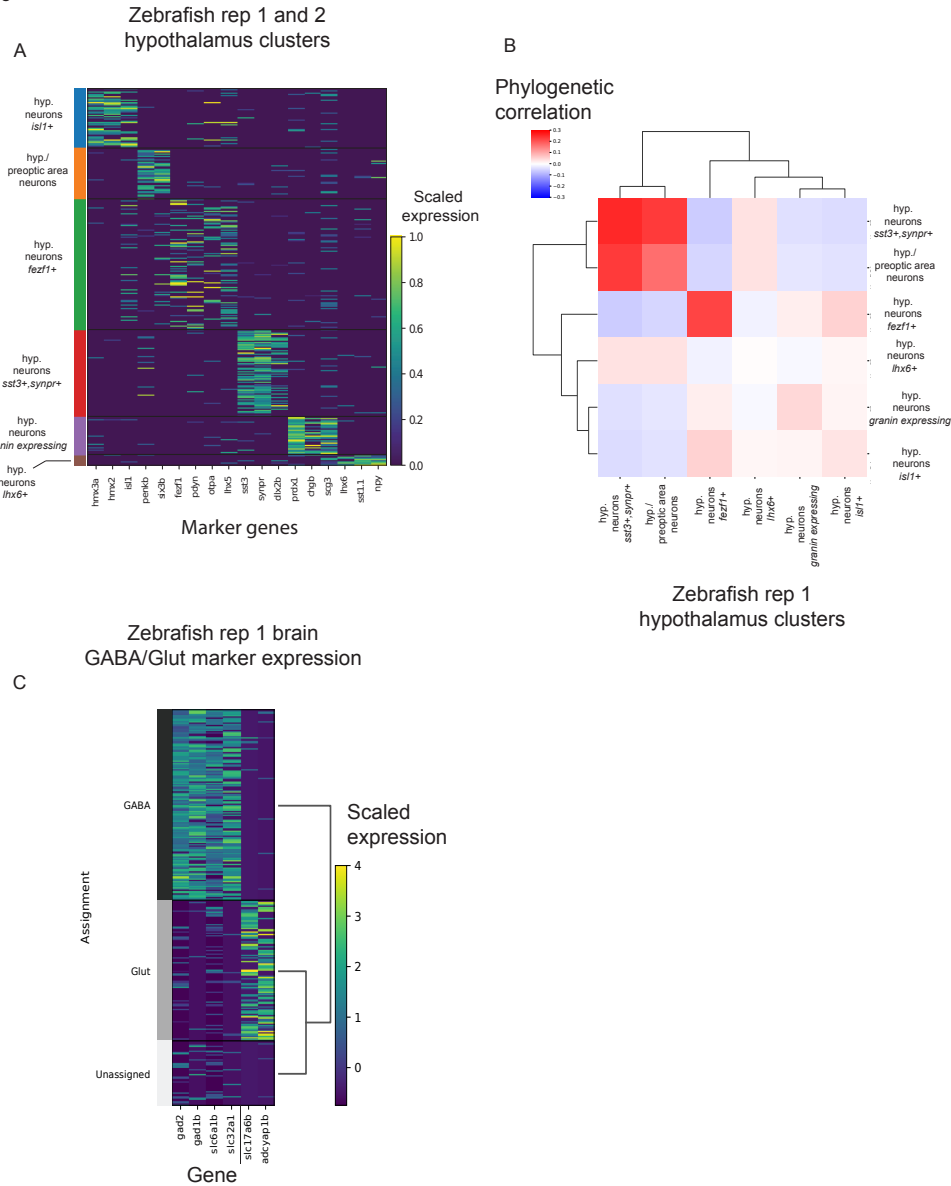
D) Ranked pairwise cell type phylogenetic correlations (z scores) for embryo 6. Pairs with z scores > 3 highlighted. Text colored by germ layer as in **B**.

2464

2465

2466

Figure S4



2467

2468

Figure S4: PATH identifies cell fate-determining factors across anatomical, defined tissue and gene expression layers during neurogenesis in zebrafish

2469

2470

A) Heat map of scaled expression of representative marker genes across hypothalamus clusters. Marker genes and clusters were defined by Raj et al. [2018].

2471

2472

B) Hypothalamus cluster (from **A**) phylogenetic correlations.

2473

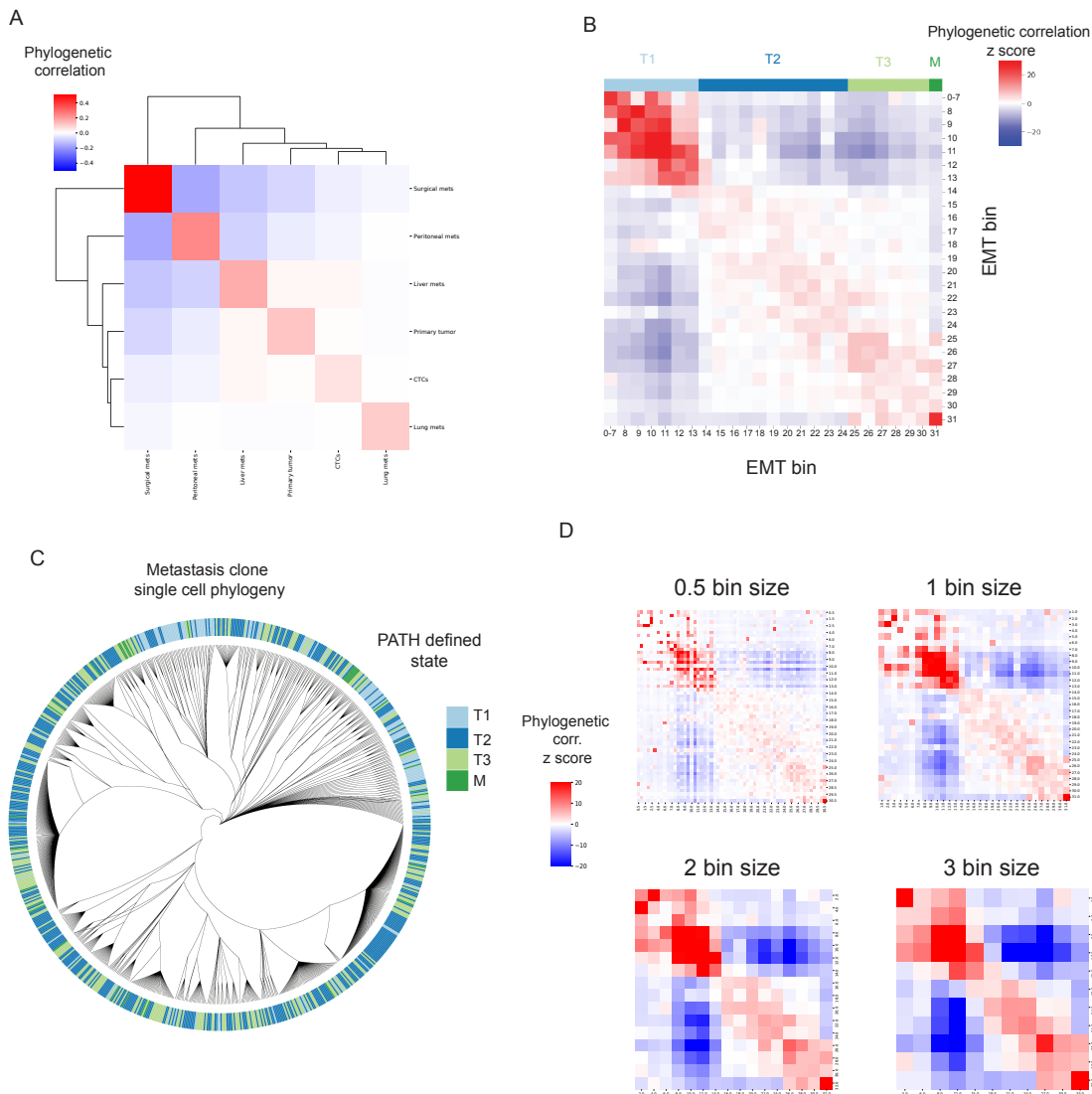
C) Heat map of GABA markers (*gad2*, *gad1b*, *slc6a1b*, *slc32a1*) and Glut (*slc17a6b*, *adcypap1b*) signaling in forebrain neurons of zebrafish replicate 1 (see **Methods** for assignment of cells into GABA, Glutamatergic (Glut) and Unassigned categories).

2474

2475

2476

Figure S5



2477

2478

Figure S5: Quantifying the heritability versus plasticity of EMT transcriptional states

2479

A) Tumor cell harvest site phylogenetic correlations.

2480

B) EMT bin phylogenetic correlations (z scores). Colors represent putative states. Full table of EMT bin phylogenetic correlations of can be found in **Table S3**.

2481

2482

C) Single-cell phylogeny from mouse 1, clone 1 from [Simeonov et al. \[2021\]](#), containing 700 of 7,968 randomly chosen cells for visualization. Each leaf represents a single cell. Cells are colored by PATH-defined states (T1, T2, T3, M).

2483

2484

D) EMT bin phylogenetic correlation (z score) heat maps using different bin sizes (0.5, 1, 2, 3).

2485

2486

Figure S6

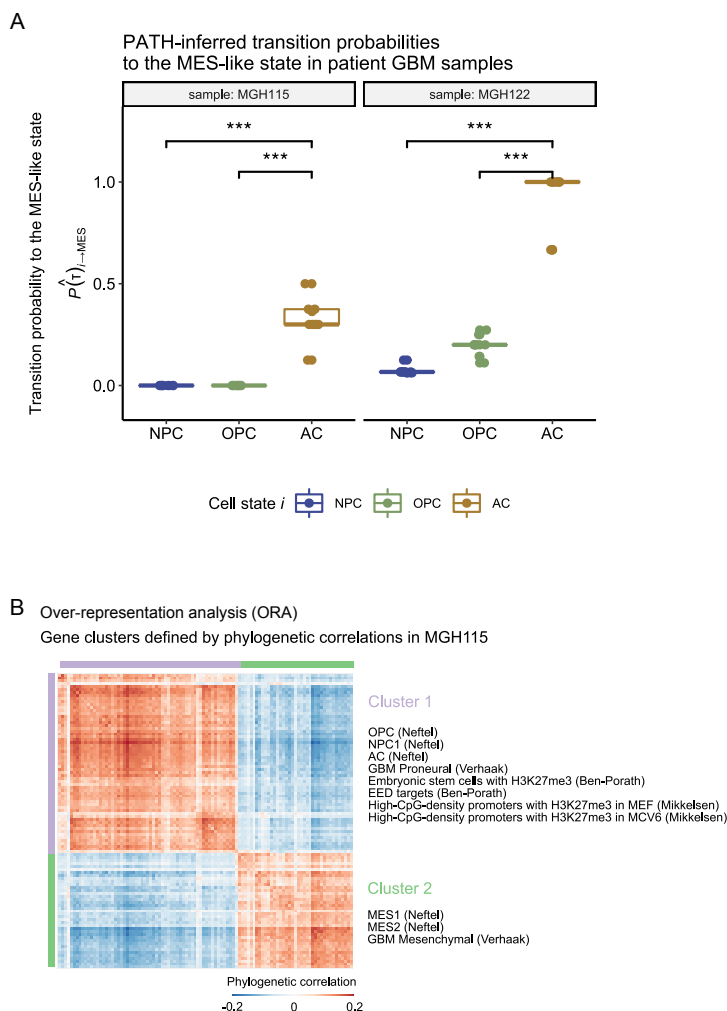


Figure S6: PATH inferred cell state transitions and gene set enrichment in human glioblastoma

A) PATH-inferred transition probabilities $\hat{P}(\tau)$ (**Methods: Inferring cell state transitions from phylogenetic correlations**) from neurodevelopmental-like (NPC-/OPC-/AC-like) cell states to the MES-like cell state in human patient-derived GBM samples MGH115 and MGH122 (**Methods: Human patient glioblastoma**). Points correspond to PATH inferences for each sample phylogeny-replicate per sample. Significance determined by two-sided t-test ($p < 9.7e-6$ and $p < 8.2e-9$ for NPC-like vs AC-like in MGH115 and MGH122 respectively; $p < 9.7e-6$ and $p < 7.8e-9$ for OPC-like vs AC-like in MGH115 and MGH122, respectively). Colors correspond to cell state.

B) Heat map of the phylogeny-replicate mean phylogenetic correlations (**Methods: Phylogenetic correlations**) for the top 100 most heritable genes (determined by phylogeny-replicate mean gene phylogenetic auto-correlation z scores) in MGH115. Over-representation analysis (ORA) performed on the genes in each of the two clusters, defined by hierarchical clustering using Ward's method, separately. Phylogenetic correlations were computed using an inverse-node-distance weighting (**Methods: Human patient glioblastoma**). Only select gene sets are depicted for Cluster 2; remaining significantly enriched gene sets are in **Table S6**.

GBM gene modules (NPC-/OPC-/AC-/MES-like) were shortened to (NPC/OPC/AC/MES).

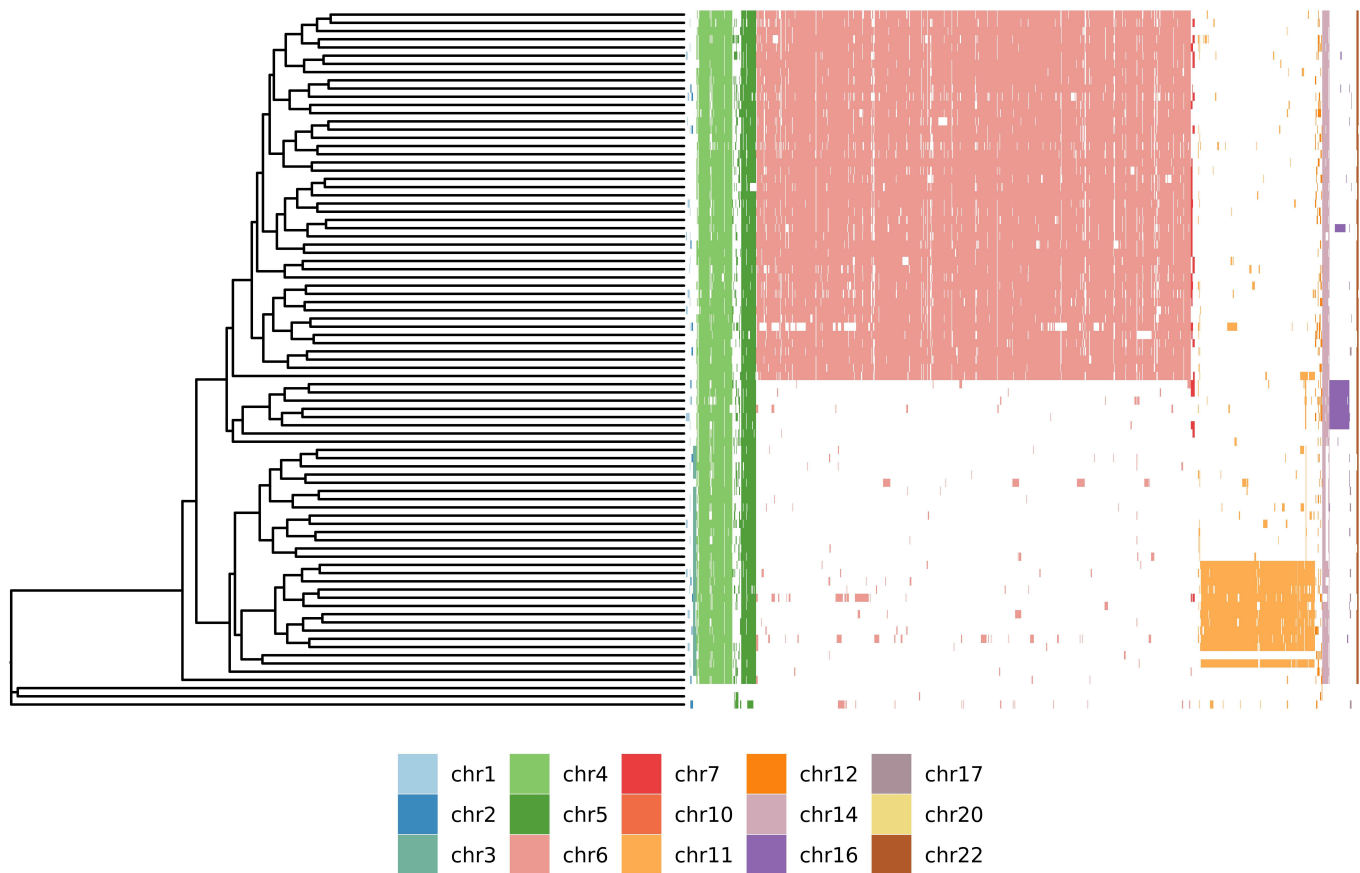


Figure S7: Quantifying cell state heterogeneity in B-ALL using single-cell whole genome sequencing

Genome-wide copy-number deletion annotations projected onto the B-ALL single-cell phylogeny from **Fig. 7A**.

References

- Rene C Adam and Elaine Fuchs. The yin and yang of chromatin dynamics in stem cell fate selection. *Trends Genet.*, 32(2):89–100, February 2016. 2510
- Arianna Baggiolini, Scott J Callahan, Emily Montal, Joshua M Weiss, Tuan Trieu, Mohita M Tagore, Sam E Tischfield, Ryan M Walsh, Shruthy Suresh, Yujie Fan, Nathaniel R Campbell, Sarah C Perlee, Nathalie Saurat, Miranda V Hunter, Theresa Simon-Vermot, Ting-Hsiang Huang, Yilun Ma, Travis Hollmann, Satish K Tickoo, Barry S Taylor, Ekta Khurana, Richard P Koche, Lorenz Studer, and Richard M White. Developmental chromatin programs determine oncogenic competence in melanoma. *Science*, 373(6559):eabc1048, September 2021. 2514
- Rachel C Bandler, Ilaria Vitali, Ryan N Delgado, May C Ho, Elena Dvoretzskova, Josue S Ibarra Molinas, Paul W Frazel, Maesoumeh Mohammadkhani, Robert Machold, Sophia Maedler, Shane A Liddelow, Tomasz J Nowakowski, Gord Fishell, and Christian Mayer. Single-cell delineation of lineage and genetic identity in the mouse brain. *Nature*, 601(7893):404–409, December 2021. 2519
- Chloé S Baron and Alexander van Oudenaarden. Unravelling cellular relationships during development and regeneration using genetic lineage tracing. *Nat. Rev. Mol. Cell Biol.*, 20(12):753–765, December 2019. 2522
- Charles C Bell, Katie A Fennell, Yih-Chih Chan, Florian Rambow, Miriam M Yeung, Dane Vassiliadis, Luis Lara, Paul Yeh, Luciano G Martelotto, Aljosja Rogiers, Brandon E Kremer, Olena Barbash, Helai P Mohammad, Timothy M Johanson, Marian L Burr, Arindam Dhar, Natalie Karpnich, Luyi Tian, Dean S Tyler, Laura MacPherson, Junwei Shi, Nathan Pinnawala, Chun Yew Fong, Anthony T Papenfuss, Sean M Grimmond, Sarah-Jane Dawson, Rhys S Allan, Ryan G Kruger, Christopher R Vakoc, David L Goode, Shalin H Naik, Omer Gilan, Enid Y N Lam, Jean-Christophe Marine, Rab K Prinjha, and Mark A Dawson. Targeting enhancer switching overcomes non-genetic drug resistance in acute myeloid leukaemia. *Nat. Commun.*, 10(1):2723, June 2019. 2524
- Marco Benevento, Tomas Hökfelt, and Tibor Harkany. Ontogenetic rules for the molecular diversification of hypothalamic neurons. *Nat. Rev. Neurosci.*, 23(10):611–627, October 2022. 2531
- Brent A Bidby, Wenjun Kong, Kenji Kamimoto, Chuner Guo, Sarah E Waye, Tao Sun, and Samantha A Morris. Single-cell mapping of lineage and identity in direct reprogramming. *Nature*, 564(7735):219–224, December 2018. 2533
- Lacramioara Bintu, John Yong, Yaron E Antebi, Kayla McCue, Yasuhiro Kazuki, Narumi Uno, Mitsuo Oshimura, and Michael B Elowitz. Dynamics of epigenetic regulation at the single-cell level. *Science*, 351(6274):720–724, February 2016. 2536
- S P Blomberg and T Garland. Tempo and mode in evolution: phylogenetic inertia, adaptation and comparative methods, 2002. 2538
- Simon P Blomberg, Theodore Garland, Jr, and Anthony R Ives. Testing for phylogenetic signal in comparative data: behavioral traits are more labile. *Evolution*, 57(4):717–745, April 2003. 2541
- Laurie A Boyer, Kathrin Plath, Julia Zeitlinger, Tobias Brambrink, Lea A Medeiros, Tong Ihn Lee, Stuart S Levine, Marius Wernig, Adriana Tajonar, Mridula K Ray, George W Bell, Arie P Otte, Miguel Vidal, David K Gifford, Richard A Young, and Rudolf Jaenisch. Polycomb complexes repress developmental regulators in murine embryonic stem cells. *Nature*, 441(7091):349–353, May 2006. 2543
- Samuel W Brady, Kathryn G Roberts, Zhaohui Gu, Lei Shi, Stanley Pounds, Deqing Pei, Cheng Cheng, Yunfeng Dai, Meenakshi Devidas, Chunxu Qu, et al. The genomic landscape of pediatric acute lymphoblastic leukemia. *Nature genetics*, 54(9):1376–1389, 2022. 2547
- Yehuda Brody, Robert J Kimmerling, Yosef E Maruvka, David Benjamin, Joshua J Elacqua, Nicholas J Haradhvala, Jaegil Kim, Kent W Mouw, Kristjana Frangaj, Amnon Koren, Gad Getz, Scott R Manalis, and Paul C Blainey. Quantification of somatic mutation flow across individual cell division events by lineage sequencing. *Genome Res.*, 28(12):1901–1918, December 2018. 2549
- Ronan Chaligne, Federico Gaiti, Dana Silverbush, Joshua S Schiffman, Hannah R Weisman, Lloyd Kluegel, Simon Gritsch, Sunil D Deochand, L Nicolas Gonzalez Castro, Alyssa R Richman, Johanna Klughammer, Tommaso Biancalani, Christoph Muus, Caroline Sheridan, Alicia Alonso, Franco Izzo, Jane Park, Orit Rozenblatt-Rosen, Aviv Regev, Mario L Suvà, and Dan A Landau. Epigenetic encoding, heritability and plasticity of glioma transcriptional cell states. *Nat. Genet.*, 53(10):1469–1479, October 2021. 2554

- Joseph M Chan, Samir Zaidi, Jillian R Love, Jimmy L Zhao, Manu Setty, Kristine M Wadosky, Anuradha Gopalan, Zi-Ning Choo, Sitara Persad, Jungmin Choi, Justin LaClair, Kayla E Lawrence, Ojasvi Chaudhary, Tianhao Xu, Ignas Masilionis, Irina Linkov, Shangqian Wang, Cindy Lee, Afsar Barlas, Michael J Morris, Linas Mazutis, Ronan Chaligne, Yu Chen, David W Goodrich, Wouter R Karthaus, Dana Pe'er, and Charles L Sawyers. Lineage plasticity in prostate cancer depends on JAK/STAT inflammatory signaling. *Science*, 377(6611):1180–1191, September 2022. 2558–2562
- Michelle M Chan, Zachary D Smith, Stefanie Grosswendt, Helene Kretzmer, Thomas M Norman, Britt Adamson, Marco Jost, Jeffrey J Quinn, Dian Yang, Matthew G Jones, Alex Khodaverdian, Nir Yosef, Alexander Meissner, and Jonathan S Weissman. Molecular recording of mammalian embryogenesis. *Nature*, 570(7759):77–82, June 2019. 2563–2565
- Rony Chanoch-Myers, Adi Wider, Mario L Suva, and Itay Tirosch. Elucidating the diversity of malignant mesenchymal states in glioblastoma by integrative analysis. *Genome Med.*, 14(1):106, September 2022. 2566–2567
- Bjoern Chapuy, Chip Stewart, Andrew J Dunford, Jaegil Kim, Atanas Kamburov, Robert A Redd, Mike S Lawrence, Margaretha GM Roemer, Amy J Li, Marita Ziepert, et al. Molecular subtypes of diffuse large b cell lymphoma are associated with distinct pathogenic mechanisms and outcomes. *Nature medicine*, 24(5):679–690, 2018. 2568–2570
- Yanguang Chen. A new methodology of spatial cross-correlation analysis. *PLoS One*, 10(5):e0126158, May 2015. 2571
- James M Cheverud and Malcolm M Dow. An autocorrelation analysis of genetic variation due to lineal fission in social groups of rhesus macaques, 1985. 2572–2573
- Kristian Cibulskis, Michael S Lawrence, Scott L Carter, Andrey Sivachenko, David Jaffe, Carrie Sougnez, Stacey Gabriel, Matthew Meyerson, Eric S Lander, and Gad Getz. Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nature biotechnology*, 31(3):213–219, 2013. 2574–2576
- Raymond L Czaplewski and Robin M Reich. *Expected Value and Variance of Moran's Bivariate Spatial Autocorrelation Statistic for a Permutation Test*. 1993. 2577–2578
- Marco L Davila and Renier J Brentjens. Cd19-targeted car t cells as novel cancer immunotherapy for relapsed or refractory b-cell acute lymphoblastic leukemia. *Clinical advances in hematology & oncology: H&O*, 14(10):802, 2016. 2579–2580
- Ugo Del Monte. Does the cell number 10(9) still really fit one gram of tumor tissue? *Cell Cycle*, 8(3):505–506, February 2009. 2581–2582
- Olivier Delaneau, Jean-François Zagury, Matthew R Robinson, Jonathan L Marchini, and Emmanouil T Dermitzakis. Accurate, scalable and integrative haplotype estimation. *Nature communications*, 10(1):5436, 2019. 2583–2584
- David Detomaso and Nir Yosef. Hotspot identifies informative gene modules across modalities of single-cell genomics. *Cell Systems*, 12(5):446–456.e9, May 2021. 2585–2586
- José Alexandre F Diniz-Filho, Thiago Santos, Thiago Fernando Rangel, and Luis Mauricio Bini. A comparison of metrics for estimating phylogenetic signal under alternative evolutionary models. *Genet. Mol. Biol.*, 35(3):673–679, July 2012. 2587–2588
- Giacomo Donati, Emanuel Rognoni, Toru Hiratsuka, Kifayathullah Liakath-Ali, Esther Hoste, Gozde Kar, Melis Kayikci, Roslin Russell, Kai Kretzschmar, Klaas W Mulder, Sarah A Teichmann, and Fiona M Watt. Wounding induces dedifferentiation of epidermal gata6 cells and acquisition of stem cell properties. *Nat. Cell Biol.*, 19(6):603–613, June 2017. 2589–2592
- Anushka Dongre and Robert A Weinberg. New insights into the mechanisms of epithelial-mesenchymal transition and implications for cancer. *Nat. Rev. Mol. Cell Biol.*, 20(2):69–84, February 2019. 2593–2594
- Mohammed El-Kebir, Gryte Satas, and Benjamin J Raphael. Inferring parsimonious migration histories for metastatic cancers. *Nat. Genet.*, 50(5):718–726, May 2018. 2595–2596
- Weixiang Fang, Claire M Bell, Abel Sapirstein, Soichiro Asami, Kathleen Leeper, Donald J Zack, Hongkai Ji, and Reza Kalhor. Quantitative fate mapping: A general framework for analyzing progenitor state dynamics via retrospective lineage barcoding. *Cell*, 185(24):4604–4620, 2022. 2597–2599
- Katie A Fennell, Dane Vassiliadis, Enid Y N Lam, Luciano G Martelotto, Jesse J Balic, Sebastian Hollizeck, Tom S Weber, Timothy Semple, Qing Wang, Denise C Miles, Laura MacPherson, Yih-Chih Chan, Andrew A Guirguis, Lev M Kats, Emily S Wong, Sarah-Jane Dawson, Shalin H Naik, and Mark A Dawson. Non-genetic determinants of malignant clonal fitness at single-cell resolution. *Nature*, 601(7891):125–131, January 2022. 2600–2603

- Federico Gaiti, Ronan Chaligne, Hongcang Gu, Ryan M Brand, Steven Kothén-Hill, Rafael C Schulman, Kirill Grigorev, Davide Risso, Kyu-Tae Kim, Alessandro Pastore, Kevin Y Huang, Alicia Alonso, Caroline Sheridan, Nathaniel D Omans, Evan Biederstedt, Kendell Clement, Lili Wang, Joshua A Felsenfeld, Erica B Bhavsar, Martin J Aryee, John N Allan, Richard Furman, Andreas Gnirke, Catherine J Wu, Alexander Meissner, and Dan A Landau. Epigenetic evolution and lineage histories of chronic lymphocytic leukaemia. *Nature*, 569(7757):576–580, May 2019. 2604–2608
- John H Gillespie. *Population Genetics: A Concise Guide*. Johns Hopkins University Press, July 2004. 2609
- John L Gittleman and Mark Kot. Adaptation: Statistics and a null model for estimating phylogenetic effects. *Syst. Biol.*, 39(3):227–241, September 1990. 2610–2611
- Anita Gola and Elaine Fuchs. Environmental control of lineage plasticity and stem cell memory. *Curr. Opin. Cell Biol.*, 69:88–95, April 2021. 2612–2613
- Veronica Gonzalez-Pena, Sivaraman Natarajan, Yuntao Xia, David Klein, Robert Carter, Yakun Pang, Bridget Shaner, Kavya Annu, Daniel Putnam, Wenan Chen, et al. Accurate genomic variant detection in single cells with primary template-directed amplification. *Proceedings of the National Academy of Sciences*, 118(24):e2024176118, 2021. 2614–2616
- Geoffrey Grimmett and David Stirzaker. *Probability and random processes*. Oxford university press, 2020. 2617
- Stefanie Grosswendt, Helene Kretzmer, Zachary D Smith, Abhishek Sampath Kumar, Sara Hetzel, Lars Wittler, Sven Klages, Bernd Timmermann, Shankar Mukherji, and Alexander Meissner. Epigenetic regulator function through mouse gastrulation. *Nature*, 584(7819):102–108, August 2020. 2618–2620
- Gunes Gundem, Peter Van Loo, Barbara Kremeyer, Ludmil B Alexandrov, Jose M C Tubio, Elli Papaemmanuil, Daniel S Brewer, Heini M L Kallio, Gunilla Högnäs, Matti Annala, Kati Kivinummi, Victoria Goody, Calli Latimer, Sarah O’Meara, Kevin J Dawson, William Isaacs, Michael R Emmert-Buck, Matti Nykter, Christopher Foster, Zsofia Kote-Jarai, Douglas Easton, Hayley C Whitaker, ICGC Prostate Group, David E Neal, Colin S Cooper, Rosalind A Eeles, Tapio Visakorpi, Peter J Campbell, Ultan McDermott, David C Wedge, and G Steven Bova. The evolutionary history of lethal metastatic prostate cancer. *Nature*, 520(7547):353–357, April 2015. 2621–2626
- Richard P Halley-Stott and John B Gurdon. Epigenetic memory in the context of nuclear reprogramming and cancer. *Brief. Funct. Genomics*, 12(3):164–173, May 2013. 2627–2628
- Timothy R Hammond, Connor Dufort, Lasse Dissing-Olesen, Stefanie Giera, Adam Young, Alec Wysoker, Alec J Walker, Frederick Gergits, Michael Segel, James Nemesh, Samuel E Marsh, Arpiar Saunders, Evan Macosko, Florent Ginhoux, Jimmiao Chen, Robin J M Franklin, Xianhua Piao, Steven A McCarroll, and Beth Stevens. Single-Cell RNA sequencing of microglia throughout the mouse lifespan and in the injured brain reveals complex Cell-State changes, 2019. 2629–2632
- Douglas Hanahan. Hallmarks of cancer: New dimensions. *Cancer Discov.*, 12(1):31–46, January 2022. 2633
- Thomas F Hansen and Emília P Martins. Translating between microevolutionary process and macroevolutionary patterns: the correlation structure of interspecific data. *Evolution*, 50(4):1404–1417, August 1996. 2634–2635
- Yuhan Hao, Stephanie Hao, Erica Andersen-Nissen, William M Mauck, 3rd, Shiwei Zheng, Andrew Butler, Maddie J Lee, Aaron J Wilk, Charlotte Darby, Michael Zager, Paul Hoffman, Marlon Stoeckius, Eftymia Papalexi, Eleni P Mimitou, Jaison Jain, Avi Srivastava, Tim Stuart, Lamar M Fleming, Bertrand Yeung, Angela J Rogers, Juliana M McElrath, Catherine A Blish, Raphael Gottardo, Peter Smibert, and Rahul Satija. Integrated analysis of multimodal single-cell data. *Cell*, 184(13):3573–3587.e29, June 2021. 2636–2640
- Toshiro Hara, Rony Chanoch-Myers, Nathan D Mathewson, Chad Myskiw, Lyla Atta, Lillian Bussema, Stephen W Eichhorn, Alissa C Greenwald, Gabriela S Kinker, Christopher Rodman, L Nicolas Gonzalez Castro, Hiroaki Wakimoto, Orit Rozenblatt-Rosen, Xiaowei Zhuang, Jean Fan, Tony Hunter, Inder M Verma, Kai W Wucherpfennig, Aviv Regev, Mario L Suvà, and Itay Tirosh. Interactions between cancer cells and immune cells drive transitions to mesenchymal-like states in glioblastoma. *Cancer Cell*, 39(6):779–792.e11, June 2021. 2641–2645
- Sahand Hormoz, Nicolas Desprat, and Boris I Shraiman. Inferring epigenetic dynamics from kin correlations. *Proc. Natl. Acad. Sci. U. S. A.*, 112(18):E2281–9, May 2015. 2646–2647
- Sahand Hormoz, Zakary S Singer, James M Linton, Yaron E Antebi, Boris I Shraiman, and Michael B Elowitz. Inferring Cell-State transition dynamics from lineage trees and endpoint Single-Cell measurements. *Cell Syst*, 3(5):419–433.e8, November 2016. 2648–2650

- Bahram Houchmandzadeh, Eric Wieschaus, and Stanislas Leibler. Establishment of developmental precision and proportions in the early drosophila embryo. *Nature*, 415(6873):798–802, February 2002. 2651
2652
- Jacob Househam, Timon Heide, George D Cresswell, Inmaculada Spiteri, Chris Kimberley, Luis Zapata, Claire Lynn, Chela James, Maximilian Mossner, Javier Fernandez-Mateos, Alessandro Vinceti, Ann-Marie Baker, Calum Gabbutt, Alison Berner, Melissa Schmidt, Bingjie Chen, Eszter Lakatos, Vinaya Gunasri, Daniel Nichol, Helena Costa, Miriam Mitchinson, Daniele Ramazzotti, Benjamin Werner, Francesco Iorio, Marnix Jansen, Giulio Caravagna, Chris P Barnes, Darryl Shibata, John Bridgewater, Manuel Rodriguez-Justo, Luca Magnani, Andrea Sottoriva, and Trevor A Graham. Phenotypic plasticity and genetic control in colorectal cancer evolution. *Nature*, 611(7937):744–753, November 2022. 2653
2654
2655
2656
2657
2658
- Zheng Hu, Jie Ding, Zhicheng Ma, Ruping Sun, Jose A Seoane, J Scott Shaffer, Carlos J Suarez, Anna S Berghoff, Chiara Cremolini, Alfredo Falcone, Fotios Loupakakis, Peter Birner, Matthias Preusser, Heinz-Josef Lenz, and Christina Curtis. Quantitative evidence for early metastatic seeding in colorectal cancer. *Nat. Genet.*, 51(7):1113–1122, July 2019. 2659
2660
2661
- Max Jan, Thomas M Snyder, M Ryan Corces-Zimmerman, Paresh Vyas, Irving L Weissman, Stephen R Quake, and Ravindra Majeti. Clonal evolution of preleukemic hematopoietic stem cells precedes human acute myeloid leukemia. *Sci. Transl. Med.*, 4(149):149ra118, August 2012. 2662
2663
2664
- Hans Erik Johnsen, Kim Steve Bergkvist, Alexander Schmitz, Malene Krag Kjeldsen, Steen Møller Hansen, Michael Gaihede, Martin Agge Nørgaard, John Bæch, Marie-Louise Grønholdt, Frank Svendsen Jensen, et al. Cell of origin associated classification of b-cell malignancies by gene signatures of the normal b-cell hierarchy. *Leukemia & lymphoma*, 55(6):1251–1260, 2014. 2665
2666
2667
2668
- Matthew G Jones, Alex Khodaverdian, Jeffrey J Quinn, Michelle M Chan, Jeffrey A Hussmann, Robert Wang, Chenling Xu, Jonathan S Weissman, and Nir Yosef. Inference of single-cell phylogenies from lineage tracing data using cassiopeia. *Genome Biol.*, 21(1):92, April 2020. 2669
2670
2671
- Matthew G Jones, Yanay Rosen, and Nir Yosef. Interactive, integrated analysis of single-cell transcriptomic and phylogenetic data with PhyloVision. *Cell Rep Methods*, 2(4):100200, April 2022. 2672
2673
- Panagiotis Karras, Ignacio Bordeu, Joanna Pozniak, Ada Nowosad, Cecilia Pazzi, Nina Van Raemdonck, Ewout Landeloos, Yannick Van Herck, Dennis Pedri, Greet Bervoets, Samira Makhzami, Jia Hui Khoo, Benjamin Pavie, Jochen Lamote, Oskar Marin-Bejar, Michael Dewaele, Han Liang, Xingju Zhang, Yichao Hua, Jasper Wouters, Robin Browaeys, Gabriele Bergers, Yvan Saey, Francesca Bosisio, Joost van den Oord, Diether Lambrechts, Anil K Rustgi, Oliver Bechter, Cedric Blanpain, Benjamin D Simons, Florian Rambow, and Jean-Christophe Marine. A cellular hierarchy in melanoma uncouples growth and metastasis. *Nature*, 610(7930):190–198, October 2022. 2674
2675
2676
2677
2678
2679
- Gennady Korotkevich, Vladimir Sukhov, Nikolay Budin, Boris Shpak, Maxim N Artyomov, and Alexey Sergushichev. Fast gene set enrichment analysis. 2021. 2680
2681
- Alexey Kozlov, Joao M Alves, Alexandros Stamatakis, and David Posada. Cellphy: accurate and fast probabilistic inference of single-cell phylogenies from scdna-seq data. *Genome biology*, 23(1):1–30, 2022. 2682
2683
- Gloria S Kwon, Manuel Viotti, and Anna-Katerina Hadjantonakis. The endoderm of the mouse embryo arises by dynamic widespread intercalation of embryonic and extraembryonic lineages. *Dev. Cell*, 15(4):509–520, October 2008. 2684
2685
- Dan R Laks, Thomas J Crisman, Michelle Y S Shih, Jack Mottahedeh, Fuying Gao, Jantzen Sperry, Matthew C Garrett, William H Yong, Timothy F Cloughesy, Linda M Liao, Albert Lai, Giovanni Coppola, and Harley I Kornblum. Large-scale assessment of the gliomasphere model system. *Neuro. Oncol.*, 18(10):1367–1378, October 2016. 2686
2687
2688
- Arthur W Lambert, Diwakar R Pattabiraman, and Robert A Weinberg. Emerging biological principles of metastasis. *Cell*, 168(4):670–691, February 2017. 2689
2690
- Samantha B Larsen, Christopher J Cowley, Sairaj M Sajjath, Douglas Barrows, Yihao Yang, Thomas S Carroll, and Elaine Fuchs. Establishment, maintenance, and recall of inflammatory memory. *Cell Stem Cell*, 28(10):1758–1774.e8, October 2021. 2691
2692
2693
- Henry Lee-Six, Nina Friesgaard Øbro, Mairi S Shepherd, Sebastian Grossmann, Kevin Dawson, Miriam Belmonte, Robert J Osborne, Brian JP Huntly, Inigo Martincorena, Elizabeth Anderson, et al. Population dynamics of normal human blood inferred from somatic mutations. *Nature*, 561(7724):473–478, 2018. 2694
2695
2696

- P O Lewis. A likelihood approach to estimating phylogeny from discrete morphological character data. *Syst. Biol.*, 50(6): 913–925, 2001. 2697
2698
- Ruoyan Li, Yiqing Du, Zhanghua Chen, Deshu Xu, Tianxin Lin, Shanzhao Jin, Gongwei Wang, Ziyang Liu, Min Lu, Xu Chen, Tao Xu, and Fan Bai. Macroscopic somatic clonal expansion in morphologically normal human urothelium. *Science*, 370(6512):82–89, October 2020. 2699
2700
2701
- Agnese Loda, Samuel Collombet, and Edith Heard. Gene regulation in time and space during x-chromosome inactivation. *Nat. Rev. Mol. Cell Biol.*, 23(4):231–249, April 2022. 2702
2703
- Michael A Lodato, Mollie B Woodworth, Semin Lee, Gilad D Evrony, Bhaven K Mehta, Amir Karger, Soohyun Lee, Thomas W Chittenden, Alissa M D’Gama, Xuyu Cai, Lovelace J Luquette, Eunjung Lee, Peter J Park, and Christopher A Walsh. Somatic mutation in single human neurons tracks developmental and transcriptional history. *Science*, 350(6256):94–98, October 2015. 2704
2705
2706
2707
- Stilianos Louca. Simulating trees with millions of species. *Bioinformatics*, 36(9):2907–2908, May 2020. 2708
- Stilianos Louca and Michael Doebeli. Efficient comparative phylogenetics on large trees. *Bioinformatics*, 34(6):1053–1055, March 2018. 2709
2710
- Stilianos Louca and Matthew W Pennell. A general and efficient algorithm for the likelihood of diversification and Discrete-Trait evolutionary models. *Syst. Biol.*, 69(3):545–556, September 2019. 2711
2712
- Leif S Ludwig, Caleb A Lareau, Jacob C Ulirsch, Elena Christian, Christoph Muus, Lauren H Li, Karin Pelka, Will Ge, Yaara Oren, Alison Brack, Travis Law, Christopher Rodman, Jonathan H Chen, Genevieve M Boland, Nir Hacohen, Orit Rozenblatt-Rosen, Martin J Aryee, Jason D Buenrostro, Aviv Regev, and Vijay G Sankaran. Lineage tracing in humans enabled by mitochondrial mutations and Single-Cell genomics. *Cell*, 176(6):1325–1339.e22, March 2019. 2713
2714
2715
2716
- Lucie Marhounová, Alexander Kotrschal, Kristina Kverková, Niclas Kolm, and Pavel Němec. Artificial selection on brain size leads to matching changes in overall number of neurons. *Evolution*, 73(9):2003–2012, September 2019. 2717
2718
- Iñigo Martincorena, Amit Roshan, Moritz Gerstung, Peter Ellis, Peter Van Loo, Stuart McLaren, David C Wedge, Anthony Fullam, Ludmil B Alexandrov, Jose M Tubio, Lucy Stebbings, Andrew Menzies, Sara Widaa, Michael R Stratton, Philip H Jones, and Peter J Campbell. Tumor evolution. high burden and pervasive positive selection of somatic mutations in normal human skin. *Science*, 348(6237):880–886, May 2015. 2719
2720
2721
2722
- Iñigo Martincorena, Joanna C Fowler, Agnieszka Wabik, Andrew R J Lawson, Federico Abascal, Michael W J Hall, Alex Cagan, Kasumi Murai, Krishnaa Mahbubani, Michael R Stratton, Rebecca C Fitzgerald, Penny A Handford, Peter J Campbell, Kouros Saeb-Parsy, and Philip H Jones. Somatic mutant clones colonize the human esophagus with age. *Science*, 362(6417):911–917, November 2018. 2723
2724
2725
2726
- Shannon L Maude, Noelle Frey, Pamela A Shaw, Richard Aplenc, David M Barrett, Nancy J Bunin, Anne Chew, Vanessa E Gonzalez, Zhaohui Zheng, Simon F Lacey, et al. Chimeric antigen receptor t cells for sustained remissions in leukemia. *New England Journal of Medicine*, 371(16):1507–1517, 2014. 2727
2728
2729
- Tali Mazor, Aleksandr Pankov, Jun S Song, and Joseph F Costello. Intratumoral heterogeneity of the epigenome, 2016. 2730
- José L McFaline-Figueroa, Andrew J Hill, Xiaojie Qiu, Dana Jackson, Jay Shendure, and Cole Trapnell. A pooled single-cell genetic screen identifies regulatory checkpoints in the continuum of the epithelial-to-mesenchymal transition. *Nat. Genet.*, 51(9):1389–1398, September 2019. 2731
2732
2733
- Linde A Miles, Robert L Bowman, Tiffany R Merlinsky, Isabelle S Csete, Aik T Ooi, Robert Durruthy-Durruthy, Michael Bowman, Christopher Famulare, Minal A Patel, Pedro Mendez, et al. Single-cell mutation analysis of clonal evolution in myeloid malignancies. *Nature*, 587(7834):477–482, 2020. 2734
2735
2736
- Anna Minkina, Junyue Cao, and Jay Shendure. Tethering distinct molecular profiles of single cells by their lineage histories to investigate sources of cell state heterogeneity. May 2022. 2737
2738
- P A P Moran. Notes on continuous stochastic phenomena. *Biometrika*, 37(1-2):17–23, June 1950. 2739
- Tamara Münkemüller, Sébastien Lavergne, Bruno Bzeznik, Stéphane Dray, Thibaut Jombart, Katja Schiffrers, and Wilfried Thuiller. How to measure and test phylogenetic signal, 2012. 2740
2741

- Anna S Nam, Kyu-Tae Kim, Ronan Chaligne, Franco Izzo, Chelston Ang, Justin Taylor, Robert M Myers, Ghaith Abu-Zeinah, Ryan Brand, Nathaniel D Omans, Alicia Alonso, Caroline Sheridan, Marisa Mariani, Xiaoguang Dai, Eoghan Harrington, Alessandro Pastore, Juan R Cubillos-Ruiz, Wayne Tam, Ronald Hoffman, Raul Rabadan, Joseph M Scandura, Omar Abdel-Wahab, Peter Smibert, and Dan A Landau. Somatic mutations and cell identity linked by genotyping of transcriptomes. *Nature*, 571(7765):355–360, July 2019. 2742–2746
- 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.*, 22(1):3–18, January 2021. 2747–2748
- Atsushi Natsume, Motokazu Ito, Keisuke Katsushima, Fumiharu Ohka, Akira Hatanaka, Keiko Shinjo, Shinya Sato, Satoru Takahashi, Yuta Ishikawa, Ichiro Takeuchi, Hiroki Shimogawa, Motonari Uesugi, Hideyuki Okano, Seung U Kim, Toshihiko Wakabayashi, Jean-Pierre J Issa, Yoshitaka Sekido, and Yutaka Kondo. Chromatin regulator PRC2 is a key regulator of epigenetic plasticity in glioblastoma. *Cancer Res.*, 73(14):4559–4570, July 2013. 2749–2752
- Nicholas Navin, Jude Kendall, Jennifer Troge, Peter Andrews, Linda Rodgers, Jeanne McIndoo, Kerry Cook, Asya Stepansky, Dan Levy, Diane Esposito, Lakshmi Muthuswamy, Alex Krasnitz, W Richard McCombie, James Hicks, and Michael Wigler. Tumour evolution inferred by single-cell sequencing. *Nature*, 472(7341):90–94, April 2011. 2753–2755
- S Nee, R M May, and P H Harvey. The reconstructed evolutionary process. *Philos. Trans. R. Soc. Lond. B Biol. Sci.*, 344(1309):305–311, May 1994. 2756–2757
- Cyril Neftel, Julie Laffy, Mariella G Filbin, Toshiro Hara, Marni E Shore, Gilbert J Rahme, Alyssa R Richman, Dana Silverbush, Mckenzie L Shaw, Christine M Hebert, John Dewitt, Simon Gritsch, Elizabeth M Perez, L Nicolas Gonzalez Castro, Xiaoyang Lan, Nicholas Druck, Christopher Rodman, Danielle Dionne, Alexander Kaplan, Mia S Bertalan, Julia Small, Kristine Pelton, Sarah Becker, Dennis Bonal, Quang-De Nguyen, Rachel L Servis, Jeremy M Fung, Ravindra Mylvaganam, Lisa Mayr, Johannes Gojo, Christine Haberler, Rene Geyeregger, Thomas Czech, Irene Slavic, Brian V Nahed, William T Curry, Bob S Carter, Hiroaki Wakimoto, Priscilla K Brastianos, Tracy T Batchelor, Anat Stemmer-Rachamimov, Maria Martinez-Lage, Matthew P Frosch, Ivan Stamenkovic, Nicolo Riggi, Esther Rheinbay, Michelle Monje, Orit Rozenblatt-Rosen, Daniel P Cahill, Anoop P Patel, Tony Hunter, Inder M Verma, Keith L Ligon, David N Louis, Aviv Regev, Bradley E Bernstein, Itay Tirosh, and Mario L Suvà. An integrative model of cellular states, plasticity, and genetics for glioblastoma. *Cell*, 178(4):835–849.e21, August 2019. 2758–2767
- James G Nicholson and Howard A Fine. Diffuse glioma heterogeneity and its therapeutic implications. *Cancer Discov.*, 11(3):575–590, March 2021. 2768–2769
- Sonja Nowotschin, Manu Setty, Ying-Yi Kuo, Vincent Liu, Vidur Garg, Roshan Sharma, Claire S Simon, Nestor Saiz, Rui Gardner, Stéphane C Boutet, Deanna M Church, Pamela A Hoodless, Anna-Katerina Hadjantonakis, and Dana Pe'er. The emergent landscape of the mouse gut endoderm at single-cell resolution. *Nature*, 569(7756):361–367, May 2019. 2770–2772
- Yaara Oren, Michael Tsabar, Michael S Cuoco, Liat Amir-Zilberstein, Heidie F Cabanos, Jan-Christian Hütter, Bomiao Hu, Pratiksha I Thakore, Marcin Tabaka, Charles P Fulco, William Colgan, Brandon M Cuevas, Sara A Hurvitz, Dennis J Slamon, Amy Deik, Kerry A Pierce, Clary Clish, Aaron N Hata, Elma Zaganjor, Galit Lahav, Katerina Politi, Joan S Brugge, and Aviv Regev. Cycling cancer persister cells arise from lineages with distinct programs. *Nature*, 596(7873):576–582, August 2021. 2773–2777
- M Pagel. Inferring the historical patterns of biological evolution. *Nature*, 401(6756):877–884, October 1999. 2778
- Mark Pagel. Detecting correlated evolution on phylogenies: a general method for the comparative analysis of discrete characters. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 255(1342):37–45, 1994. 2779–2780
- Efthymia Papalexli and Rahul Satija. Single-cell RNA sequencing to explore immune cell heterogeneity. *Nat. Rev. Immunol.*, 18(1):35–45, January 2018. 2781–2782
- Ievgenia Pastushenko and Cédric Blanpain. EMT transition states during tumor progression and metastasis. *Trends Cell Biol.*, 29(3):212–226, March 2019. 2783–2784
- Ievgenia Pastushenko, Audrey Brisebarre, Alejandro Sifrim, Marco Fioramonti, Tatiana Revenco, Soufiane Boumahdi, Alexandra Van Keymeulen, Daniel Brown, Virginie Moers, Sophie Lemaire, Sarah De Clercq, Esmeralda Minguijón, Cédric Balsat, Youri Sokolow, Christine Dubois, Florian De Cock, Samuel Scozzaro, Federico Sopena, Angel Lanas, Nicky D’Haene, Isabelle Salmon, Jean-Christophe Marine, Thierry Voet, Panagiota A Sotiropoulou, and Cédric Blanpain. Identification of the tumour transition states occurring during EMT. *Nature*, 556(7702):463–468, April 2018. 2785–2789

- Weike Pei, Fuwei Shang, Xi Wang, Ann-Kathrin Fanti, Alessandro Greco, Katrin Busch, Kay Klapproth, Qin Zhang, Claudia Quedenau, Sascha Sauer, Thorsten B Feyerabend, Thomas Höfer, and Hans-Reimer Rodewald. Resolving fates and Single-Cell transcriptomes of hematopoietic stem cell clones by PolyloxExpress barcoding. *Cell Stem Cell*, 27(3):383–395.e8, September 2020. 2790–2793
- Blanca Pijuan-Sala, Jonathan A Griffiths, Carolina Guibentif, Tom W Hiscock, Wajid Jawaid, Fernando J Calero-Nieto, Carla Mulas, Ximena Ibarra-Soria, Richard C V Tyser, Debbie Lee Lian Ho, Wolf Reik, Shankar Srinivas, Benjamin D Simons, Jennifer Nichols, John C Marioni, and Berthold Göttgens. A single-cell molecular map of mouse gastrulation and early organogenesis. *Nature*, 566(7745):490–495, February 2019. 2794–2797
- Lindsey W Plasschaert, Rapolas Žilionis, Rayman Choo-Wing, Virginia Savova, Judith Knehr, Guglielmo Roma, Allon M Klein, and Aron B Jaffe. A single-cell atlas of the airway epithelium reveals the CFTR-rich pulmonary ionocyte, 2018. 2798–2799
- 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, David Roazen, et al. Scaling accurate genetic variant discovery to tens of thousands of samples. *BioRxiv*, page 2011178, 2017. 2800–2802
- Jeffrey J Quinn, Matthew G Jones, Ross A Okimoto, Shigeki Nanjo, Michelle M Chan, Nir Yosef, Trever G Bivona, and Jonathan S Weissman. Single-cell lineages reveal the rates, routes, and drivers of metastasis in cancer xenografts. *Science*, 371(6532), February 2021. 2803–2805
- Bushra Raj, Daniel E Wagner, Aaron McKenna, Shristi Pandey, Allon M Klein, Jay Shendure, James A Gagnon, and Alexander F Schier. Simultaneous single-cell profiling of lineages and cell types in the vertebrate brain. *Nat. Biotechnol.*, 36(5):442–450, June 2018. 2806–2808
- Klaus Rehe, Kerrie Wilson, Simon Bomken, Daniel Williamson, Julie Irving, Monique L den Boer, Martin Stanulla, Martin Schrappe, Andrew G Hall, Olaf Heidenreich, et al. Acute b lymphoblastic leukaemia-propagating cells are present at high frequency in diverse lymphoblast populations. *EMBO molecular medicine*, 5(1):38–51, 2013. 2809–2811
- Liam J Revell. phytools: an R package for phylogenetic comparative biology (and other things), 2012. 2812
- David F Robinson and Leslie R Foulds. Comparison of phylogenetic trees. *Math. Biosci.*, 53(1-2):131–147, February 1981. 2813
- Alejo E Rodriguez-Fraticelli, Caleb Weinreb, Shou-Wen Wang, Rosa P Migueles, Maja Jankovic, Marc Usart, Allon M Klein, Sally Lowell, and Fernando D Camargo. Single-cell lineage tracing unveils a role for TCF15 in haematopoiesis. *Nature*, 583(7817):585–589, July 2020. 2814–2816
- Aljosja Rogiers, Irene Lobon, Lavinia Spain, and Samra Turajlic. The genetic evolution of metastasis. *Cancer Res.*, 82(10):1849–1857, May 2022. 2817–2818
- Michaela Mrugala Rothová, Alexander Valentin Nielsen, Martin Proks, Yan Fung Wong, Alba Redo Riveiro, Madeleine Linneberg-Agerholm, Eyal David, Ido Amit, Ala Trusina, and Joshua Mark Brickman. Identification of the central intermediate in the extra-embryonic to embryonic endoderm transition through single-cell transcriptomics. *Nat. Cell Biol.*, 24(6):833–844, June 2022. 2819–2822
- Sohrab Salehi, Fatemeh Dorri, Kevin Chern, Farhia Kabeer, Nicole Rusk, Tyler Funnell, Marc J Williams, Daniel Lai, Mirela Andronescu, Kieran R Campbell, Andrew McPherson, Samuel Aparicio, Andrew Roth, Sohrab Shah, and Alexandre Bouchard-Côté. Cancer phylogenetic tree inference at scale from 1000s of single cell genomes. November 2022. 2823–2825
- Vijay G Sankaran, Jonathan S Weissman, and Leonard I Zon. Cellular barcoding to decipher clonal dynamics in disease. *Science*, 378(6616):eabm5874, October 2022. 2826–2827
- Bechara Saykali, Navrita Mathiah, Wallis Nahaboo, Marie-Lucie Racu, Latifa Hammou, Matthieu Defrance, and Isabelle Migeotte. Distinct mesoderm migration phenotypes in extra-embryonic and embryonic regions of the early mouse embryo. *Elife*, 8, April 2019. 2828–2830
- Brett A Schroeder, Jennifer Jess, Hari Sankaran, and Nirali N Shah. Clinical trials for chimeric antigen receptor t-cell therapy: lessons learned and future directions. *Current Opinion in Hematology*, 29(4):225–232, 2022. 2831–2832

- Sydney M Shaffer, Margaret C Dunagin, Stefan R Torborg, Eduardo A Torre, Benjamin Emert, Clemens Krepler, Marilda Beqiri, Katrin Sproesser, Patricia A Brafford, Min Xiao, Elliott Eggan, Ioannis N Anastopoulos, Cesar A Vargas-Garcia, Abhyudai Singh, Katherine L Nathanson, Meenhard Herlyn, and Arjun Raj. Rare cell variability and drug-induced reprogramming as a mode of cancer drug resistance. *Nature*, 546(7658):431–435, June 2017. 2833–2836
- Sydney M Shaffer, Benjamin L Emert, Raúl A Reyes Hueros, Christopher Cote, Guillaume Harmange, Dylan L Schaff, Ann E Sizemore, Rohit Gupte, Eduardo Torre, Abhyudai Singh, Danielle S Bassett, and Arjun Raj. Memory sequencing reveals heritable Single-Cell gene expression programs associated with distinct cellular behaviors. *Cell*, 182(4):947–959.e17, August 2020. 2837–2840
- Kamen P Simeonov, China N Byrns, Megan L Clark, Robert J Norgard, Beth Martin, Ben Z Stanger, Jay Shendure, Aaron McKenna, and Christopher J Lengner. Single-cell lineage tracing of metastatic cancer reveals selection of hybrid EMT states. *Cancer Cell*, 39(8):1150–1162.e9, August 2021. 2841–2843
- Bastiaan Spanjaard, Bo Hu, Nina Mitic, Pedro Olivares-Chauvet, Sharan Janjuha, Nikolay Ninov, and Jan Philipp Junker. Simultaneous lineage tracing and cell-type identification using CRISPR-Cas9-induced genetic scars. *Nat. Biotechnol.*, 36(5):469–473, June 2018. 2844–2846
- Tanja Stadler and Mike Steel. Distribution of branch lengths and phylogenetic diversity under homogeneous speciation models. *J. Theor. Biol.*, 297:33–40, March 2012. 2847–2848
- Mike A Steel and David Penny. Distributions of tree comparison Metrics—Some new results. *Syst. Biol.*, 42(2):126–141, June 1993. 2849–2850
- Aravind Subramanian, Pablo Tamayo, Vamsi K Mootha, Sayan Mukherjee, Benjamin L Ebert, Michael A Gillette, Amanda Paulovich, Scott L Pomeroy, Todd R Golub, Eric S Lander, and Jill P Mesirov. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U. S. A.*, 102(43):15545–15550, October 2005. 2851–2854
- Mario-Luca Suvà, Nicolò Riggi, Michalina Janiszewska, Ivan Radovanovic, Paolo Provero, Jean-Christophe Stehle, Karine Baumer, Marie-Aude Le Bitoux, Denis Marino, Luisa Cironi, Victor E Marquez, Virginie Clément, and Ivan Stamenkovic. EZH2 is essential for glioblastoma cancer stem cell maintenance. *Cancer Res.*, 69(24):9211–9218, December 2009. 2855–2858
- N Takagi and M Sasaki. Preferential inactivation of the paternally derived X chromosome in the extraembryonic membranes of the mouse. *Nature*, 256(5519):640–642, August 1975. 2859–2860
- Jean Paul Thiery. Epithelial-mesenchymal transitions in tumour progression. *Nat. Rev. Cancer*, 2(6):442–454, June 2002. 2861
- Shifaan Thowfeequ and Shankar Srinivas. Embryonic and extraembryonic tissues during mammalian development: shifting boundaries in time and space. *Philos. Trans. R. Soc. Lond. B Biol. Sci.*, 377(1865):20210255, December 2022. 2862–2863
- Samra Turajlic, Hang Xu, Kevin Litchfield, Andrew Rowan, Tim Chambers, Jose I Lopez, David Nicol, Tim O’Brien, James Larkin, Stuart Horswell, Mark Stares, Lewis Au, Mariam Jamal-Hanjani, Ben Challacombe, Ashish Chandra, Steve Hazell, Claudia Eichler-Jonsson, Aspasia Soultati, Simon Chowdhury, Sarah Rudman, Joanna Lynch, Archana Fernando, Gordon Stamp, Emma Nye, Faiz Jabbar, Lavinia Spain, Sharanpreet Lall, Rosa Guarch, Mary Falzon, Ian Proctor, Lisa Pickering, Martin Gore, Thomas B K Watkins, Sophia Ward, Aengus Stewart, Renzo DiNatale, Maria F Becerra, Ed Reznik, James J Hsieh, Todd A Richmond, George F Mayhew, Samantha M Hill, Catherine D McNally, Carol Jones, Heidi Rosenbaum, Stacey Stanislaw, Daniel L Burgess, Nelson R Alexander, Charles Swanton, PEACE, and TRACERx Renal Consortium. Tracking cancer evolution reveals constrained routes to metastases: TRACERx renal. *Cell*, 173(3):581–594.e12, April 2018. 2864–2872
- Samra Turajlic, Andrea Sottoriva, Trevor Graham, and Charles Swanton. Resolving genetic heterogeneity in cancer, 2019. 2873
- Geraldine A Van der Auwera and Brian D O’Connor. *Genomics in the cloud: using Docker, GATK, and WDL in Terra*. O’Reilly Media, 2020. 2874–2875
- David van Dijk, Roshan Sharma, Juozas Nainys, Kristina Yim, Pooja Kathail, Ambrose J Carr, Cassandra Burdziak, Kevin R Moon, Christine L Chaffer, Diwakar Pattabiraman, Brian Bierie, Linas Mazutis, Guy Wolf, Smita Krishnaswamy, and Dana Pe’er. Recovering gene interactions from Single-Cell data using data diffusion. *Cell*, 174(3):716–729.e27, July 2018. 2876–2879

- Bert Vogelstein, Nickolas Papadopoulos, Victor E Velculescu, Shibin Zhou, Luis A Diaz, Jr, and Kenneth W Kinzler. 2880
Cancer genome landscapes. *Science*, 339(6127):1546–1558, March 2013. 2881
- Hiroaki Wakimoto, Gayatry Mohapatra, Ryuichi Kanai, William T Curry, Stephen Yip, Mai Nitta, Anoop P Patel, 2882
Zachary R Barnard, Anat O Stemmer-Rachamimov, David N Louis, Robert L Martuza, and Samuel D Rabkin. Main- 2883
tenance of primary tumor phenotype and genotype in glioblastoma stem cells. *Neuro. Oncol.*, 14(2):132–144, November 2884
2011. 2885
- Fang Wang, Qihan Wang, Vakul Mohanty, Shaoheng Liang, Jinzhuang Dou, Jincheng Han, Darlan Conterno Minussi, 2886
Ruli Gao, Li Ding, Nicholas Navin, and Ken Chen. MEDALT: single-cell copy number lineage tracing enabling gene 2887
discovery. *Genome Biol.*, 22(1):70, February 2021. 2888
- Shou-Wen Wang, Michael J Herges, Kilian Hurley, Darrell N Kotton, and Allon M Klein. Cospar identifies early cell 2889
fate biases from single-cell transcriptomic and lineage information. *Nature Biotechnology*, 40(7):1066–1074, 2022. 2890
- Daniel Wartenberg. Multivariate spatial correlation: A method for exploratory geographical analysis, 1985. 2891
- Robert S Welner, Rosana Pelayo, and Paul W Kincade. Evolving views on the genealogy of b cells. *Nature Reviews 2892
Immunology*, 8(2):95–106, 2008. 2893
- Warren A Whyte, David A Orlando, Denes Hnisz, Brian J Abraham, Charles Y Lin, Michael H Kagey, Peter B Rahl, 2894
Tong Ihn Lee, and Richard A Young. Master transcription factors and mediator establish super-enhancers at key cell 2895
identity genes. *Cell*, 153(2):307–319, April 2013. 2896
- F Alexander Wolf, Philipp Angerer, and Fabian J Theis. SCANPY: large-scale single-cell gene expression data analysis. 2897
Genome Biol., 19(1):15, February 2018. 2898
- Fengying Wu, Jue Fan, Yayi He, Anwen Xiong, Jia Yu, Yixin Li, Yan Zhang, Wencheng Zhao, Fei Zhou, Wei Li, Jie 2899
Zhang, Xiaosheng Zhang, Meng Qiao, Guanghui Gao, Shanbao Chen, Xiaoxia Chen, Xuefei Li, Likun Hou, Chunyan 2900
Wu, Chunxia Su, Shengxiang Ren, Margarete Odenthal, Reinhard Buettner, Nan Fang, and Caicun Zhou. Single-cell 2901
profiling of tumor heterogeneity and the microenvironment in advanced non-small cell lung cancer. *Nat. Commun.*, 12 2902
(1):2540, May 2021. 2903
- Xinjie Xu, Qihang Sun, Xiaoqian Liang, Zitong Chen, Xiaoli Zhang, Xuan Zhou, Meifang Li, Huilin Tu, YU Liu, Sanfang 2904
Tu, et al. Mechanisms of relapse after cd19 car t-cell therapy for acute lymphoblastic leukemia and its prevention and 2905
treatment strategies. *Frontiers in immunology*, 10:2664, 2019. 2906
- Dian Yang, Matthew G Jones, Santiago Naranjo, William M Rideout, 3rd, Kyung Hoi Joseph Min, Raymond Ho, Wei 2907
Wu, Joseph M Replogle, Jennifer L Page, Jeffrey J Quinn, Felix Horns, Xiaojie Qiu, Michael Z Chen, William A Freed- 2908
Pastor, Christopher S McGinnis, David M Patterson, Zev J Gartner, Eric D Chow, Trevor G Bivona, Michelle M Chan, 2909
Nir Yosef, Tyler Jacks, and Jonathan S Weissman. Lineage tracing reveals the phylodynamics, plasticity, and paths of 2910
tumor evolution. *Cell*, 185(11):1905–1923.e25, May 2022. 2911
- Z Yang and S Kumar. Approximate methods for estimating the pattern of nucleotide substitution and the variation of 2912
substitution rates among sites. *Mol. Biol. Evol.*, 13(5):650–659, May 1996. 2913
- Nevin Yusufova, Andreas Kloetgen, Matt Teater, Adewola Osunsade, Jeannie M Camarillo, Christopher R Chin, Ashley S 2914
Doane, Bryan J Venters, Stephanie Portillo-Ledesma, Joseph Conway, et al. Histone h1 loss drives lymphoma by 2915
disrupting 3d chromatin architecture. *Nature*, 589(7841):299–305, 2021. 2916
- Andy GX Zeng, Suraj Bansal, Liqing Jin, Amanda Mitchell, Weihsu Claire Chen, Hussein A Abbas, Michelle Chan- 2917
Seng-Yue, Veronique Voisin, Peter van Galen, Anne Tierens, et al. A cellular hierarchy framework for understanding 2918
heterogeneity and predicting drug response in acute myeloid leukemia. *Nature medicine*, 28(6):1212–1223, 2022. 2919
- Hongkui Zeng. What is a cell type and how to define it? *Cell*, 185(15):2739–2755, July 2022. 2920