

## Clonal dynamics and somatic evolution of haematopoiesis in mouse

Chiraag D. Kapadia<sup>1,2</sup>, Nicholas Williams<sup>3</sup>, Kevin J. Dawson<sup>3</sup>, Caroline Watson<sup>9</sup>, Matthew J. Yousefzadeh<sup>6,10</sup>, Duy Le<sup>2,7</sup>, Kudzai Nyamondo<sup>3,4</sup>, Alex Cagan<sup>3,11</sup>, Sarah Waldvogel<sup>1,2</sup>, Josephine De La Fuente<sup>1,2</sup>, Daniel Leongamornlert<sup>3</sup>, Emily Mitchell<sup>3,4,5</sup>, Marcus A. Florez<sup>2,7</sup>, Rogelio Aguilar<sup>1,2</sup>, Alejandra Martell<sup>1,2</sup>, Anna Guzman<sup>1,2</sup>, David Harrison<sup>8</sup>, Laura J. Niedernhofer<sup>6</sup>, Katherine Y. King<sup>2,7</sup>, Peter J. Campbell<sup>3</sup>, Jamie Blundell<sup>9</sup>, Margaret A. Goodell<sup>\*1,2</sup>, Jyoti Nangalia<sup>\*3,4,5</sup>

1. Department of Molecular and Cellular Biology, Baylor College of Medicine, Houston, TX, USA.
2. Stem Cells and Regenerative Medicine Center, Baylor College of Medicine, Houston, TX, USA.
3. Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK
4. Wellcome-MRC Cambridge Stem Cell Institute, Jeffrey Cheah Biomedical Centre, Cambridge, UK.
5. Department of Haematology, University of Cambridge, Cambridge, UK.
6. Institute on the Biology of Aging and Metabolism, Department of Biochemistry, Molecular Biology and Biophysics, University of Minnesota, Minneapolis, MN, USA.
7. Department of Pediatrics, Division of Infectious Diseases, Baylor College of Medicine, Houston, TX, USA.
8. The Jackson Laboratory, Bar Harbor, ME, USA.
9. Early Cancer Institute, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK.
10. Columbia Center for Translational Immunology, Columbia Center for Human Longevity, Department of Medicine, Columbia University Medical Center, New York, NY, USA.
11. Departments of Genetics, Pathology & Veterinary Medicine, University of Cambridge, Cambridge, UK.

\* **Correspondence:** Margaret A. Goodell, [goodell@bcm.edu](mailto:goodell@bcm.edu) & Jyoti Nangalia, [jn5@sanger.ac.uk](mailto:jn5@sanger.ac.uk)

Summary word count: 208

Text word count: 6582

6 Figures, 9 Extended Data Figures, 4 Supplementary Notes

1 **Abstract**

2 Haematopoietic stem cells maintain blood production throughout life. While extensively  
3 characterised using the laboratory mouse, little is known about how the population is sustained and  
4 evolves with age. We isolated stem cells and progenitors from young and old mice, identifying  
5 221,890 somatic mutations genome-wide in 1845 single cell-derived colonies, and used  
6 phylogenetic analysis to infer the ontogeny and population dynamics of the stem cell pool. Mouse  
7 stem cells and progenitors accrue ~45 somatic mutations per year, a rate only about 2-fold greater  
8 than human progenitors despite the vastly different organismal sizes and lifespans. Phylogenetic  
9 patterns reveal that stem and multipotent progenitor cell pools are both established during  
10 embryogenesis, after which they independently self-renew in parallel over life. The stem cell pool  
11 grows steadily over the mouse lifespan to approximately 70,000 cells, self-renewing about every  
12 six weeks. Aged mice did not display the profound loss of stem cell clonal diversity characteristic  
13 of human haematopoietic ageing. However, targeted sequencing revealed small, expanded clones  
14 in the context of murine ageing, which were larger and more numerous following haematological  
15 perturbations and exhibited a selection landscape similar to humans. Our data illustrate both  
16 conserved features of population dynamics of blood and distinct patterns of age-associated somatic  
17 evolution in the short-lived mouse.

## 18 Introduction

19 The haematopoietic system sustains mammalian life through the continuous generation of  
20 oxygenating red blood cells, an array of immune cells, and platelets that course through all tissues.  
21 In a large animal such as the human, blood production accounts for 86% of daily cellular turnover,  
22 generating ~280 billion cells per day<sup>1</sup>. This process relies on a hierarchy of progenitors that  
23 successively amplify cellular output towards fully differentiated blood cells. All are believed to  
24 ultimately derive from rare haematopoietic stem cells (HSCs), a heterogeneous pool<sup>2-4</sup> maintained  
25 in a relatively protected state to support blood production throughout life.

26 HSCs are the best-studied and utilised of somatic stem cells. They are the basis for life-saving  
27 bone marrow transplantation and have been purified from humans and mice for studies on their  
28 molecular regulation. HSCs are capable of being activated by various stimuli, such as infection and  
29 bleeding<sup>5-8</sup>, in order to rapidly replenish differentiated blood cells as needed, and concomitantly  
30 undergo controlled self-renewal to sustain the stem cell pool over time.

31 Like all somatic cells, HSCs accumulate somatic mutations with age<sup>9-13</sup>. In humans, some  
32 mutations promote cellular fitness, driving clonal outgrowth during normal ageing<sup>9</sup>. Such 'clonal  
33 haematopoiesis' (CH), while remaining at very low levels in younger individuals, is ubiquitous in the  
34 elderly where it results in a dramatic loss of clonal diversity<sup>9</sup>. CH is a known risk factor for blood  
35 cancers and age-associated non-cancerous disease, and may encode other ageing  
36 phenotypes<sup>9,14-17</sup>. Extensive clonal expansions have also been described across human tissues  
37 where they are associated with ageing, cancer and other diseases, reflecting the consequences of  
38 lifelong somatic evolution<sup>18-22</sup>. Whether these patterns of somatic evolution are also features of  
39 ageing in other species is unknown.

40 Within *Mammalia*, the rate of somatic mutation accrual in colonic epithelium inversely scales with  
41 lifespan; that is, species acquire a similar magnitude of mutations by the end-of-life independent of  
42 lifespan<sup>23</sup>. However, it is unclear if this pattern extends to other tissues beyond the colon and  
43 whether the consequences of somatic evolution over human life also scale to shorter-lived species.

44 The inbred laboratory mouse is used ubiquitously across biomedical research. It has been used  
45 extensively to study haematopoiesis, leading to fundamental tenets of somatic stem cell biology.  
46 The most commonly used strain, C57Bl/6J, has a median lifespan of 28 months<sup>24</sup>, 1/35th that of  
47 humans, and broadly recapitulates many phenotypic features of human ageing, with some

48 preliminary data suggesting a lower rate of CH<sup>25</sup>. Here, we study the ontogeny, clonal dynamics,  
49 and selection landscapes of murine HSC populations *in vivo* to understand the evolutionary  
50 processes shaping the maintenance and ageing of blood production.

## 51 **RESULTS**

### 52 **Whole genome sequencing of hematopoietic stem and progenitor cells**

53 To study somatic mutagenesis and clonal dynamics in the haematopoietic compartment within the  
54 laboratory mouse, we purified HSCs from three young (3-months) and three aged (30-months)  
55 healthy C57BL/6J female mice (Fig.1A, Extended Data Fig.1A), ages estimated to be human  
56 lifespan equivalents of ~20 and 85-90 years respectively (Supplementary Note 1). A longstanding  
57 consensus is that haematopoiesis is supported by long-term stem cells (LT-HSCs, henceforth,  
58 referred to as HSCs) which give rise to multipotent progenitors (MPPs, sometimes called short-  
59 term HSCs). Extensive functional analysis established that both HSCs and MPPs, distinguished on  
60 the basis of their cell-surface markers, can support haematopoiesis and produce all differentiated  
61 blood cell types, but HSCs can engraft hosts following multiple rounds of serial transplantation,  
62 whereas MPPs cannot <sup>26-31</sup>. Therefore, we examined both these populations *in vivo*. Single HSCs  
63 and MPPs were expanded *in vitro* to produce colonies (Fig.1A, Extended Data Fig.1B) for whole-  
64 genome DNA sequencing at an average depth of 14X. From individual 3-month and 30-month  
65 animals (n=6), we sequenced 61-235 HSC-derived colonies and 70-191 MPP-derived colonies  
66 (Fig.1B). We also purified fewer HSCs from 17 additional mice aged 1 day to 30 months (total 242,  
67 ranging from 9-24 cells per animal). Following exclusion of 139 colonies due to low sequencing  
68 coverage or lack of clonality (Extended Data Fig.1C and Methods), 1547 whole genomes (908  
69 HSCs, 639 MPPs) were taken forward for somatic mutation identification and phylogenetic  
70 reconstruction.

### 71 **Somatic mutation accumulation in murine haematopoietic stem cells**

72 Comparison of HSCs from young and old mice revealed a constant rate of somatic mutation  
73 accumulation with age (Fig.1C). Mice aged 3 months had an estimated 59.5 single base  
74 substitutions (SBS) (95% confidence interval, CI, 57.3-61.7), and by 30 months had acquired 161.4  
75 SBS per HSC (CI 155.1-167.8), corresponding to 45.3 SBS per year (CI 42.2-48.4) or a somatic  
76 mutation being acquired every 8-9 days within HSCs. Across the diploid mouse genome, this  
77 reflects a mutation rate of  $8.3 \cdot 10^{-9}$  bp per year (CI  $7.7-8.9 \cdot 10^{-9}$  bp/year). Few insertions-deletions  
78 were captured per colony (median 1, range 0-4) with no chromosomal changes observed. Previous

79 studies suggest that MPPs are a more rapidly cycling population<sup>5,32,33</sup> thought to amplify cell  
80 production from HSCs, which could result in a greater mutation burden. However, there was no  
81 difference in mutation burden between HSCs and MPPs (Fig.1D), consistent with observations  
82 from human blood wherein no appreciable differences in somatic mutation burdens between HSCs  
83 and more differentiated blood cells are apparent<sup>11,12</sup>.

84 The murine HSC SBS rate is about twice that of humans (14-17 SBS per year)<sup>9-11,34</sup>, given their  
85 similar genome sizes. This is consistent with the concept that somatic mutation rates are negatively  
86 correlated with lifespan<sup>23</sup> such that short-lived species have higher rates of somatic mutation  
87 accumulation than longer-lived animals. However, the ~10-fold difference in ultimate mutation  
88 burden (~150 in HSCs from 30-month-old mice vs >1,500 in human HSCs of an equivalent age of  
89 85-90 years) is much greater than expected given that total end-of-life somatic mutation numbers  
90 in mammalian intestinal crypts show low variation regardless of life-span<sup>23</sup>. Thus, we wished to  
91 validate the lower-than-expected somatic mutation burden observed in aged murine stem cells.

92 First, we compared genome-wide mutation burdens in HSCs with that of matched intestinal crypts  
93 from the same three aged mice. Following WGS of individually microdissected clonal crypts (n=16,  
94 range 5-6 per sample, Extended Data Fig.1D, Methods), we confirmed that colonic epithelium  
95 exhibited substantially higher mutation burdens, similar to that reported previously<sup>23</sup> (Fig.1E),  
96 confirming that we were not underestimating mutations in HSCs. Secondly, we undertook  
97 independent nano-error rate whole-genome duplex sequencing<sup>12</sup> of matched whole blood from the  
98 three aged animals. This method identifies mutations in single DNA molecules and, thus, can  
99 orthogonally estimate mutation burden from peripheral blood. The mutation burden was not  
100 statistically different from that of haematopoietic colonies (Fig.1E). We did note a non-significant  
101 trend towards higher mutation burden estimates from whole blood than HSC colonies – this is likely  
102 due to whole blood including lymphoid cells which have higher mutation burdens<sup>35</sup>. Despite whole  
103 blood having a mixture of mature cell types and the different sequencing technologies used, these  
104 data confirm that somatic mutation rates in blood do not inversely scale with lifespan to the same  
105 degree as observed in colon.

## 106 **Aetiology of mutational processes in haematopoietic stem cells in mouse**

107 The pattern of sharing of somatic mutations across individual colonies can be used to reconstruct  
108 a phylogenetic tree that depicts their ancestral lineage relationships (Methods). We use the term  
109 'lineage' here to represent the direct line of descent rather than different blood cell types. Figure 2

110 shows the phylogenetic trees for a 3-month- and 30-month-old mouse, with additional young and  
111 aged phylogenies in Extended Data Figure 2. At the tips of the trees are individual HSC- (blue) and  
112 MPP-derived colonies (red); the branches that trace upwards from each colony to the root of the  
113 tree reflect the somatic mutations present in that individual colony and how these mutations are  
114 shared across other colonies. Individual branchpoints (“coalescences”) represent ancestral cell  
115 divisions wherein descendants of both daughter cells were captured at sampling. Colonies that  
116 share a common ancestor on the phylogeny are referred to as a clade.

117 First, we wished to understand the aetiology of the higher rate of mutation accumulation in murine  
118 HSC and MPPs compared to human HSCs. DNA replication during cell division is one source of  
119 mutations reflecting DNA polymerase base incorporation errors. To estimate the rate of DNA  
120 replication-associated somatic mutation accumulation in mice, we studied the distribution of nodes  
121 on the phylogenies with more than two descendant lineages, termed polytomies.

122 These are evidence of ancestral cell divisions which were not associated with the acquisition of a  
123 somatic mutation and can be used to infer the average number of mutations that are acquired per  
124 cell division<sup>10</sup> (Extended Data Fig.3). We focused on the roots of the trees where we capture the  
125 greatest number of coalescences, both due to the small population size and the rapid divisions at  
126 this point in life. We observed 266 lineages by 12 mutations of molecular time in five donors that  
127 had adequate (>10 lineages) diversity. Of the 265 symmetrical self-renewing cell divisions that  
128 would have required, 44 were mutationally silent, leading to a mutation rate estimate of 1.80 (95%  
129 CI: 1.46-2.19) mutations per cell division during early life (Extended Data Fig.3). This estimate is  
130 not significantly different from that previously observed in humans (1.84 mutations/cell division,  
131  $p=0.5$ )<sup>13</sup>, suggesting that excess mutation accumulation is not occurring due to poorer fidelity during  
132 DNA replication in murine stem cells.

133 Mutagenic biological processes yield distinguishable patterns of base substitutions at trinucleotide  
134 sequence contexts, termed mutation signatures. We identified three mutational processes  
135 (Extended Data Fig.4A, Methods): i) SBS1 reflecting the spontaneous deamination of methylated  
136 cytosines, ii) SBS5 likely produced by cell-intrinsic damage and repair processes, and iii) SBS18  
137 characterised by C>A transversions potentially linked to oxidative damage. Substitutions attributed  
138 to SBS1 and SBS5 increased with age (8.64 SBS1/month and 32.52 SBS5/month), keeping with  
139 their clock-like nature across species; indeed, these processes account for most mutations in  
140 healthy human HSCs. Mutations attributed to SBS18 (mean 5.3, range 1.5-18 per colony, Extended

141 Data Fig.4B), were previously identified in murine colorectal crypts<sup>23</sup>, but did not appear to  
142 accumulate with age. To explore the timing of SBS18 mutations, we deconvoluted branch-specific  
143 mutations (Extended Data Fig.4C). SBS18 accrued before three months of life, followed by a  
144 plateau (Extended Data Fig.4D), suggesting a specific early-life vulnerability to these mutations,  
145 reminiscent of their presence in human placenta and human foetal HSCs<sup>13,36</sup>. Taken together, the  
146 higher relative somatic mutation accumulation rate in mice is underlaid by context-specific  
147 mutational processes (SBS18) and a higher rate of endogenous DNA damage and reduced repair  
148 (SBS1 and 5).

### 149 **Origin and parallel establishment of HSC and MPP populations**

150 We next sought to examine the lineage relationships between the HSCs and MPPs. Classical  
151 models of the haematopoietic differentiation hierarchy propose that MPPs derive from HSCs<sup>37,38</sup>.  
152 In recent years, a more nuanced and dynamic picture has emerged, with the identification of  
153 additional self-renewing progenitor compartments<sup>2,4</sup>. Using our phylogenetic approach, stem cell  
154 ontogeny can be retraced *in vivo* during unperturbed haematopoiesis. Working up from the  
155 phylogenetic “tip” states of HSC (blue) or MPP (red), we infer the identity of ancestral branches  
156 and coalescences based on the identity of their nearest sibling cell (detailed in Supplemental Note  
157 2). Branches where we are unable to infer the established cell type for one or more lines of descent  
158 are coloured black. We observed clear vertical bands of HSC-only (all “blue”) and MPP-only (all  
159 “red”) ancestral lineages across the trees representing independent clades (Fig.2A,B, Extended  
160 Data Fig.2), with a minority of HSCs (“blue” tips) being sampled from MPP (“red” clades) and vice  
161 versa. The clear separation of MPPs and HSCs suggests that most HSCs are derived from HSC  
162 self-renewing divisions, and most MPPs are derived from MPP self-renewing divisions, with each  
163 population independently self-renewing in parallel throughout life. If HSCs and MPPs were closely  
164 related to one another, as might be the case if MPPs were recently generated from HSCs, then  
165 one would expect the two cell types to be intermixed across the phylogenetic tree, with individual  
166 clades (cells derived from a common ancestor) containing cells of both types. However, we  
167 observed that clades were largely uniform in composition, containing more cells of the same type  
168 than would be expected if the population of HSCs and MPPs were intermixed (Fig.2C). This  
169 phylogenetic separation of HSC and MPPs provides strong evidence that these two populations  
170 independently contribute to blood production in the mouse.

171 The lack of intermixing between HSCs and MPPs on the phylogenetic trees suggests long-term  
172 inheritance of the HSC or MPP ‘state’, presumably encoded epigenetically. Therefore, we next  
173 explored when such sustained MPP and HSC cell state commitments may have occurred during  
174 life. Coalescences near the top of a phylogeny (near the root) reflect cell divisions that occur soon  
175 after conception. Most branches and coalescences here are ‘black’ (Fig.2A,B, Extended Data  
176 Fig.2), as no established HSC or MPP lineages could be inferred at this time. Stable heritable  
177 identity of either HSC or MPP appears established at similar times – by around 25 mutations of  
178 molecular time suggesting that a substantial proportion of the HSC and MPP populations appear  
179 to diverge early in life. The mixed effects regression model of mutation rate suggests that ~50  
180 somatic mutations may be acquired before birth (intercept of mixed effects model, 48.2, CI 45.61-  
181 50.8, Methods). Thus, HSCs and MPPs are likely established in parallel during foetal development.

182 To explore when *in utero* this was occurring, we evaluated somatic mutations present in both HSCs  
183 and colonic crypts from the same animals – by definition, mutations shared between these tissues  
184 arose in a common ancestor whose progeny contributed to both blood and colonic epithelium. As  
185 blood is derived from mesoderm and colonic epithelium is derived from endoderm, any shared  
186 mutations must have occurred in embryonic cells prior to gastrulation. Mutations on the  
187 haematopoietic phylogeny were observed in sampled colonic crypts (n=4-6 crypts per 30-month-  
188 old mouse) down to 9-11 mutations of molecular time (Extended Data Fig.3), with decreasing  
189 representation of mutations further from the phylogeny root, timing these shared mutations to have  
190 occurred during gastrulation. Indeed, branches with an inferred HSC or MPP identity did not share  
191 mutations with the colon, consistent with these lineages being established after germ layer  
192 specification.

193 Given the likely embryonic establishment of distinct HSC and MPP pools, we next considered the  
194 simplest series of cell state changes (eg HSC to MPP, or HSC to MPP, etc) that might be required  
195 to capture the observed cell identities. We first considered the prevailing view, that MPPs are  
196 generated from HSCs, such that HSC fate occurs prior to specification of MPP. We counted the  
197 number of cell identity changes required to reach the sampled cell identity. Surprisingly, the HSC-  
198 to-MPP model was equivalently parsimonious (requiring a number of cell state changes that was  
199 not statistically different) to a model where all cells start as MPP (with HSCs able to arise from  
200 MPPs), an ontogeny not generally considered likely (Fig.2D, also see Supplementary Note 2).  
201 Overall, our data suggest that many long-term HSC and MPP lineages are established

202 independently and in parallel during early development, and that MPPs do not always arise from  
203 HSCs, contrary to classical haematopoiesis models.

#### 204 **Modelling HSC and MPP establishment and transitions**

205 To formalise the above ideas and develop an ontogeny model for HSCs and MPPs that fits with  
206 our observed data, we developed a hidden Markov tree model. The Markov approach allows  
207 estimation of the rates at which a cell state makes transitions as it evolves through time. We defined  
208 three unobservable ancestral states: embryonic precursor cell (EMB), HSC, and MPP. We then  
209 used the observed outcomes of HSC versus MPP tip states to infer both the sequence and the  
210 transition rates between these states during life (Methods, Supplementary Note 2). We considered  
211 all cells prior to gastrulation (<10 mutations) as EMB, and then assumed that in each unit of  
212 molecular time, there is a fixed probability of transitioning out of this state to either an HSC or MPP  
213 state (with subsequent fixed probabilities of further transitions such as HSC-to-MPP). In addition to  
214 characterising the most feasible model parameters that fit the observed data using a maximum  
215 likelihood approach (Methods), we also estimated the rate of HSC to MPP (and vice-versa)  
216 transitions during life, to account for any HSC/MPP mixed clades (Supplementary Note 2).

217 We fitted the above model to i) each donor, ii) each age group, and iii) the whole cohort. Based on  
218 a nested likelihood ratio test analysis, we found that the model fitted to each age group (young and  
219 old mice separately) was most consistent with our data (Supplementary Note 2). Across the whole  
220 cohort, we found that a model in which EMB can transition to either HSC or MPP was a significantly  
221 better fit than an “HSC-first” model, where all EMB must transition to HSC prior to any MPP  
222 specification. However, when testing young and old mice separately, we were only able to reject  
223 an “HSC-first” ontogeny model in older mice. We could not reject an “HSC-first” model in younger  
224 mice as our data suggested more frequent HSC to MPP transitions earlier in life (Fig.2E). This  
225 apparent inconsistency in the results between young and old mice could perhaps be explained if  
226 the HSCs that produce the MPPs early in life are extinguished by old age, and thus could not be  
227 sampled for inclusion in the phylogeny. Alternatively, the rate of HSCs that transition to MPPs may  
228 be greater earlier in life. Further work is required to explore this. Interestingly, our model indicates  
229 that 50% of all HSC and MPP lineages in young and aged mice had committed to their cell state  
230 before 50 mutations of molecular time, likely before birth. As might be expected, HSC to MPP  
231 transitions were more frequent than MPP to HSC transitions, which were extremely rare (1 in 1000  
232 transitions) and within the plausible limitations of cell-sorting accuracies. (Supplemental Note 2).

## 233 **Haematopoietic population dynamics over life**

234 The pattern of coalescences (branch points) in a phylogenetic tree reflects the ratio ( $N/\lambda$ ) of the  
235 overall population size ( $N$ ) and the HSC self-renewing cell division rate ( $\lambda$ ) over time – both smaller  
236 populations and more frequent cell divisions decrease the interval between coalescences. Mice  
237 haematopoietic phylogenies show a different pattern of coalescences (Fig.2A-B, Extended Data  
238 Fig.2) to equivalently-aged humans<sup>9</sup>. Human haematopoietic phylogenies have a profusion of early  
239 coalescences, reflecting the period of rapid cell division during embryonic growth. Coalescences  
240 are then infrequently observed due to the presence of both a large and stable HSC population by  
241 early adulthood, reappearing only in elderly human phylogenies within clonal expansions when  
242 clonal diversity dramatically collapses.

243 By comparison, murine haematopoietic phylogenies display coalescences continuing down the tree  
244 (see Extended Data Fig.5 for side-by-side human-mouse comparisons). These time intervals  
245 between coalescences can be used to infer the HSC population trajectory ( $N/\lambda$ ) using a  
246 phylodynamic Bayesian framework (Methods). We observed an early period of exponential HSC  
247 growth followed by progressively increasing  $N/\lambda$  over the murine lifespan (Fig.3A), consistent with  
248 the observed increase in total HSC number with age by flow cytometry (Fig.3B), and other  
249 studies<sup>39,40</sup>. Our findings contrast with hematopoietic progenitor population trajectories in humans<sup>9,10</sup>  
250 which exhibit a population growth plateau during adulthood followed by stable population size for  
251 the remainder of life. Interestingly, we infer entirely overlapping  $N/\lambda$  trajectories for HSCs and  
252 MPPs. Together with their similar mutation burdens and lineage independence, these data suggest  
253 that murine HSC and MPP clonal dynamics during steady state *in vivo* haematopoiesis are  
254 indistinguishable.

255 We next developed a joint HSC/MPP population dynamics model (given our data suggests both  
256 populations contribute equivalently to haematopoiesis), in which the population of stem cells grows  
257 towards the target population size, taking into account loss of HSC and MPP cells via cell death or  
258 differentiation (Methods). We then applied approximate Bayesian computation<sup>41</sup>, which generates  
259 simulations of phylogenetic trees to estimate the most likely posterior distributions of population  
260 size and symmetrical self-renewing division rates. Using this approach, we estimate that the murine  
261 HSC-MPP population grows to around 70,000 cells (median 72,414, CI 25,510-98,540). Symmetric  
262 cell divisions occur approximately every 6 weeks (median 6.4 weeks, CI 1.8-13.2 weeks). Stem

263 cells exit the population, by either death or differentiation, about once every 18 weeks (CI 2.3-357  
264 weeks). Posterior density estimates for each mouse are shown in Fig.3 and Extended Data Fig.6.

### 265 **Stem cell contribution to progenitors and mature blood cells**

266 Given the observed lineage independence of HSC and MPP, and their overlapping growth  
267 trajectories, we wondered what difference *in vivo* might exist between the two populations.  
268 Therefore, we studied if HSCs and MPPs might differentially contribute to downstream lineage-  
269 biased progenitors and mature blood cells.

270 We isolated single cells from a mixed progenitor compartment (LSK cells) that includes  
271 granulocyte/macrophage-biased MPPs (MPP<sup>GM</sup>) and lymphoid-biased MPPs (MPP<sup>Ly</sup>) from the  
272 three aged animals. We performed whole genome sequencing on colonies from 298 LSK cells and  
273 performed phylogenetic analysis. Within the extended phylogenetic trees, we observed no  
274 discernable bias in MPP<sup>GM</sup> or MPP<sup>Ly</sup> emerging preferentially from HSC versus MPP (Extended Data  
275 Fig.7A-C). To evaluate this more formally, we separately examined clades that contained MPP<sup>GM</sup>  
276 or MPP<sup>Ly</sup> and evaluated if they were more phylogenetically linked to HSCs or MPPs than expected  
277 by chance. Neither MPP<sup>GM</sup> nor MPP<sup>Ly</sup> preferentially derived from either HSC or MPP beyond  
278 random chance (Extended Data Fig.7D-E), confirming that both HSCs and MPPs produce these  
279 downstream progenitors at seemingly similar proportions. However, these data are limited by a  
280 relatively low number of sampled MPP<sup>GM</sup> and MPP<sup>Ly</sup>.

281 We next evaluated if peripheral blood cells were preferentially derived from HSCs or MPPs. We  
282 performed deep targeted sequencing on peripheral blood DNA from the three aged mice for a  
283 subset of mutations displayed on the corresponding phylogenetic trees (Methods). The fraction of  
284 cells in peripheral blood harbouring a mutation present on the phylogenetic tree can be used to  
285 estimate how much that lineage contributes to blood production. For example, if a single cell or  
286 lineage contributed avidly to differentiated progeny, then its mutations would be seen at high  
287 proportion (variant allele frequency, VAF) in peripheral blood. We recaptured mutations in the  
288 peripheral blood that were acquired in both ancestral HSCs and ancestral MPPs, suggesting that  
289 both these cell types actively contribute to mature blood production. Mutations private to single  
290 cells on the phylogeny were subclonal, occurring below 0.1% VAF in peripheral blood (Extended  
291 Data Fig.8A) in line with each HSC/MPP contributing only a small amount of overall blood  
292 production. While both HSC and MPP ancestral lineages gave rise to peripheral blood, we  
293 observed a slight bias toward increased representation of ancestral MPP lineages compared to

294 HSCs, though this difference was subtle (Extended Data Fig.8B). This subtle difference may be  
295 due to increased proliferation of MPP descendants, or differences in compartment population size  
296 earlier in life; we cannot distinguish between these possibilities.

### 297 **Absence of large clonal expansions in aged mice**

298 A striking feature of the phylogenetic trees in aged mice is the uniform distribution of long branches  
299 with no expanded clades (Fig.3C, Extended Data Fig.5). This indicates mouse haematopoiesis  
300 maintains clonal diversity instead of collapsing into an oligoclonal state as observed in elderly  
301 humans<sup>9</sup>. Indeed, our population dynamics simulations confidently recapitulated observed  
302 phylogenies under a model of neutral growth in the absence of selection. Concordantly, no colonies  
303 (n=1305) displayed mutations in murine orthologues of genes associated with human clonal  
304 haematopoiesis (CH), which could act as potential driver events, aligning with a topology devoid of  
305 observable late-life clonal exponential growth. Among 49,849 SNVs observed across young and  
306 aged samples, the relative rate of nonsynonymous mutation acquisition also did not significantly  
307 depart from neutrality (Fig.3D), with no novel genes identified as being under selection. This  
308 indicates that positive selection does not explain the catalogue of somatic mutations observed.

309 Given that mutation entry, which furnishes a population with phenotypic variation and substrate for  
310 selection, is occurring at a higher rate in mice relative to humans, and in genomes of comparable  
311 size, we considered reasons for the lack of observable clonal expansions (on the phylogenies) and  
312 absence of selection on non-synonymous mutations (using dN/dS), both of which manifest  
313 ubiquitously over time in human haematopoiesis<sup>9</sup>. One possibility here is that there are insufficient  
314 HSC and MPP divisions within the short lifespan of mice to facilitate detectable clonal expansions  
315 of cells with fitness-inferring mutations. Secondly, as both population size and the frequency of self-  
316 renewing cell divisions (captured in  $N/\lambda$ ) determine the rate of random drift, and hence the drift  
317 threshold that selection must overcome<sup>42</sup>, the fitness (s) of newly arising mutations may also be  
318 insufficient for their carrier subclones to exceed the genetic drift threshold within a mouse lifespan  
319 ( $s=\lambda/N$  representing the drift threshold<sup>42</sup>).

320 In the first scenario, clones under selection (i.e., with necessary driver mutations) will still be  
321 present, but would just be too small to detect using a phylogenetic approach that only readily  
322 identifies larger clones (>5% clonal fraction). In the second scenario, the fitness landscape of any  
323 detectable clones would reflect the specific murine haematopoietic drift threshold. Therefore, we

324 sought to address both questions to better understand the evolutionary processes shaping somatic  
325 evolution in blood.

### 326 **Positive selection during homeostatic and perturbed murine haematopoiesis**

327 To examine murine blood for very small expanded haematopoietic clones and the presence of  
328 clonal haematopoiesis (CH), we employed targeted duplex consensus sequencing (Methods). We  
329 reasoned that clonal expansions in mice could be driven by mutations in orthologues of at least  
330 some of the same genes that drive human CH due to their evolutionary conserved biological  
331 functions. We designed a target panel covering murine orthologues of 24 genes associated with  
332 human CH (61.8 kb panel, Methods) and tested whole blood from mice aged 3 to 37 months for  
333 mutations. Median duplex consensus coverage per sample ranged from 28,000–41,000X, allowing  
334 detection of variants present at the magnitude of 1 in 10,000 cells (Methods, Supplementary Note  
335 3).

336 We observed expanded CH clones that increased in prevalence with age (Fig.4A,B). Samples from  
337 young mice (3 months) displayed infrequent or absent (range 0-1) clones, while those from the  
338 oldest mice (37 months) displayed on average 3.5 clones (range 1-6) per animal across these  
339 targeted genes. Average clone size was very small at 0.017% of nucleated blood cells (range  
340 0.0036-0.27%) – representing clonal fractions between 1 in 500 to 1 in 30,000 cells. Clonal  
341 expansions were recurrently driven by mutations in *Dnmt3a* and *Tet2*, genes frequently mutated in  
342 human CH<sup>43</sup>, but also *Bcor* and *Bcorl1*, observed in humans following bone marrow immune  
343 insult<sup>44</sup>. These data are consistent with a previous report identifying rare expanded clones in mice  
344 following transplant<sup>25</sup>.

345 Increased clonal prevalence with age was observed across different laboratory strains, including  
346 the genetically heterogeneous HET3 strain, and at similar clonal fractions (Fig.4C,D), confirming  
347 that small clones driven by known CH drivers are not specific to the C57BL/6J strain. Clones were  
348 present in biological replicates and persistent in mice sampled longitudinally over four months,  
349 though individual clonal dynamics varied (Fig.4E). Variants displayed enrichment for  
350 nonsynonymous mutations across these genes (dN/dS 2.00, CI 1.01-4.02), with per-gene positive  
351 selection evident for *Dnmt3a*, *Bcor*, and *Bcorl1* (dN/dS>1, q<0.1) (Fig.4F). These data confirm that  
352 these small clonal expansions in murine blood are being shaped by positive selection and are not  
353 the result of genetic drift.

354 Laboratory mice are maintained in exceptionally clean conditions with a controlled diet and  
355 environment, in contrast to the regular microbial exposures and systemic insults experienced by  
356 humans. We considered whether similar exposures, which may accelerate CH in humans<sup>45,46</sup>, could  
357 enhance selection and clonal expansion in mice. In humans, mutant *TP53* and *PPM1D* clones are  
358 positively selected for in the context of chemotherapy<sup>47-49</sup>, while *BCOR* mutated clones have a  
359 fitness advantage in the bone marrow environment of aplastic anaemia<sup>44</sup>. To examine whether the  
360 murine haematopoietic selective landscape can be similarly altered, we applied a series of  
361 infectious or myeloablative exposures.

362 We first subjected mice to a normalised microbial experience (NME), in which laboratory mice are  
363 infected with common mouse microbes via exposure to fomite (pet store) bedding, resulting in the  
364 transfer of bacterial, viral, and parasitic pathogens<sup>50,51</sup>. Such exposure has been shown to drive  
365 functional maturation of the murine immune system<sup>50</sup>. Aged NME-exposed mice displayed an  
366 increased burden of somatic clones, especially driven by *Trp53* (Fig.5A). As NME exposure  
367 transmits multiple types of pathogens, making it challenging to disentangle specific pathogen  
368 effects, we next performed targeted exposure to *Mycobacterium avium*, which has been shown to  
369 activate HSCs and lead to chronic inflammation<sup>52</sup>. Aged mice chronically infected with *M. avium*  
370 showed an increased frequency of *Dnmt3a*, *Bcor*, *Tet2*, and *Asx1* mutant clones (Fig.5B),  
371 suggesting that clones harbouring these mutations experience a competitive advantage in the  
372 context of infectious exposure. Differences in driver mutation prevalence between NME and *M.*  
373 *avium*-infected mice may reflect infection severity or immune response differences.

374 To observe the impact of myeloablation, aged mice were treated with commonly used  
375 chemotherapeutic agents 5-fluorouracil and cisplatin. When treated with monthly doses of cisplatin,  
376 we observed globally increased somatic clonal burden (Fig.5C,  $p=0.027$ ). Clones driven by *Trp53*,  
377 *Tet2*, and *Asx1* were enriched relative to controls, and gene-level dN/dS analysis indicated that  
378 *Trp53* was under positive selection for nonsynonymous mutations (Fig.5C), analogous to human  
379 observations<sup>47-49</sup>. Similarly, aged mice treated recurrently with the chemotherapeutic agent 5-  
380 fluorouracil displayed clones at magnitudes-greater proportions than age-matched controls  
381 (Fig.5D). Broadly, these data illustrate that haematopoietic mutation accrual and selection are  
382 sufficient to drive native CH in mice, with modulable selection landscapes.

### 383 **Fitness landscape of clonal haematopoiesis in murine haematopoiesis**

384 Having observed an evolutionarily conserved clonal selection landscapes in murine blood, we  
385 wished to understand why observed clone sizes were much smaller (median 0.017%) compared  
386 to human CH at equivalent times during lifespan. Therefore, we estimated the fitness landscape of  
387 these driver mutations in mouse. We evaluated the distribution of observed variant allele fractions  
388 (VAF) from the targeted duplex sequencing, using an established continuous time branching  
389 evolutionary framework for HSC dynamics<sup>53</sup> (Methods). How the observed distribution of VAFs,  
390 predicted by the evolutionary framework, changes with age is then used to infer the underlying  
391 effective population size ( $N/\lambda$ ), mutation rates ( $\mu$ ), and fitness effects (i.e., clonal growth  
392 percentage per year) of non-synonymous mutations. Due to increasing  $N/\lambda$  with mouse age  
393 (Fig.3A), only clones from mice of a similar age (here chosen to be 24-25 months) during steady-  
394 state haematopoiesis were included in the analyses.

395 By analysing the distribution of neutral mutation VAFs (clones at low VAF bearing synonymous or  
396 intronic mutations), we first yielded an independent orthogonal estimate of  $N/\lambda$  (Methods) of  
397 approximately 16,500 HSC-years (CI, 11,122-21,836, Fig.6A). This inference generated from  
398 targeted sequencing is consistent with that generated from whole genome sequencing and  
399 approximate Bayesian computation (ABC) ( $N/\lambda$  7,918 HSC-years, CI 2,277-20,309). Differences in  
400 the estimates for  $N/\lambda$  from ABC versus the branching evolutionary framework<sup>53</sup> are likely influenced  
401 by (i) the ABC method takes into account population growth inferred from the phylogenetic trees,  
402 whereas the branching evolutionary framework assumes a stable population size, and (ii) the  
403 branching evolutionary framework model relies on using the intronic/synonymous mutation rate as  
404 the background for identifying clonal expansions, which may not reflect the genome-wide mutation  
405 levels. Across the 61.8 kb panel the synonymous/intronic mutation rate was estimated at  $1.8 \cdot 10^{-4}$   
406 base pairs per year (CI  $1.2-2.7 \cdot 10^{-4}$ ). We estimate a nonsynonymous mutation rate of  $3.4 \cdot 10^{-4}$  base  
407 pairs per year (CI:  $2.9-3.9 \cdot 10^{-4}$ ), again only considering VAFs below the maximum observed  
408 synonymous/intronic VAF (Methods), as clones larger than this could be under the influence of  
409 positive selection. The total mutation rate within our targeted panel was thus  $5.2 \cdot 10^{-4}$  per year,  
410 which when scaled to total genome size, corresponds to a global mutation rate of  $11.77 \cdot 10^{-9}$  per  
411 base pair per year (CI  $9.28-14.94 \cdot 10^{-9}$ ). Encouragingly, this is similar to the mutation rate directly  
412 observed from whole genome sequencing of single cell-derived colonies of  $8.29 \cdot 10^{-9}$  per base pair  
413 per year (CI  $7.73-8.85 \cdot 10^{-9}$ ). This agreement suggests that even these low-VAF clones detected  
414 from duplex consensus sequencing represent *bona fide* clonal expansions.

415 Having inferred  $N/\lambda$  and the non-synonymous mutation rate, we could estimate the distribution of  
416 fitness effects driven by non-synonymous mutations (Methods). Our analysis suggests that ~7%  
417 (CI 5-21%) have strong fitness effects (50-200% growth per year) (Fig.6B). Considering that we  
418 infer mouse stem cells to be self-renewing roughly every 6 weeks (CI 2.3-12.5 weeks), an *annual*  
419 growth rate of 200% translates to a *per* symmetrical self-renewing division selective advantage of  
420 ~15% (5-30%), in line with reported selection coefficients of mutated genes associated with CH in  
421 humans<sup>9,34,53</sup>. Indeed, in the short-lived mouse, variants with weaker fitness (<50%) might have  
422 insufficient time to enter exponential, deterministic growth within the population, given that clones  
423 are not established until  $t_{years} > \frac{1}{s}$  (ref. <sup>54</sup>), although any background growth in population size  
424 could circumvent this, allowing for weaker variants to fix in the population. This may also explain  
425 why some of the low-VAF clones identified by duplex consensus sequencing did not increase in  
426 clone size over time (Fig.4E).

## 427 **DISCUSSION**

428 Here, we study the ontogeny, population dynamics and somatic evolution of haematopoietic stem  
429 cells in the most widely used mammalian model organism, the laboratory mouse. Classical models  
430 of blood production depict HSCs at the very top of the haematopoietic differentiation hierarchy,  
431 beneath which all blood cell types emanate. Recent studies suggest additional heterogeneity at the  
432 top of this haematopoietic hierarchy and nuanced self-renewing dynamics<sup>2-4</sup>. Our phylogenetic  
433 data suggest that MPPs (distinguished by their lack of expression of the CD150 marker<sup>55</sup>) do not  
434 always arise from HSCs, and that both populations are established during embryogenesis,  
435 following which they independently self-renew throughout murine life. That MPPs are noticeably  
436 generated from HSCs in a transplant setting may underscore the difference between their *potential*  
437 in an experimental setting and their steady-state *in vivo* function. Moreover, lymphoid and myeloid  
438 progenitors appear to equally derive from HSC and MPP lineages. Recapture of shared variants  
439 indicate both MPPs and HSCs contribute to differentiated peripheral blood production, with a slight  
440 bias toward production from ancestral MPP lineages. These data are aligned with lineage tracing  
441 that showed both populations are capable of making all cell types during normal life<sup>2,32,56</sup>, and with  
442 recent reports of MPPs derived from the embryo (eMPPs) that contribute to lifelong  
443 haematopoiesis<sup>2,57</sup>.

444 We show that HSCs and MPPs grow in lockstep over life, with indistinguishable clonal dynamics  
445 and proliferation rates, to reach a combined population of 25,000-100,000 cells, remarkably close

446 to estimates of the human HSC pool size (20,000-200,000 stem cells)<sup>9,10</sup> and reminiscent of  
447 suggestions of conservation of stem cell numbers across mammalian species<sup>58</sup>. Considering the  
448 log-fold difference in body mass and consequent demands on blood production, this similarity may  
449 be surprising. In both organisms, but especially the mouse, the number of stem cells far exceeds  
450 the apparent lifetime need; the stem cell compartment of a single mouse can be used to fully  
451 reconstitute the blood of ~50 transplant recipients<sup>59</sup>. Perhaps a large stem cell pool confers an  
452 evolutionary advantage in the face of naturally occurring exposures to environmental pathogens<sup>60</sup>  
453 and tissue injury, through both increased tolerance of stem cell losses and improved adaptation  
454 afforded by somatically acquired genomic and epigenetic diversity.

455 Somatic mutation rates have recently been shown to scale inversely with mammalian lifespan. In  
456 colonic epithelium, mice accumulate mutations 20 times faster than humans, aligned with the  
457 difference between their lifespans<sup>23</sup>. This observation raises the intriguing possibility that somatic  
458 mutation rates are visible to selection through their effects on ageing and lifespan. Our data show  
459 this pattern does not extend to blood - the murine HSC mutation rate is only two-fold higher than  
460 human<sup>9-11</sup> despite a 35-fold shorter lifespan, suggesting that mutation accrual patterns across  
461 species are under tissue-specific evolutionary constraints. Indeed, somatic mutation rates in  
462 germline cells are lower in mouse than in human<sup>61</sup> and under the influence of distinct factors such  
463 as effective population size and age of reproductive maturity<sup>62</sup>. In blood, it is plausible that a low  
464 somatic mutation rate is required to minimise the entry of detrimental disease-causing mutations,  
465 which when combined with a large stem cell pool, may also reduce the fixation probability of any  
466 such mutations. Alternatively, it is also possible that the mutation rate may not reflect  
467 haematopoietic adaptation in the mouse but rather a historical evolutionary constraint or a feature  
468 of phylogenetic legacy<sup>63</sup>.

469 Patterns of somatic evolution in humans provide one plausible mechanism by which ageing  
470 phenotypes occur. The presence of clonal expansions in elderly human blood driven by somatic  
471 mutations is associated with diseases of ageing. However, in the laboratory mouse, which also  
472 displays phenotypes of ageing including increased cancer incidence<sup>64</sup>, we only observe small  
473 mutation-driven clonal expansions in blood by the end of life, suggesting that any role age-  
474 associated haematopoietic oligoclonality plays in human ageing is unlikely to be shared by the  
475 laboratory mouse. The dramatically different population structures of haematopoiesis in the old  
476 mouse versus old human, together with the small clones (necessitating sensitive detection

477 methods) are crucial factors to be considered when using murine models for future studies of  
478 natural CH or haematopoietic ageing. Alternative model organisms, such as non-human primates,  
479 display similar stem cell cycling behaviour<sup>65,66</sup> to humans and larger age-related clonal  
480 expansions<sup>67,68</sup>, and thus may be suited to evaluate aspects of native hematopoietic dynamics  
481 across the lifespan.

482 Native murine clones do expand upon systemic exposures and recapitulate patterns previously  
483 observed in correlative humans studies<sup>47,48</sup> and in exposures administered following murine  
484 transplant<sup>69-72</sup> (reviewed in depth in refs. <sup>17,45,46</sup>). We postulate that the size of clonal expansions is  
485 constrained in mouse due to infrequent HSC self-renewing divisions during homeostatic conditions.  
486 Our data fit with mouse stem cells self-renewing every six weeks (1.8-13.2 weeks), within the broad  
487 range of previous estimates from once in 4 to 24 weeks<sup>73-75</sup>. Whilst this is more frequent than  
488 human HSCs (estimated to divide at 1-2 times a year), for patterns of oligoclonality in humans to  
489 be recapitulated in the much shorter-lived mouse via genes conferring similar fitness advantages,  
490 stem cells would need to self-renew much more frequently. It is possible that mouse strains thought  
491 to have higher HSC turnover<sup>76</sup>, or maintained for longer periods in more “wild”-like microbial  
492 environments, would exhibit higher levels of native CH. Additional studies to characterise such  
493 strains and environments would be of interest.

494 Nevertheless, our data highlight conserved selection landscapes in mouse with detectable CH in  
495 both homeostatic haematopoiesis and under stress when using highly sensitive sequencing. With  
496 our observation of evolutionarily conserved constraints on population dynamics of blood, together  
497 these drive a distinct pattern of somatic evolution over the murine lifespan. These data provide a  
498 framework for the interpretation of future studies of haematopoietic stem cell biology and ageing  
499 using the laboratory mouse.

## 500 **Figure Legends**

### 501 **Figure 1. Somatic mutations in murine stem cell-derived haematopoietic colonies**

502 **A)** Study approach. Single-cell derived colony whole genome sequencing (WGS) of long-term  
503 haematopoietic stem cells (HSC) and multipotent progenitors (MPP) to study somatic mutations,  
504 lineage relationships and population dynamics, top; targeted duplex-sequencing of peripheral  
505 blood to identify small clonal expansions and fitness landscapes, bottom. **B)** Number of whole  
506 genomes (n=1305) of HSC- and MPP-derived colonies that underwent phylogenetic construction  
507 for each female mouse (n=6). Plots are coloured according to HSC- or MPP-derived colonies,  
508 darker and lighter shades, respectively. **C)** Burden of individual single base substitutions (SBS)  
509 observed in HSCs (n=908) from each donor. Points are coloured as in panel B. Line shows linear  
510 mixed-effect regression of mutation burden observed in colonies. Shaded areas indicate the 95%  
511 confidence interval. **D)** Comparison of SBS burden between HSC- and MPP-derived colonies from  
512 the same mice. SBS burden from HSCs are shown as circles and burden from MPPs are shown  
513 as squares. H, HSC; M, MPP, shown above animal ID. **E)** SBS burden across HSCs (data as in  
514 panel C), whole blood, and individual colonic crypts in the three aged mice. Error bars denote 95%  
515 confidence interval. Peripheral blood and colonic crypt somatic mutation burdens were measured  
516 with nanorate sequencing and WGS, respectively.

### 517 **Figure 2. Phylogenetic trees of HSCs and MPPs from a young and old mouse**

518 **A-B)** Phylogenies were constructed from young (3-months, A) and aged (30-months, B) female  
519 mice using the pattern of sharing of somatic mutations among HSC (blue) and MPP (red) colonies.  
520 Each tip represents a single colony. Branch lengths represent mutation number, corrected for  
521 varying sequencing depth of descendant colonies. Branches and coalescence colours reflect the  
522 identity of descendent colonies with HSCs in blue and MPPs in red, respectively. Branches where  
523 we are unable to infer the established cell type for one or more lines of descent are coloured black.  
524 **C)** To determine the degree of phylogenetic relatedness between HSC and MPP, we measured  
525 the amount of HSC-MPP mixing within clades. If an MPP had a recent HSC ancestor, clades should  
526 contain both cell types. We thus compared the “observed” versus “expected-by-chance” clade  
527 mixing behaviour. The mixing metric for a clade is the absolute difference between the proportion  
528 of HSCs in a clade and the expected value under equal sampling, 0.5; this metric is then averaged  
529 for all clades in a phylogeny. The vertical bar reflects the observed average clade mixing metric  
530 within the constructed phylogenies. The filled distributions reflect average clade mixing metrics that

531 would be expected by random chance or more frequent intermixing of HSCs and MPPs, and were  
532 generated by reshuffling the tip cell identities within the tree. HSC or MPP colonies are designated  
533 as being in the same clade if they share a most recent common ancestor after 25 mutations,  
534 corresponding to early foetal development. Only clades with more than 3 colonies are considered.  
535 **D)** Distributions of the number of cell identity changes required per colony to capture the observed  
536 tip states. The number of cell identity changes assuming an ‘HSC-first’ model (HSCs first give rise  
537 to MPPs) is shown in blue. The required cell identity changes for the opposite ‘MPP-first’ model, in  
538 which MPPs first give rise to HSCs, is shown in red. The null distribution, in which tip states are  
539 randomly reshuffled is shown in grey. **E)** Cell-type probability trajectories displaying specification  
540 to HSC or MPP states under a simple 3-state ontogeny model (Methods). In 30-month donors  
541 (right), we observe equal generation of HSC and MPP from embryonic progenitors (EMB) and can  
542 reject an “HSC-first” model. In 3-month donors (left), we observe relatively increased generation of  
543 HSCs from EMB and can not reject an “HSC”-first model. The displayed trajectories are based on  
544 iterating the maximum likelihood based Markov chain starting at the embryonic state. Thickness of  
545 arrows reflect the proportion of overall transitions from the EMB state to HSC and MPP states, and  
546 between HSC and MPP states. The cell identity transition rates are derived in Supplementary Note  
547 2.

### 548 **Figure 3. Population dynamics and selection in the murine stem cells**

549 **A)** Population trajectories estimated separately in HSCs and MPPs using Bayesian phylodynamics  
550 for the six samples shown in Fig 2.A-B and Extended Data Fig.2. The dark blue (HSC) and red  
551 lines (MPP) indicate the mean effective population trajectory; shaded areas are 95% confidence  
552 intervals. Vertical dashed lines separate trajectories into early life and adulthood age periods, in  
553 which different population size behaviour are observed. Inset values indicate posterior density  
554 estimates of population size ( $N$ ), symmetric cell division rate per week ( $\lambda$ ), and their ratio in ( $N/\lambda$ ) in  
555 HSC-years, as derived from approximate Bayesian computations. **B)** Haematopoietic stem and  
556 progenitor cell (HSPC) prevalence during murine ageing. The relative abundance of total HSPCs  
557 (left, defined as the LSK compartment) and individual HSPC subpopulations (right) are compared.  
558  $MPP^{L_y}$  are lymphoid-biased progenitors,  $MPP^{GM}$  are myeloid-biased progenitors, based on current  
559 immunophenotypic definitions<sup>55</sup>. **C)** Shannon diversity index for each phylogeny calculated using  
560 the number and size of unique clades present at 50 mutations molecular time. Mouse points are  
561 coloured as in Fig.1B. Grey dots depict results from data published in Mitchell et al<sup>9</sup>. **D)** Normalised

562 ratio of non-synonymous to synonymous somatic mutations (dN/dS) for somatic mutations  
563 observed across aged and young animals overlaps with 1 suggesting no departure from neutrality.

#### 564 **Figure 4. Clonal haematopoiesis during normal ageing in mouse**

565 **A)** Dot-plot describing incidence of clonal haematopoiesis in mice at increasing age. Each vertical  
566 column represents a single mouse sample with detected clone size and consequence indicated by  
567 dot size and colour. Strain is C57BL/6J. **B)** Barplot summarising clone count per sample as  
568 illustrated in A. Differences in clone incidence were quantified by the Kruskal-Wallis test. **C-D)**  
569 Murine clonal haematopoiesis incidence in the laboratory strains **C)** B6FVBF1/J (F1 hybrid from  
570 crossing inbred C57BL/6J x FVB/NJ), and **D)** HET3 (a four-way cross between C57BL/6J,  
571 BALB/cByJ, C3H/HeJ, and DBA/2J). **E)** Clone size changes in samples collected serially over 4  
572 months. Clones are coloured by mutation. **F)** dN/dS ratios for targeted genes mutated in murine  
573 clonal haematopoiesis. Variants from all donors in A were used to determine gene level dN/dS  
574 ratios. \* represents dN/dS >1 with q-value <0.1.

#### 575 **Figure 5. Haematopoietic perturbation modulates selection landscapes**

576 Clonal haematopoiesis prevalence in aged mice following **A)** normalised microbial experience  
577 (NME), **B)** *M. avium* infection, **C)** cisplatin treatment, and **D)** 5-FU myeloablation. At final sampling,  
578 aged mice were 30-months-old for the NME experiments in panel A), and were 25-months-old for  
579 the perturbation experiments in panels B), C), and D). Enrichment of clonal prevalence and dN/dS  
580 ratios departing from parity following treatment are shown for each gene. Survival curves and  
581 experimental endpoint blood counts are displayed for B) and C), using log-rank and two-sided t  
582 tests, respectively. Treatment schedules are as displayed or described in Methods.

#### 583 **Figure 6. The fitness landscape of known drivers of clonal haematopoiesis**

584 **A)** Reverse cumulative density for all synonymous (including flanking intronic regions in targeted  
585 bait set) and nonsynonymous somatic variants detected using duplex sequencing from mice aged  
586 24-25 months, arranged by increased variant allele fraction (VAF). The relative density of  
587 synonymous (and flanking intronic) variants, which are assumed to have neutral fitness, yields an  
588 estimate for  $N/\lambda$ , the ratio of population size and symmetric cell division rate (per year). The  
589 synonymous and nonsynonymous mutation rates ( $\mu$ , base pairs per year) can then be estimated  
590 using a maximum likelihood approach. **B)** Distribution of fitness effects for nonsynonymous  
591 mutations.

592 **Extended Data Figure 1. Cell isolation strategy and quality control**

593 **A)** Sorting strategy for single HSCs and MPPs from young and aged mice. Progenitor-enriched  
594 bone marrow was stained as described in the Methods, and then single cells were sorted into  
595 individual wells for *in vitro* expansion. **B)** Colony-forming efficiency of sorted HSCs and MPPs for  
596 each sample. Each bar represents the listed cell type and underlying sample ID. **C)** Variant allele  
597 fraction (VAF) distribution of all variants within a colony that pass filtering, shown for a  
598 representative clonal colony that passed sample QC (left) and a non-clonal colony that passed  
599 sample QC (right). After variant filtration, the VAF distribution of a colony's variants is centred  
600 around 50% in clonal colonies, but less than 50% in non-clonal colonies. **D)** Representative image  
601 of two colonic crypts isolated by laser capture microdissection. **E)** Correlation between total single  
602 base substitution burden and depth, for all colonies from sample M7180, shown before (left) and  
603 after (right) sequencing depth correction. **F)** Trinucleotide spectra from aggregated somatic  
604 mutations mapped to shared (truncal) or private branches of phylogenetic trees. Signatures are  
605 highly similar, suggesting artefacts are not relatively enriched in either portion of reconstructed  
606 trees.

607 **Extended Data Figure 2. Additional phylogenetic trees from young and aged mice**

608 Phylogenies for **A-B)** 2 additional young (3-month) mice and **C-D)** 2 additional aged mice (30-  
609 month), presented as described in Figure 2.

610 **Extended Data Figure 3. Early-in-life phylogenetic patterns and cross-tissue mutations**

611 Phylogenies from aged (left) and young (right) HSCs zoomed into the first 12 mutations molecular  
612 time. Polytomies in the branching structure, which represent cell division without mutation  
613 acquisition, are enriched among early-in-life cell divisions at the tops of the phylogenies. Variants  
614 shared with matched colonic crypts are layered onto the trees as pie charts. Pie chart fullness  
615 represents the proportion of colonic crypts in which the mutation present on the haematopoietic  
616 phylogeny was observed. Sample M7183 lacked sufficient early life diversity (<10 unique lineages  
617 within 12 mutations molecular time) and thus was excluded.

618 **Extended Data Figure 4. Mutational processes in murine stem cells**

619 **A)** Signature extraction overview. Trinucleotide spectra from all single-base substitutions (SBS)  
620 (top), were used for signature extraction as described in the Methods. Three signatures identified  
621 as SBS1, SBS5, and SBS18 best described the catalogue of mutations observed (cosine

622 similarity=0.997). **B)** Linear mixed-effect regression of signature-specific mutation burdens  
623 observed in colonies. Shaded areas indicate the 95% confidence interval. **C)** Signature attribution  
624 in phylogenies. Individual branches of HSC phylogenies are overlaid with signature contribution  
625 proportions. SBSs assigned to each branch were fit to SBS1, SBS5 or SBS18. **D)** Signature-  
626 specific mutation accumulation in all branches across phylogenies. Early-life branchpoints, located  
627 at the top of a given phylogeny, and shown as an inset.

#### 628 **Extended Data Figure 5. Phylogeny comparison between aged human and mouse**

629 **A)** Representative ultrametric phylogenies from the three oldest humans described in Mitchell *et*  
630 *al.*<sup>9</sup> The published trees have been randomly downsampled to 100 colonies (tips). **B)** Aged mouse  
631 phylogenies, also downsampled to 100 colonies, to allow comparison of topological structure. The  
632 median lifespan for human and mouse species are labelled and were derived as described in  
633 Supplementary Note 1. Full murine phylogenetic trees are shown in Figure 2A-B and Extended  
634 Data Figure 2).

#### 635 **Extended Data Figure 6. Approximate Bayesian inferences**

636 Results from approximate Bayesian computation (ABC) inference of **A)** population size ( $N$ ), **B)**  
637 symmetric division rate per week ( $\lambda$ ), and **C)** death rate per week ( $\nu$ ) for the three 30-month-old  
638 mice. Blue lines represent the prior density of parameters; red lines represent the posterior  
639 densities. Median posterior density estimates and 95% credibility intervals are displayed for each  
640 parameter per sample. The prior density for the death rate was bounded to ensure the growth rate  
641 ( $\lambda - \nu$ ) remained positive, as observed in *phyloDYN* trajectories in Figure 3. **D)** Joint density  
642 distributions indicating optimal parameters of population size and division rates that explain  
643 observed phylogenetic trees. The estimated  $N/\lambda$ , in HSC-years, is shown with 95% credibility  
644 intervals. Data from the three aged mice are shown.

#### 645 **Extended Data Figure 7: Extended phylogenetic trees (HSC, MPP, and early progenitor).**

646 **A-C)** Extended phylogenies were created for three 30-month mice using the pattern of sharing of  
647 somatic mutations among HSCs (blue), MPPs (red), and the mixed LSK (Lineage-, Sca1+, c-kit+)  
648 hematopoietic progenitor compartment. The LSK compartment contains HSCs and MPP, and  
649 additionally contains the myeloid-biased MPP<sup>GM</sup> (orange) and lymphoid-biased MPP<sup>LY</sup> populations  
650 (green). LSK subcompartments were identified at time of single cell sorting using a consensus  
651 definition<sup>55</sup>. Each tip represents a single colony. Branch lengths represent mutation numbers. **D-E)**

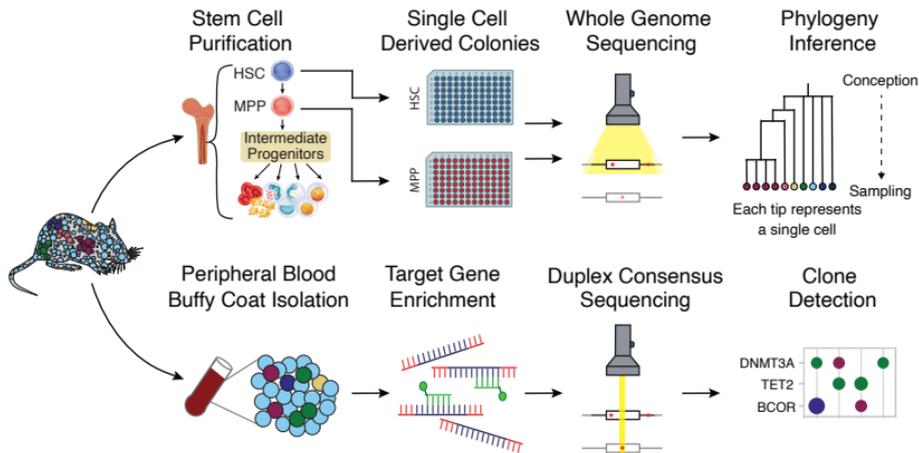
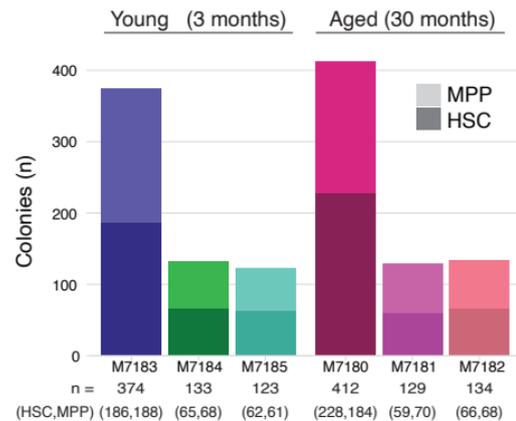
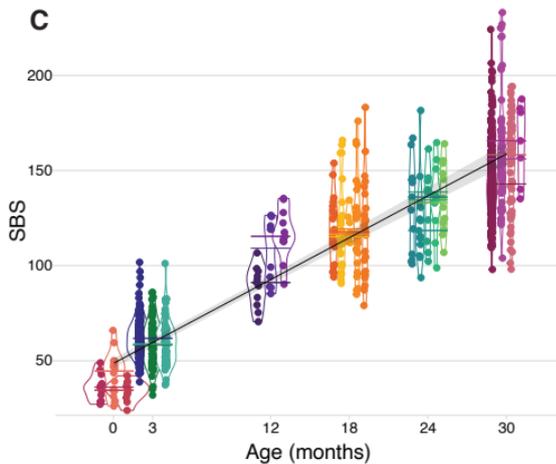
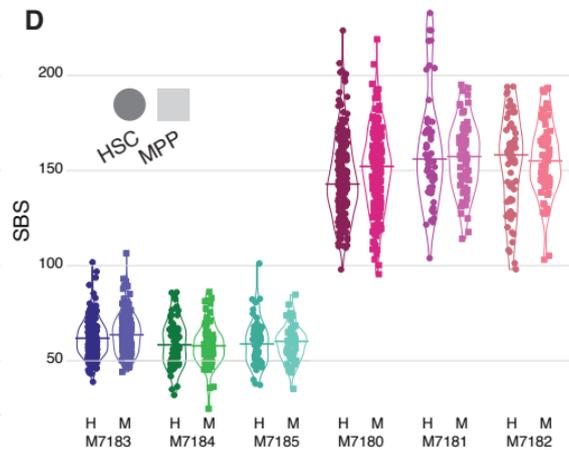
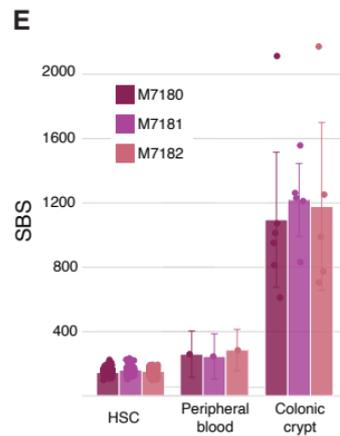
652 Clade mixing metrics for MPP<sup>GM</sup> and MPP<sup>Ly</sup> colonies used to evaluate interrelatedness with HSC  
653 and MPP. HSC, MPP and MPP<sup>GM</sup> or MPP<sup>Ly</sup> were designated as being in the same clade if they  
654 share a most recent common ancestor after 25 mutations, corresponding to early foetal  
655 development. Only clades with more than 3 colonies are considered. The vertical bar reflects the  
656 average clade mixing metric observed in the constructed phylogenies, while distributions reflect the  
657 average clade mixing metric expected random chance, estimated by reshuffling the tip states. If  
658 the observed value (vertical bar) significantly deviated from random chance (filled distribution), then  
659 there would be minimal overlap between the observed data and the random reshuffling distribution.  
660 The average clade mixing metric for MPP<sup>GM</sup> compared to HSCs (blue) and MPPs (red) is shown  
661 in **D**). The similar measure of interrelatedness of MPP<sup>Ly</sup> to HSCs and MPPs is shown in **E**).

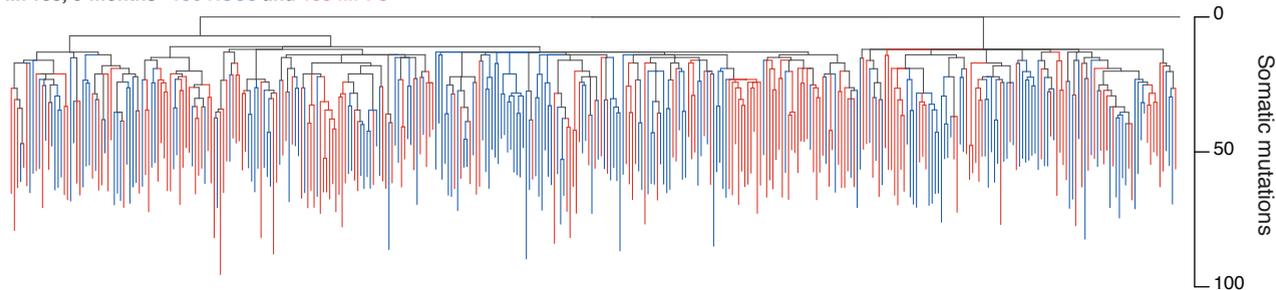
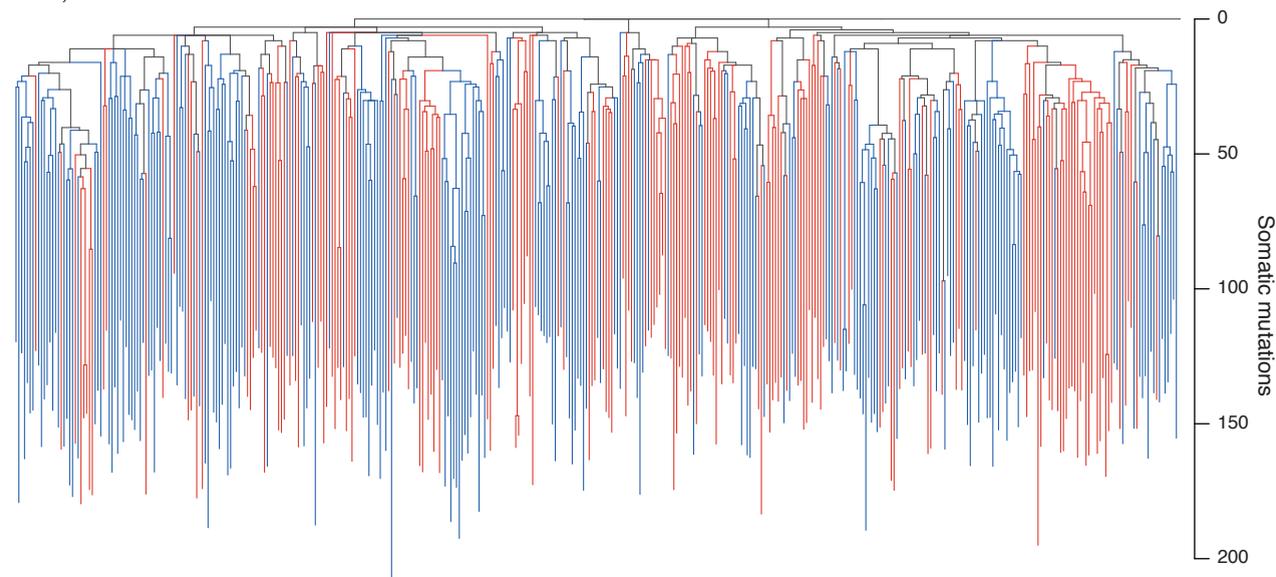
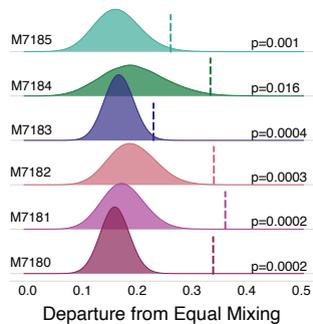
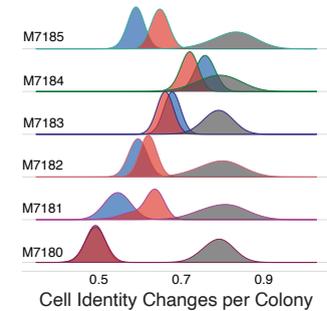
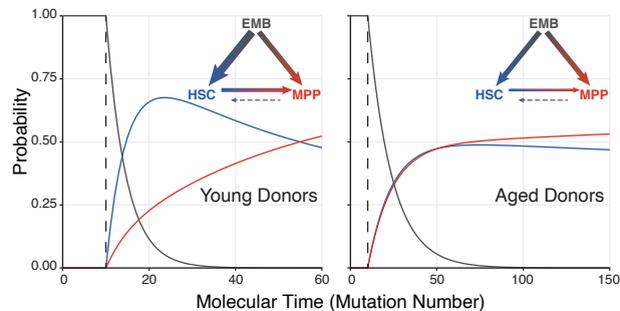
662 **Extended Data Figure 8: Mutation overlap between phylogenies and peripheral blood.**

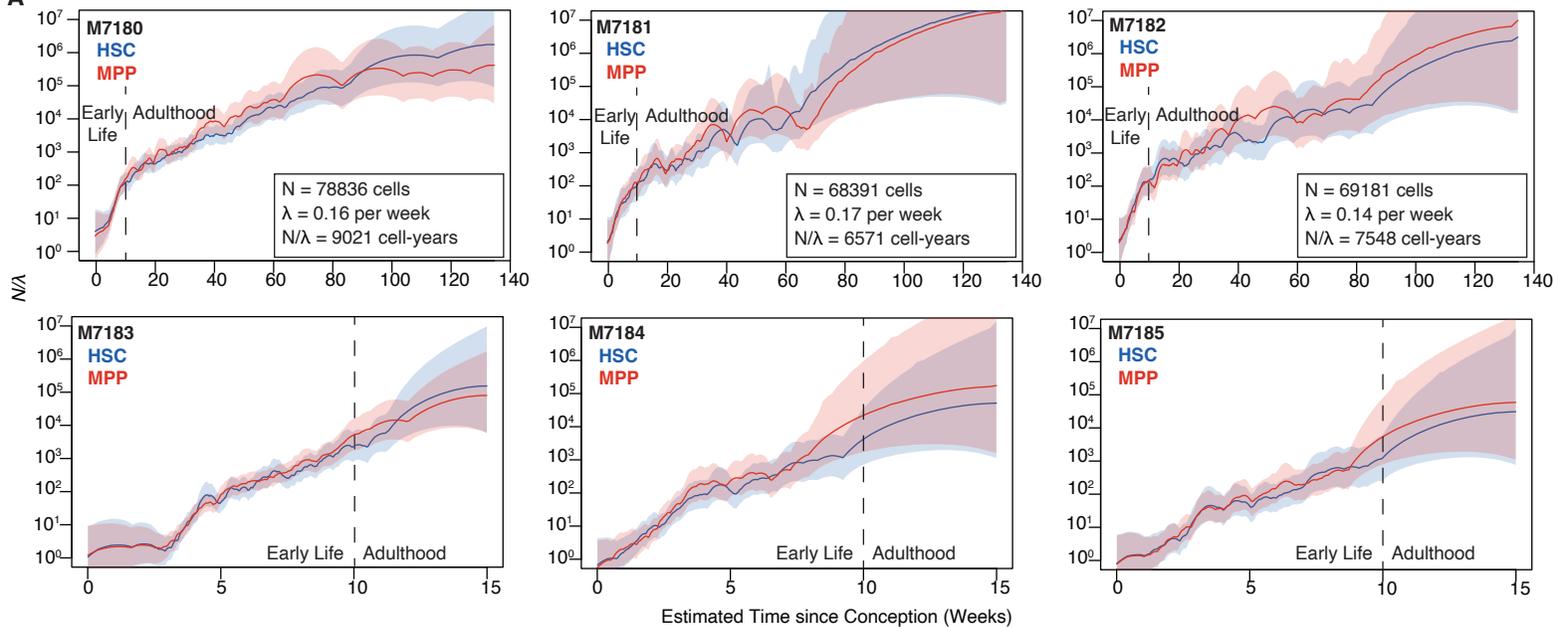
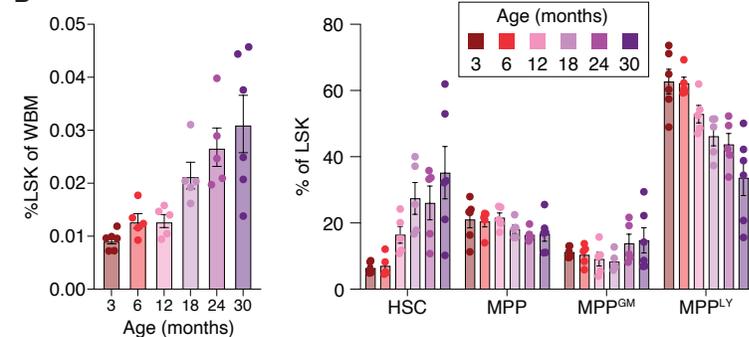
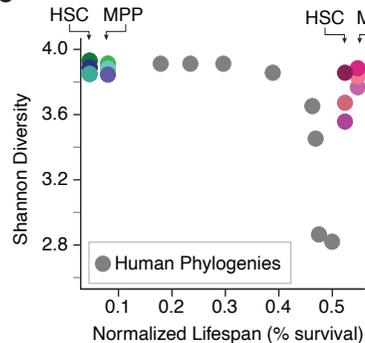
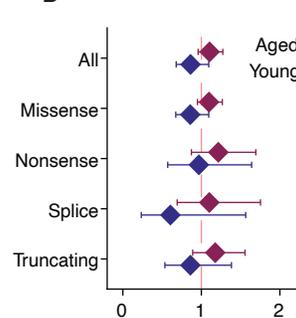
663 **A**) Phylogenies for three aged mice (as described in Extended Data Figure 7A-C) constructed to  
664 only include private branches targeted with the peripheral blood baitset. Branch shading indicates  
665 the maximum VAF among branch-specific variants captured in peripheral blood. The sampled cell  
666 immunophenotype is indicated by dot colour at the bottom of each private branch. **B**) VAF  
667 trajectories of HSC and MPP variants shared in peripheral blood. The aggregate VAF across  
668 molecular time is calculated using Gibbs sampling (Methods). Earlier molecular time corresponds  
669 to further in the ancestral past. Shaded regions denote 95% confidence intervals of the VAF  
670 estimates.

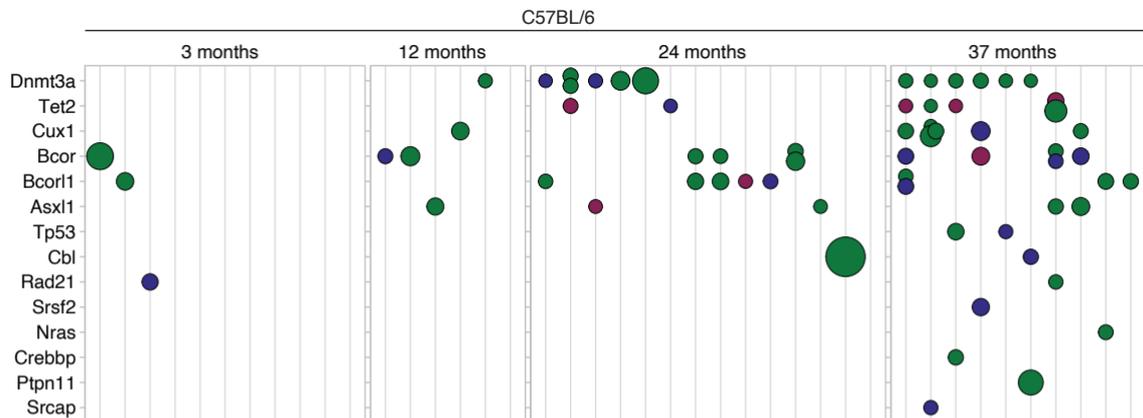
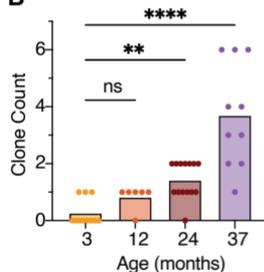
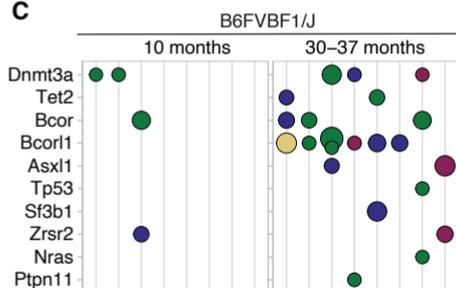
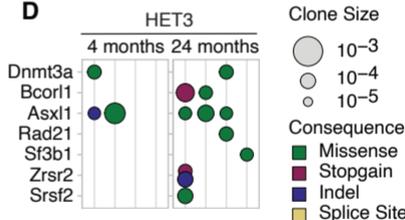
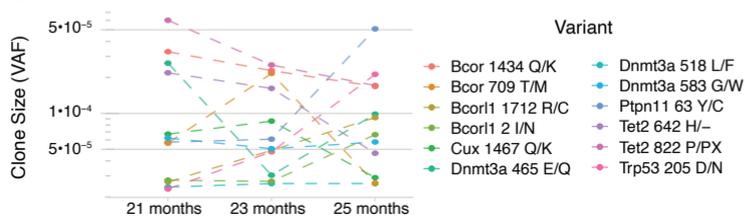
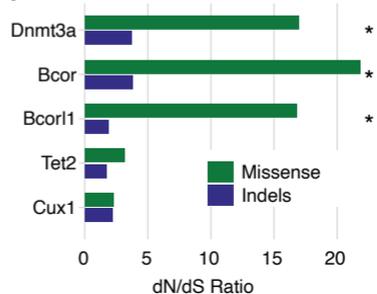
671 **Extended Data Figure 9: Peripheral blood VAF of variants shared with HSCs and MPPs.**

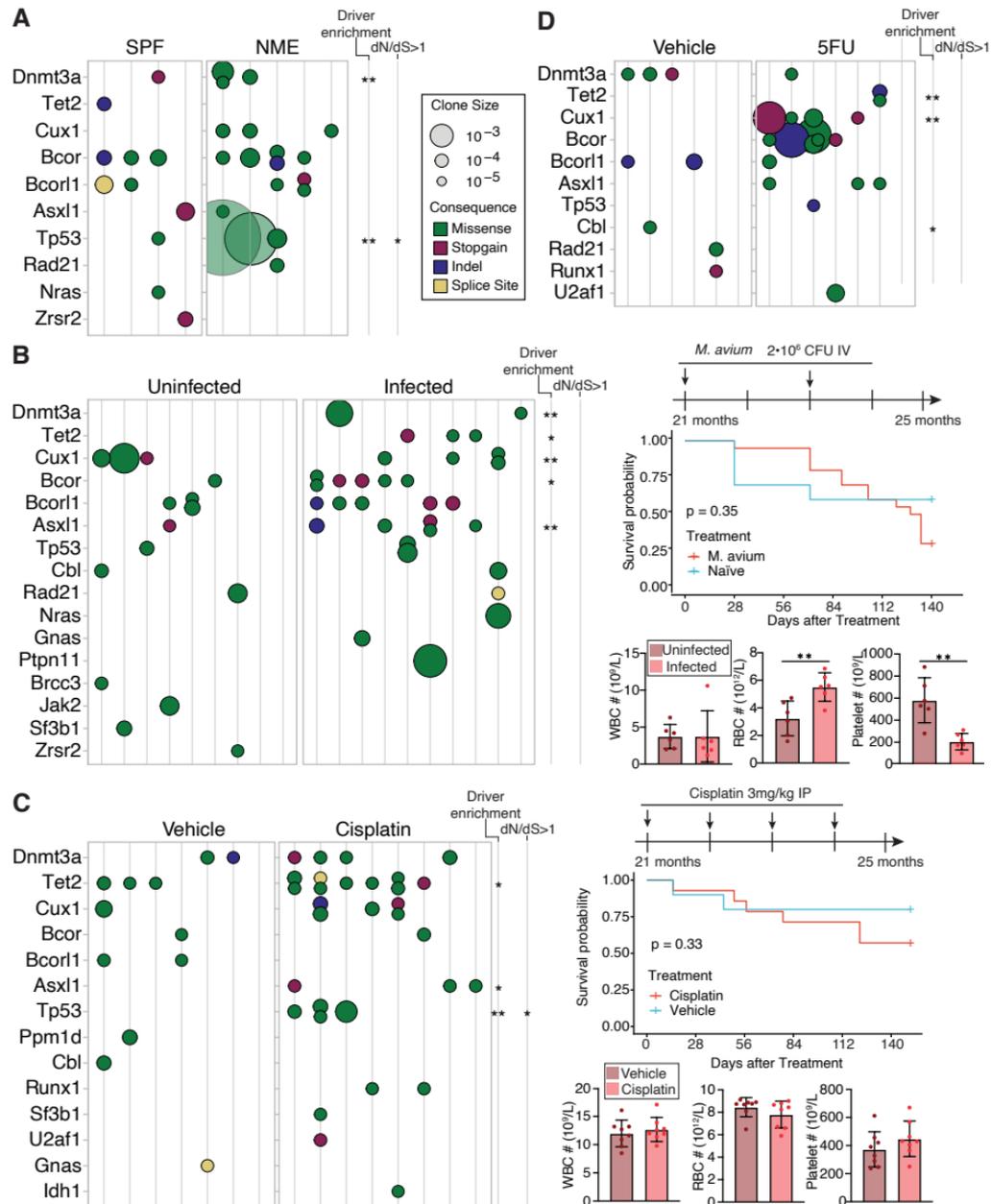
672 Baitset mutation-specific HSC and MPP phylogenies are shown for each 30-month mouse. Each  
673 branch shows mutations that were detected in peripheral blood in descending VAF order. On each  
674 branch, a row denotes a single variant mapped to that specific branch. Red fill denotes the  
675 peripheral blood VAF for the variant. VAF is denoted on a log scale from 10<sup>-5</sup> to 1; internal divisions  
676 are marked from left to right at VAF 0.0001, 0.001, 0.01, and 0.1. HSC trees are shown on the left  
677 with blue dots at terminal branches; MPP trees are shown on the right with red dots. Trees are  
678 downsampled to allow equivalent comparison between HSC and MPP branches. Only variants  
679 seen in peripheral blood with a depth > 100X are shown.

**Figure 1****A****B****C****D****E**

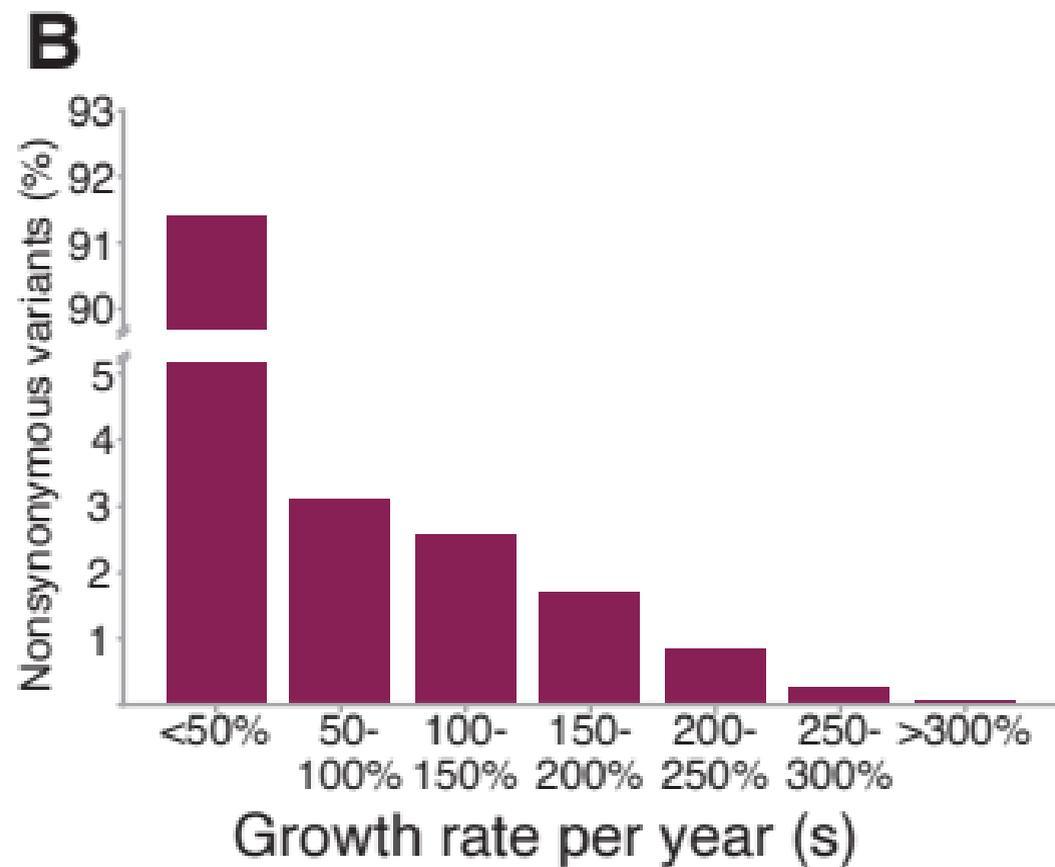
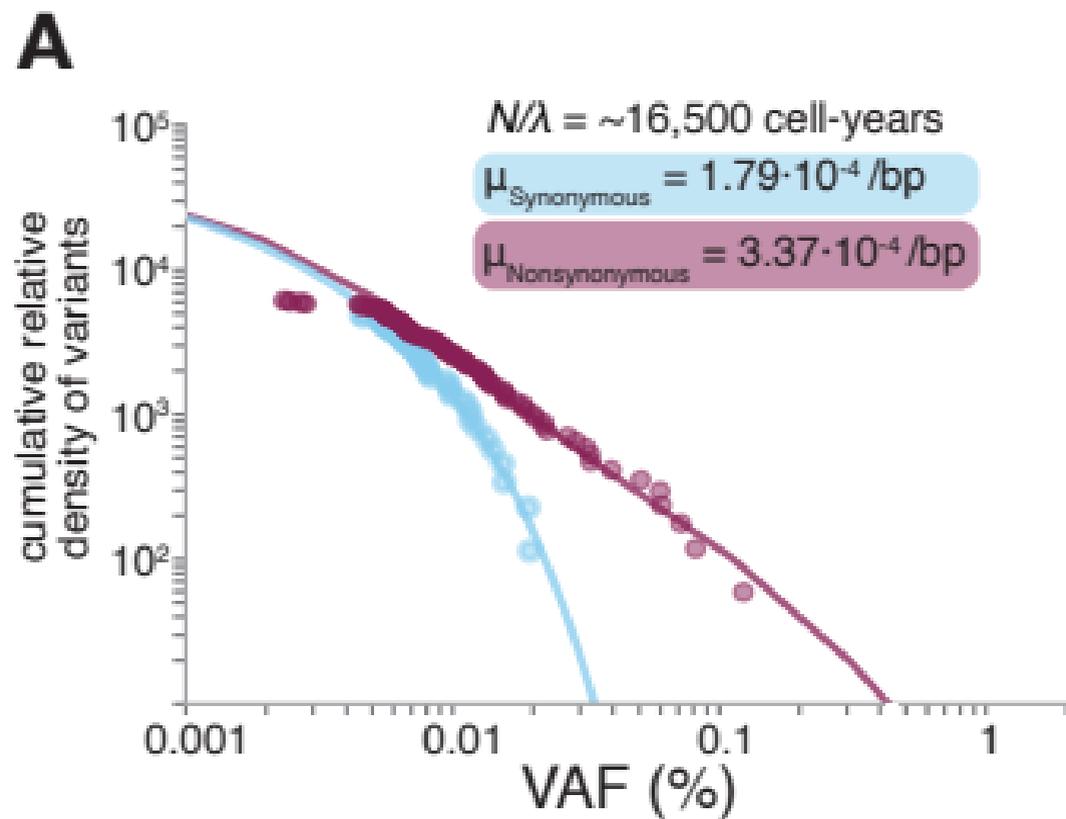
**Figure 2****A**M7183, 3-months **186 HSCs** and **188 MPPs****B**M7180, 30-months **228 HSCs** and **184 MPPs****C****D****E**

**Figure 3**
**A**

**B**

**C**

**D**


**Figure 4****A****B****C****D****E****F**

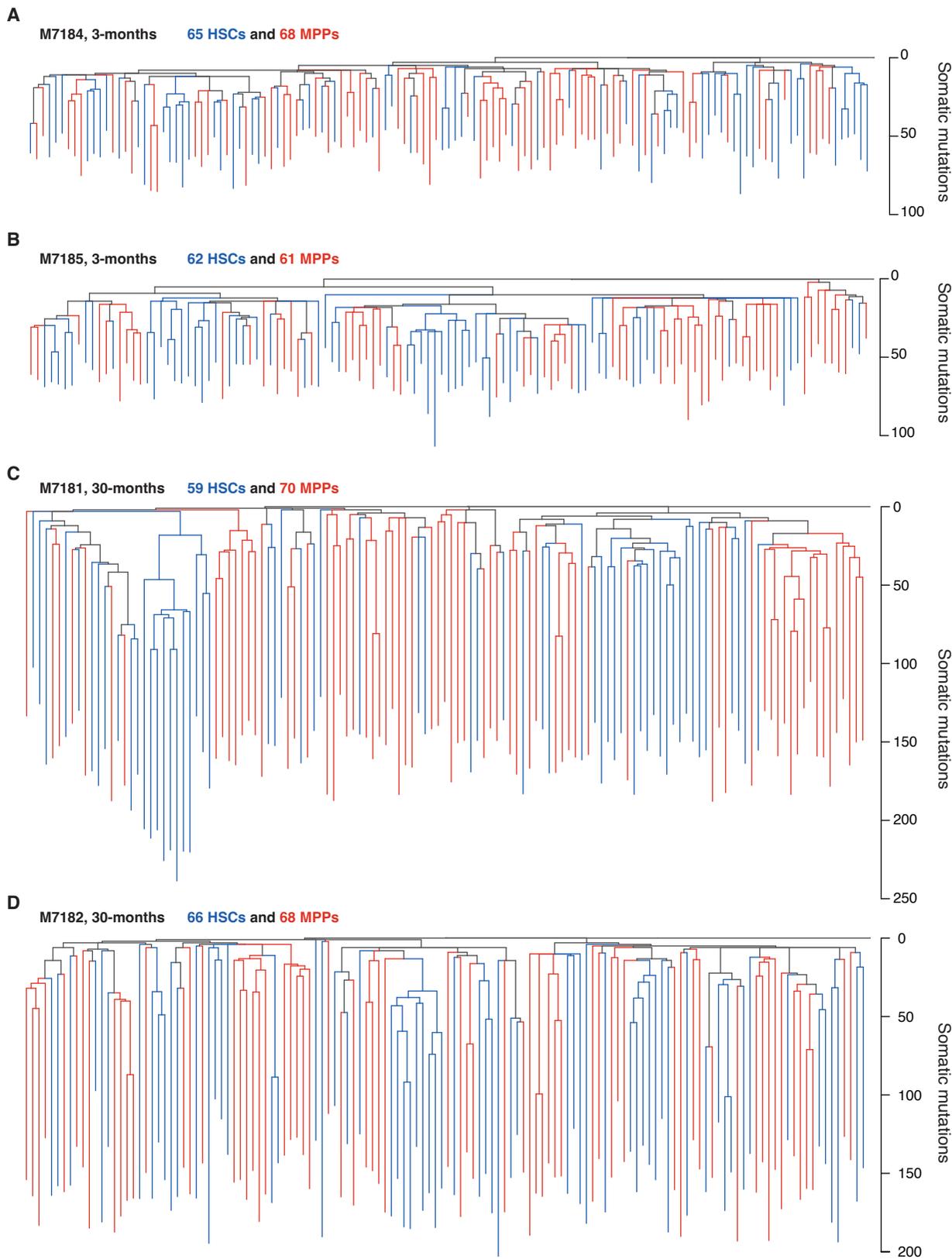
**Figure 5**

# Figure 6





## Extended Data Figure 2

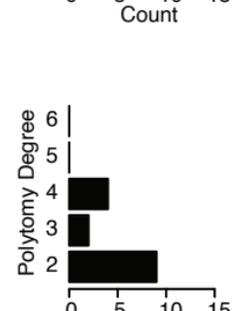
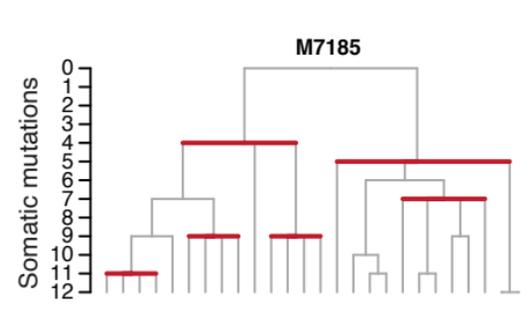
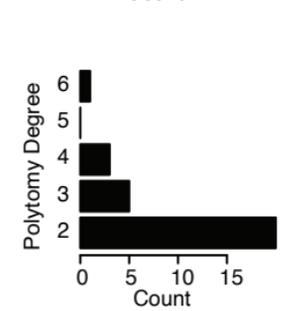
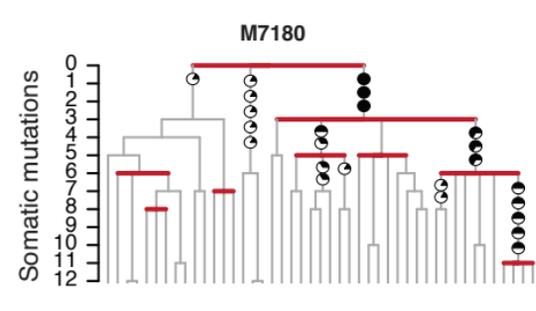
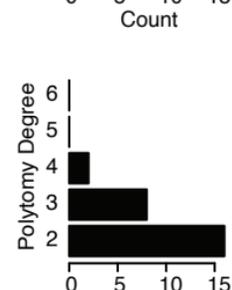
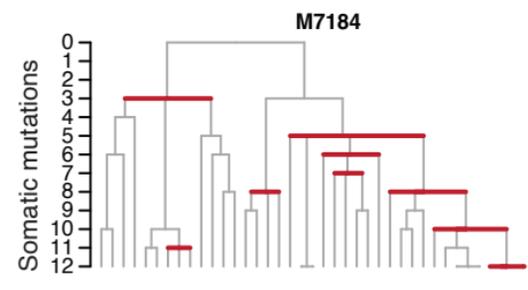
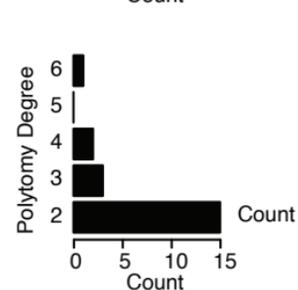
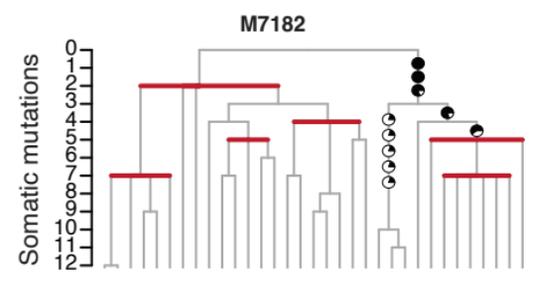
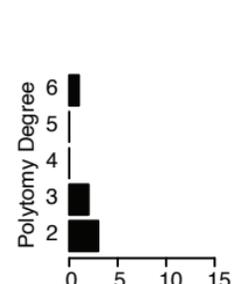
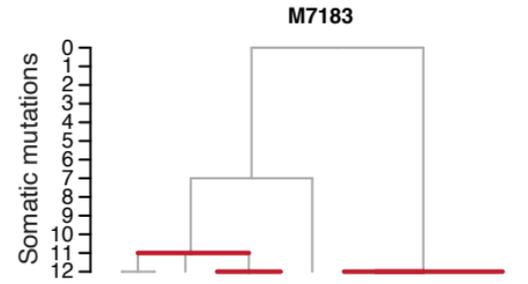
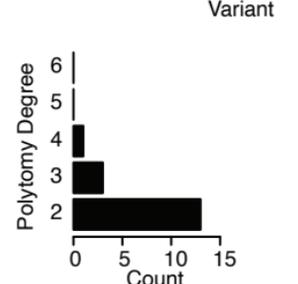
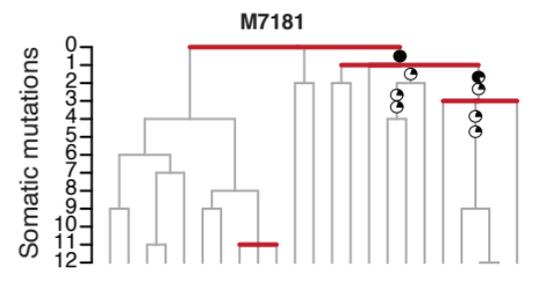


# Extended Data Figure 3

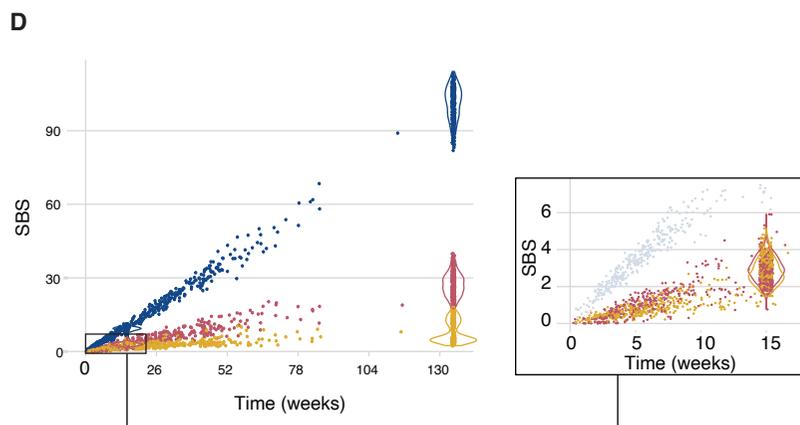
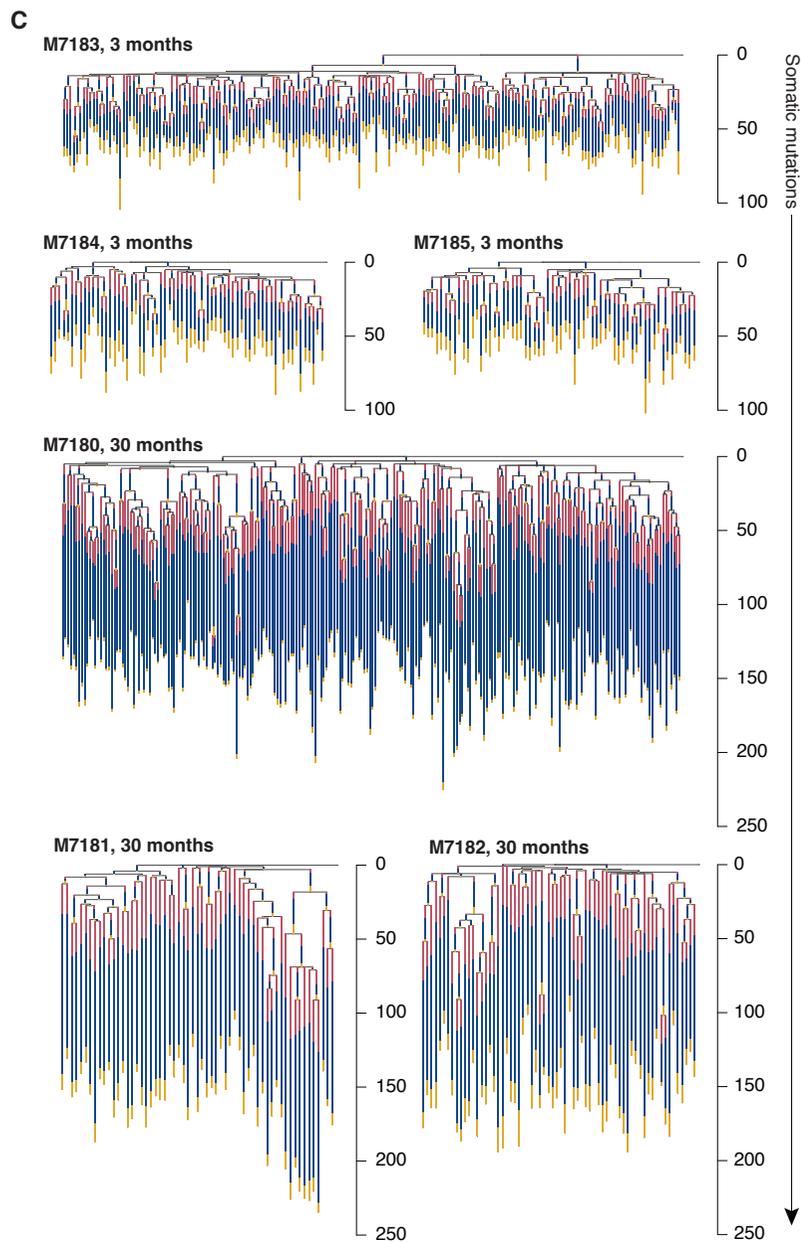
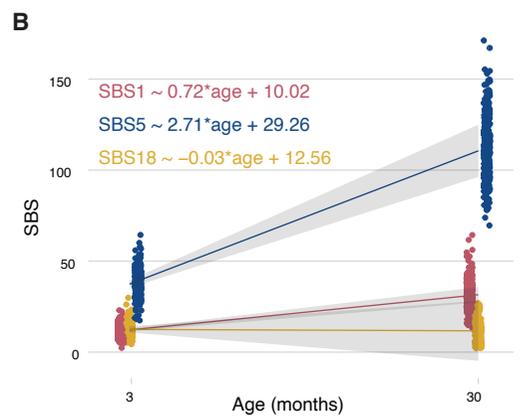
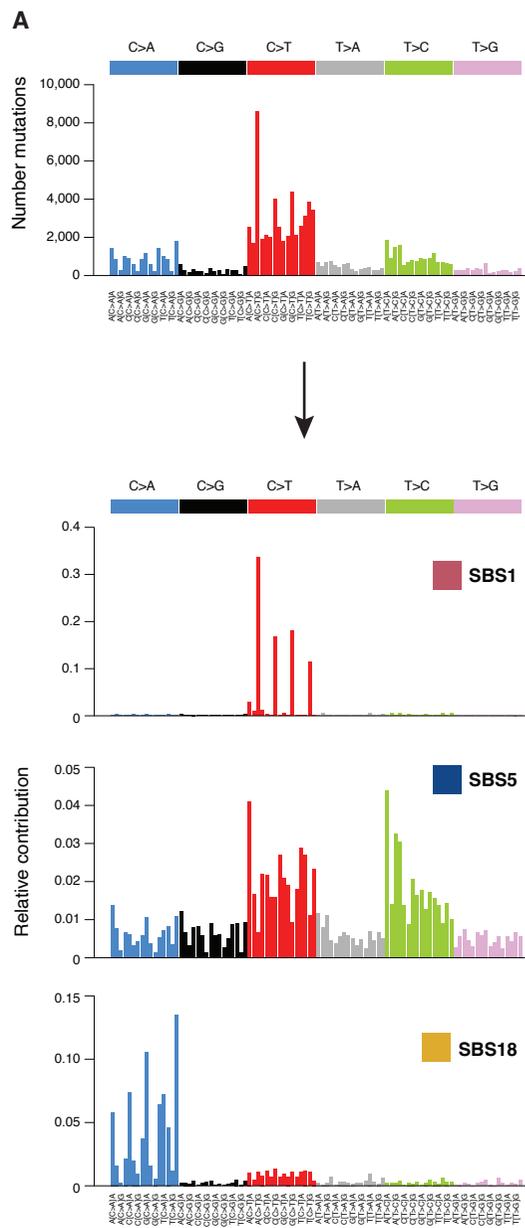
Aged (30 months)

— Polytomy  
● Shared Colonic Variant

Young (3 months)



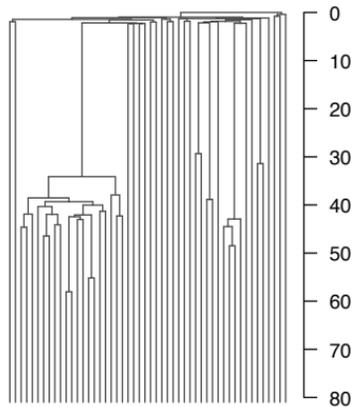
# Extended Data Figure 4



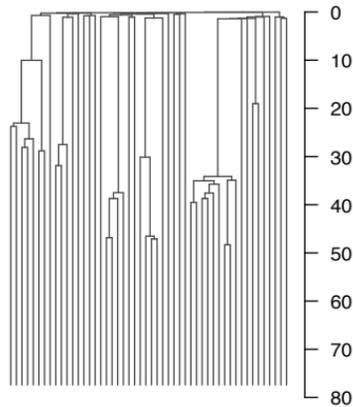
# Extended Data Figure 5

**A**

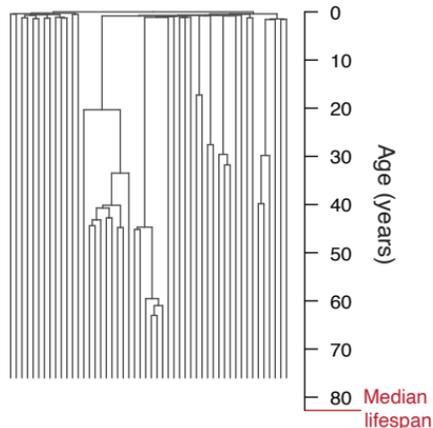
**KX003, 81-years**



**KX004, 77-years**

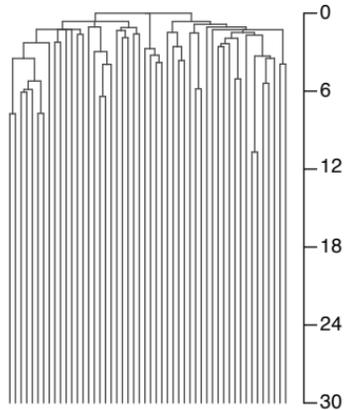


**KX008, 76-years**

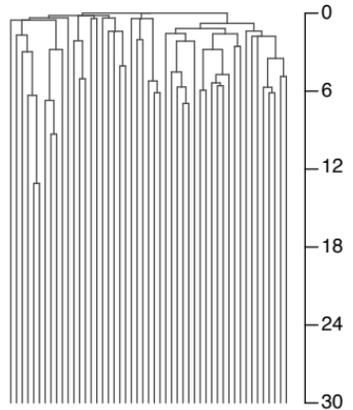


**B**

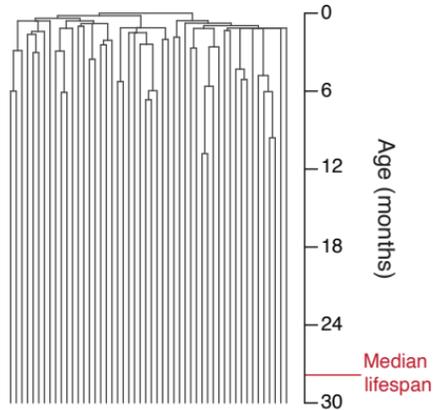
**M7180, 30-months**



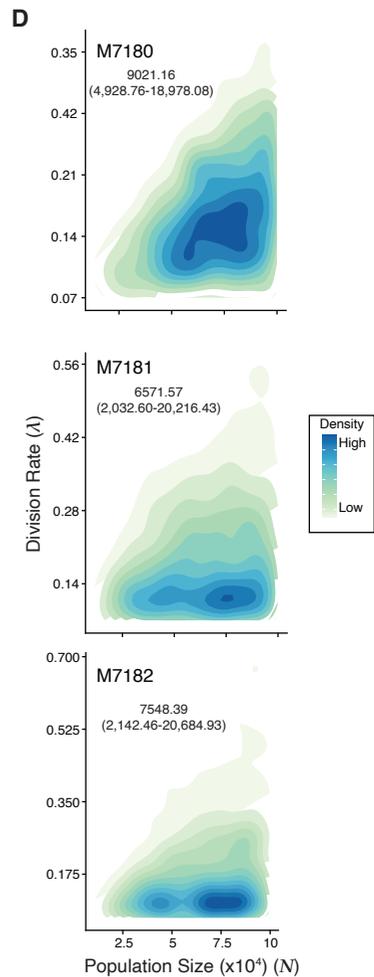
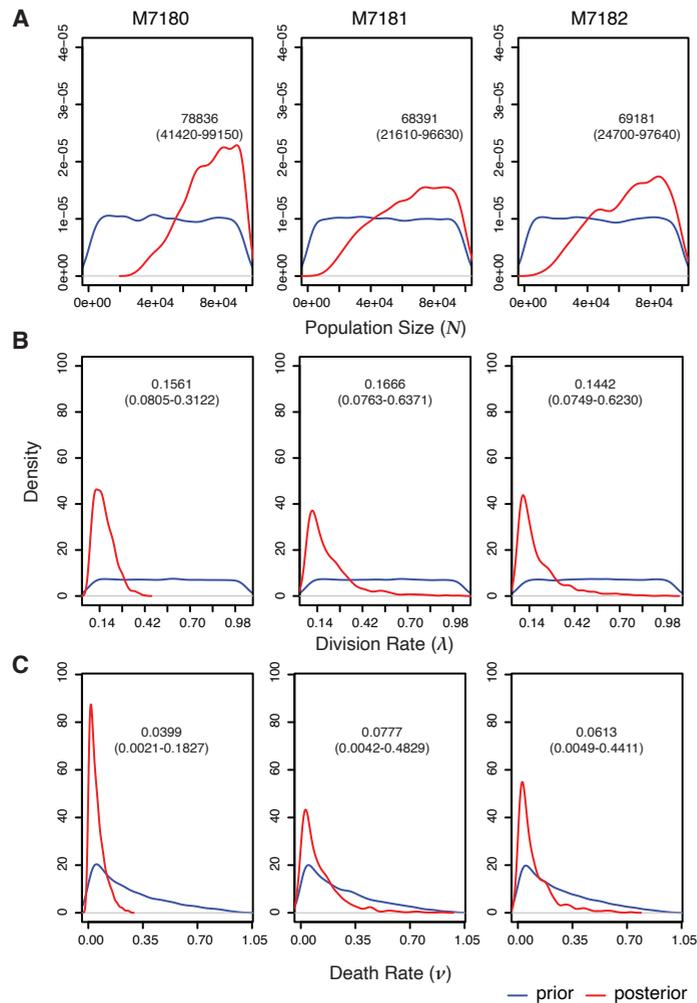
**M7182, 30-months**



**M7183, 30-months**



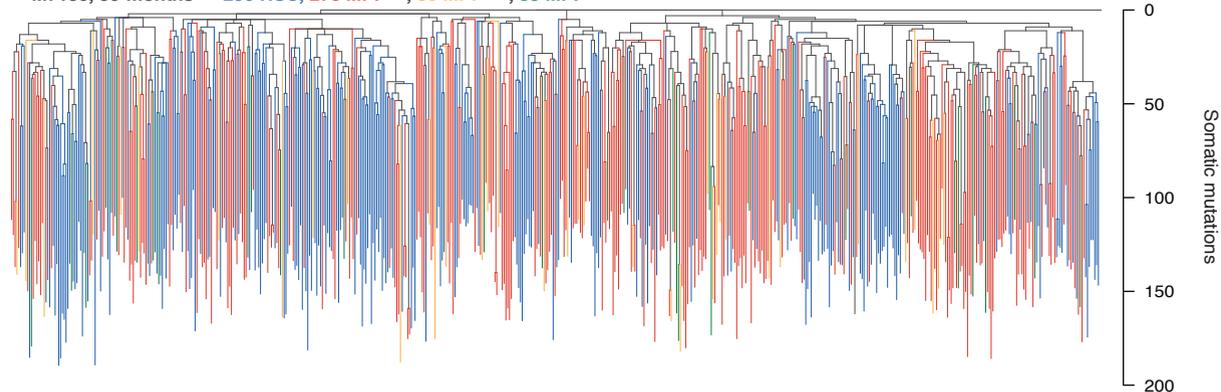
# Extended Data Figure 6



# Extended Data Figure 7

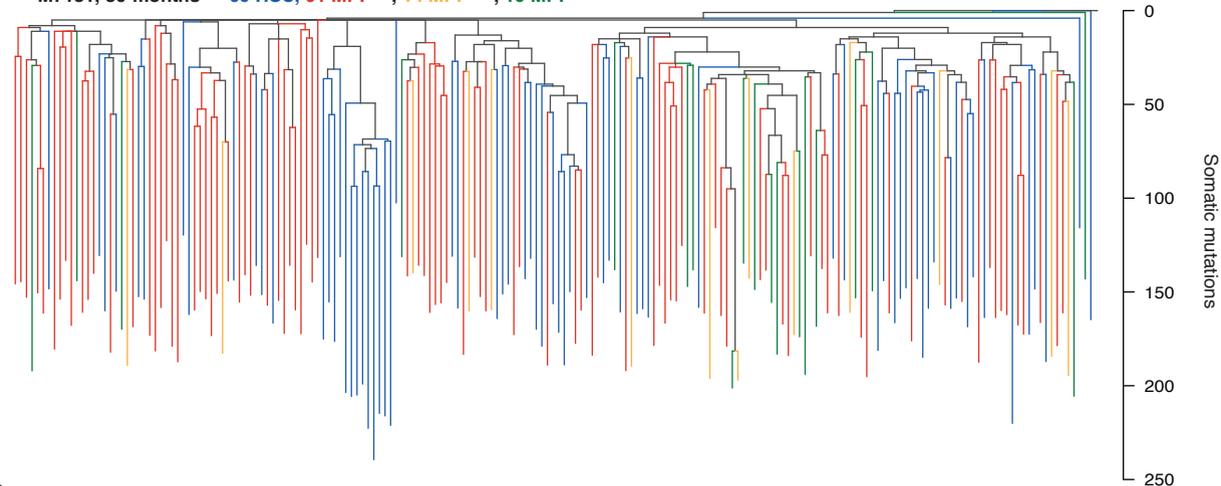
**A**

M7180, 30-months 256 HSC, 270 MPP<sup>All</sup>, 39 MPP<sup>GM</sup>, 33 MPP<sup>Ly</sup>



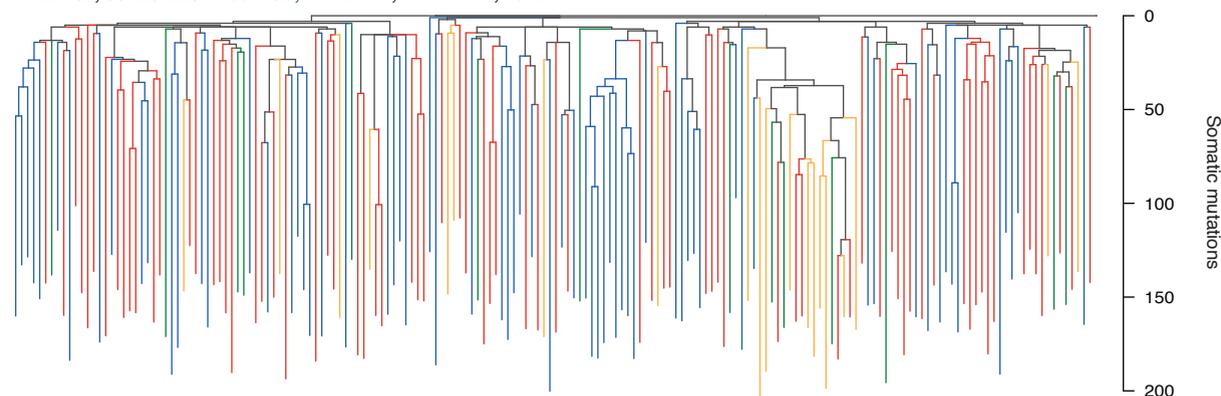
**B**

M7181, 30-months 69 HSC, 91 MPP<sup>All</sup>, 14 MPP<sup>GM</sup>, 19 MPP<sup>Ly</sup>



**C**

M7182, 30-months 69 HSC, 77 MPP<sup>All</sup>, 19 MPP<sup>GM</sup>, 15 MPP<sup>Ly</sup>



**D**

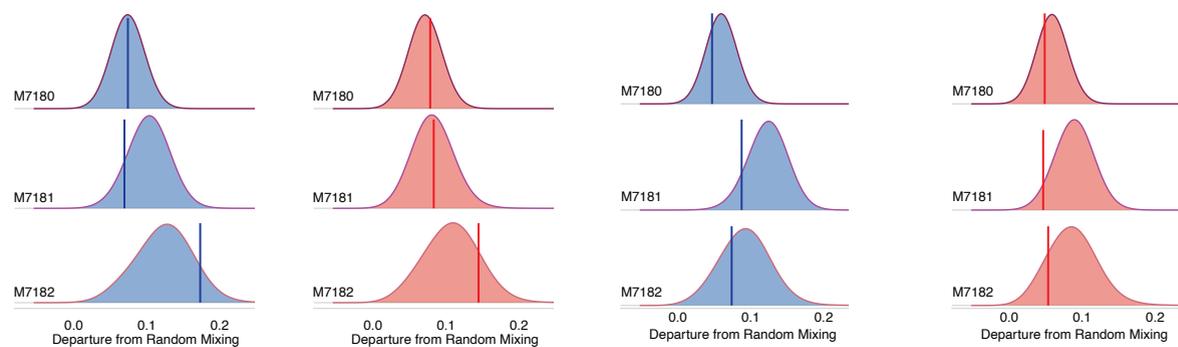
MPP<sup>GM</sup> derived from HSC

MPP<sup>GM</sup> derived from MPP

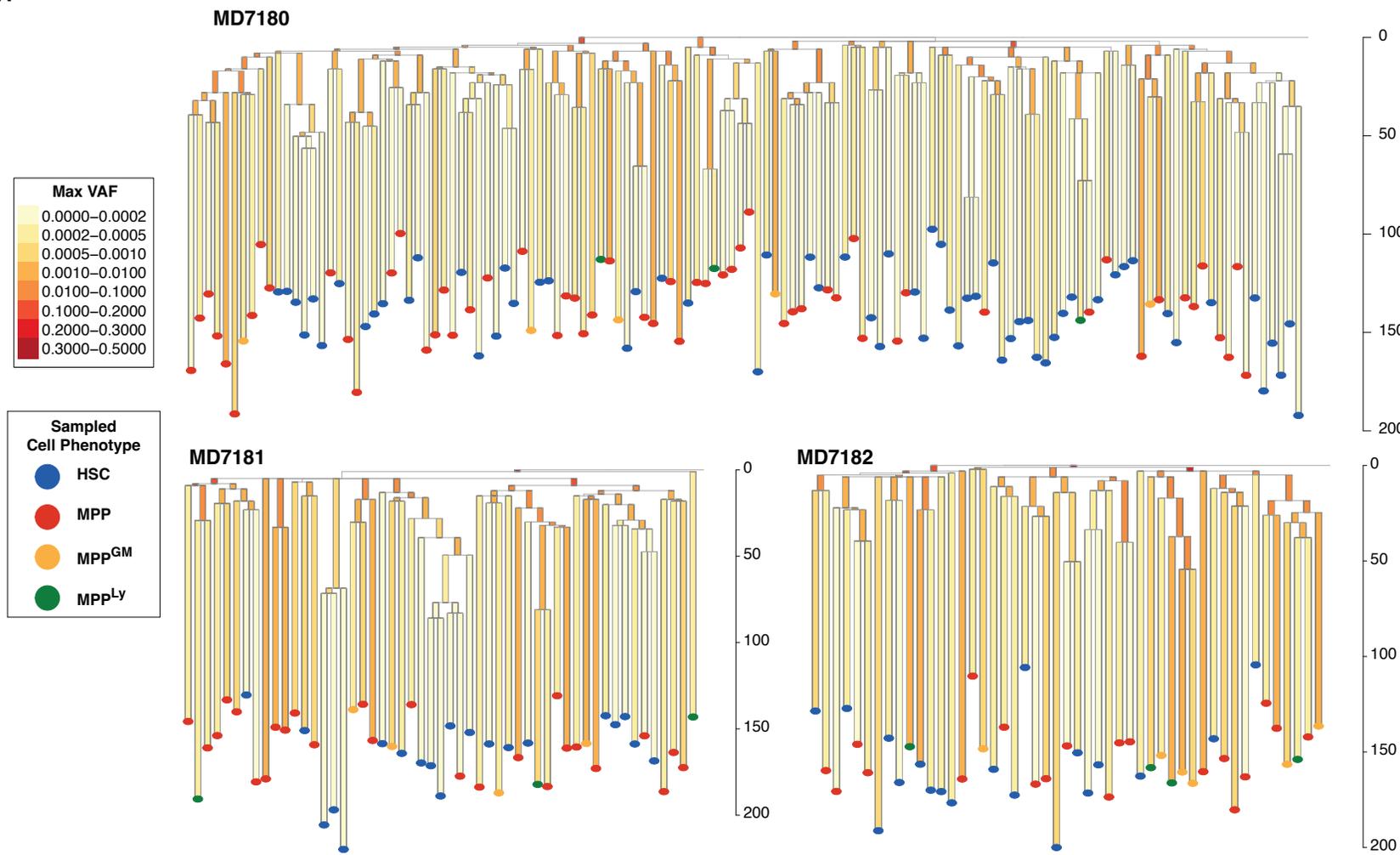
**E**

MPP<sup>Ly</sup> derived from HSC

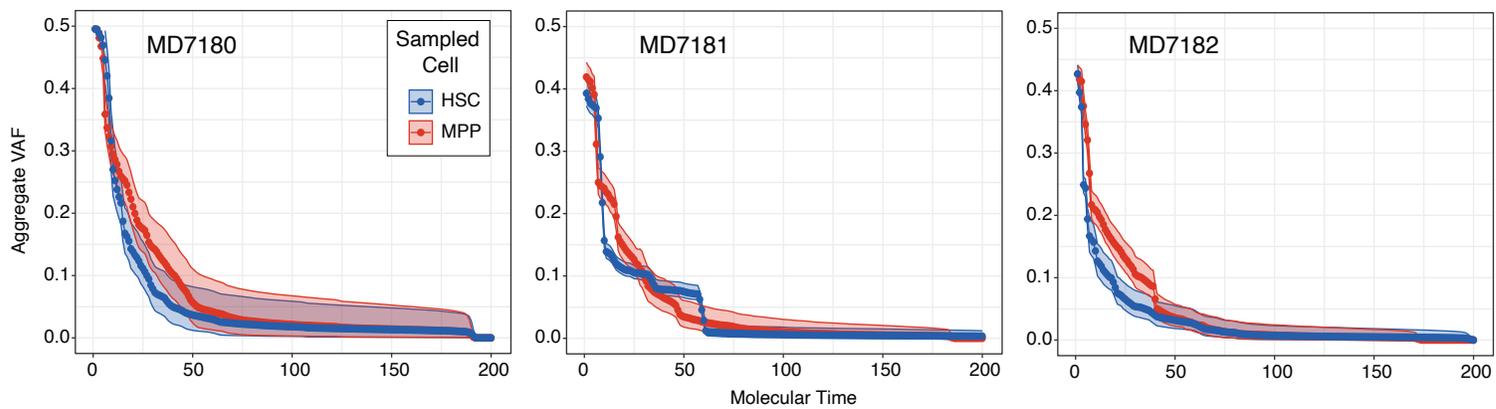
MPP<sup>Ly</sup> derived from MPP



A



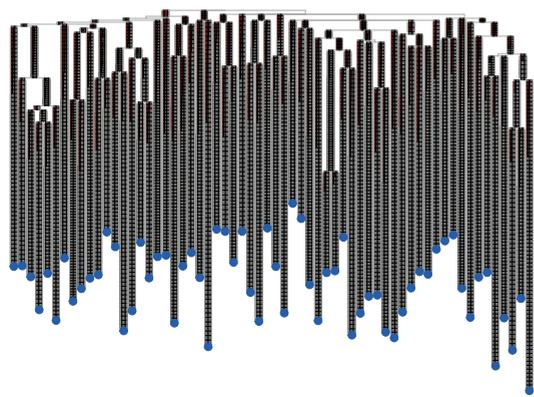
B



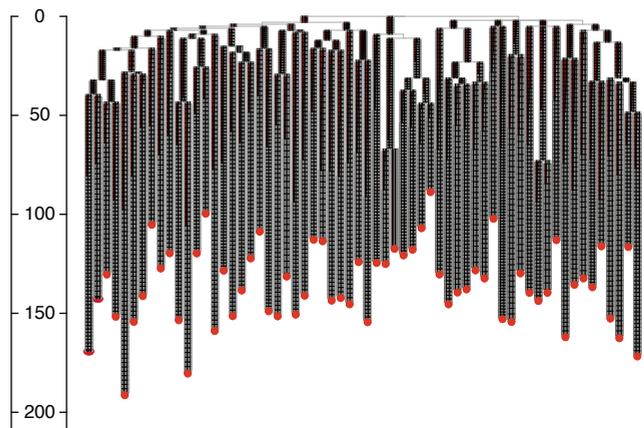
# Extended Data Figure 9

MD7180

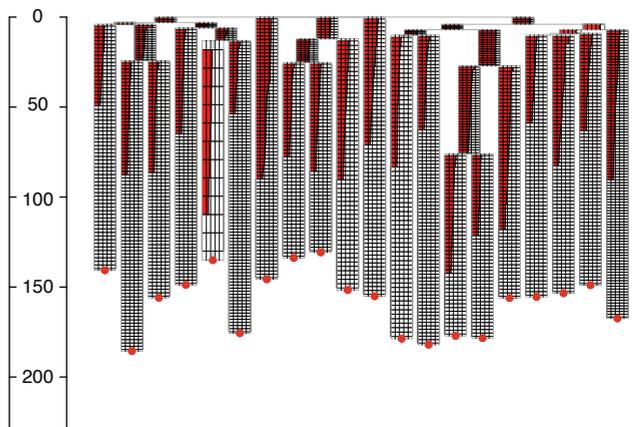
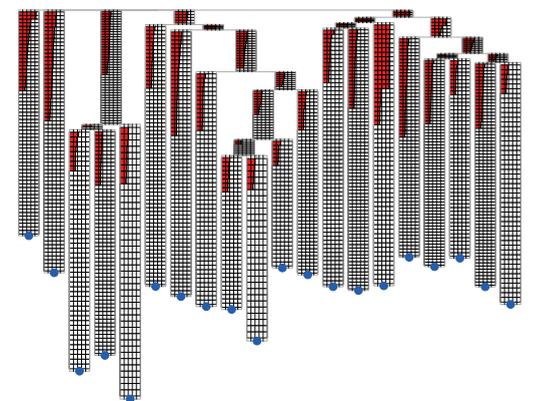
HSC



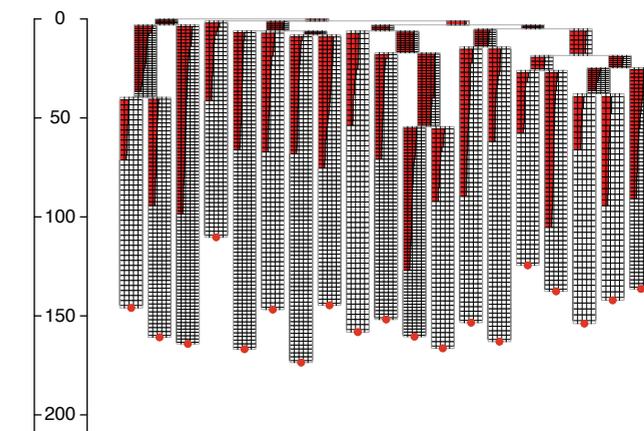
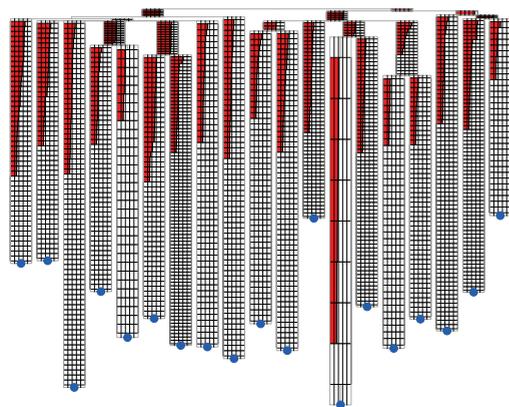
MPP



MD7181



MD7182



## 680 **Acknowledgements**

681 C.D.K. was supported by F30DK131638. The Goodell lab is supported by grants from the National  
682 Institutes of Health, including AG036695, CA183252, CA237291, DK092883, 1P01CA265748,  
683 F30HD111129 (SW), and the Milky Way Research Foundation. J.N. is supported by a Cancer  
684 Research UK Advanced Clinical Fellowship and work in the Nangalia lab is supported by Wellcome,  
685 Cancer Research UK, Alborada Trust, Rosetrees Trust and the MPN Research Foundation. L.J.N. is  
686 supported by NIH AG063543 and AG056278. The Niedernhofer lab is supported by AG063543  
687 and AG063543-02S1. The Harrison lab is supported by 5U01AG022308. DL, MAF, and KYK were  
688 supported by R35 HL155672 (KYK), R01 AI141716 (KYK), F31 HL154661 (DL), F31 HL156500  
689 (MF), and a minority graduate fellowship from the American Society of Hematology (MF). KN is  
690 supported by the Wellcome Trust and CRUK. We are grateful for the assistance of Ryan D. O'Kelly  
691 and Mark Pierson in conducting NME experiments. We thank Elisa Laurenti, Stephen Loughran  
692 and Tony Green for constructive discussions and J. Thomas Gebert and Hilda L. Chan for the  
693 critical feedback.

## 694 **Author contributions**

695 CDK, JN and MAG designed the experiments. JN and MAG supervised the project. CDK performed  
696 cell sorting and *in vitro* culture with support from SW, RA, AM, AG, EM. CDK performed genomic,  
697 phylogenetic, signature, and population dynamics analyses with support from NW, KJD, DL, JF,  
698 EM, PJC, PG, JN. NK developed hidden Markov modelling and KJD performed population dynamic  
699 inferences. AC and KN prepared colonic crypt microdissections. CDK performed mouse  
700 experiments with advice and assistance from MJY, SW, KN, AC, DL, MAF, RA, AM, AG, DH, KYK,  
701 L.J.N. CJW and JRB developed the population genetic analyses of clone sizes and parameter  
702 inferences in Figure 6. CDK, JN and MAG wrote and edited the manuscript. All authors reviewed  
703 and edited the manuscript.

## 704 **Competing interests**

705 The authors declare no competing interests.

## 706 **Database accession**

707 Whole genome sequencing data will be deposited at the European Nucleotide Archive at  
708 accession numbers ERP138320 and ERP144323. Targeted duplex sequencing data will be  
709 deposited at NCBI BioProject PRJNA1033340.

## 710 **METHODS**

### 711 **Cohort**

712 Wild-type C57BL/6 mice were bred at Baylor College of Medicine or received from the Aged Rodent  
713 Colony at the National Institute of Aging (Baltimore, MD). C57BL/6J:FVB/NJ F<sub>1</sub> hybrid mice were  
714 bred in the Niedernhofer laboratory at the University of Minnesota as previously described<sup>77</sup>. HET3  
715 mice were bred at the Jackson Laboratories as described<sup>78</sup>. C57BL/6 were housed at the AALAC-  
716 approved Center for Comparative Medicine in BSL-2 suites. Experimental procedures were  
717 approved by the Baylor College of Medicine or University of Minnesota Institutional Animal Care  
718 and Use Committees and performed following the Office of Laboratory Animal Welfare guidelines  
719 and PHS Policy on Use of Laboratory Animals.

### 720 **Hematopoietic progenitor purification**

721 Whole bone marrow (WBM) cells were isolated from murine hindlimbs and enriched for c-Kit+  
722 hematopoietic progenitors prior to fluorescence-activated cell sorting (FACS) using a BD Aria II.  
723 WBM was incubated with anti-CD117 microbeads (Miltenyi Biotec) for 30 minutes at 4C following  
724 my magnetic column enrichment (LS Columns, Miltenyi Biotec). Progenitor-enriched cells were  
725 stained with an antibody cocktail to identify specific progenitor populations using a recent  
726 consensus definition<sup>55</sup>. LSKs, containing a mixture of stem and progenitor cells, were defined as  
727 Lineage<sup>-</sup>ckit<sup>+</sup>Sca-1<sup>+</sup> (Lineage<sup>-</sup> refers to being negative for expression of a set of lineage-defining  
728 markers indicated below). HSCs were defined as LSK<sup>+</sup>FLT-3<sup>-</sup>CD48<sup>-</sup>CD150<sup>+</sup>; MPPs were defined  
729 as LSK<sup>+</sup>FLT-3<sup>-</sup>CD48<sup>-</sup>CD150<sup>-</sup>. MPP<sup>GM</sup> was defined as LSK<sup>+</sup>FLT-3<sup>-</sup>CD48<sup>+</sup>CD150<sup>-</sup> and MPP<sup>Ly</sup> was  
730 defined as LSK<sup>+</sup>FLT-3<sup>+</sup>CD150<sup>-</sup>. The gating strategy is illustrated in Extended Data Fig.1A. This  
731 immunophenotypic HSC population includes long-term stem cells with serial repopulating ability,  
732 while the MPP population is limited to short-term repopulation, as demonstrated in transplantation  
733 assays<sup>26-31</sup>. For sorting HSCs from newborn pups, the lineage marker Mac1 was excluded because  
734 it is known to be highly expressed on foetal HSCs<sup>79</sup>. Antibodies were c-kit/APC, Sca1/Pe-Cy7, Flt-  
735 3/PE, CD48/FITC, CD150/BV711, Lineage (CD4, CD8, Gr1, Mac1, Ter119)/Pacific Blue and  
736 purchased from BD Biosciences or eBioscience.

### 737 **Single-cell haematopoietic colony expansion *in vitro***

738 Cell sorting was performed on a BD Arial II in two stages. First, HSCs and MPPs were sorted into  
739 separate tubes containing ice-cold FBS using the “yield” sort purity setting to maximise positive  
740 cells. Second, the cell populations from stage one were single-cell index-sorted into individual wells

741 of a 96-well flat bottom tissue culture plates containing 100uL of Methocult M3434 medium (Stem  
742 Cell) supplemented with 1% penicillin-streptomycin (ThermoFisher). No cytokine supplements  
743 were added to the base methylcellulose medium. Cells were incubated at 37°C and 5% CO<sub>2</sub> for  
744 14±2 days, followed by manual assessment of colony growth. Colonies (>200 cells) were  
745 transferred to a fresh 96-well plate, washed once with ice-cold PBS, then centrifuged at 800xg for  
746 10 minutes. Supernatant was removed to 10-15 µL prior to DNA extraction on the fresh pellet. The  
747 Arcturus Picopure DNA Extraction kit (ThermoFisher) was used to purify DNA from individual  
748 colonies according to the manufacturer's instructions. 62-88% HSCs and MPPs produced colonies,  
749 indicating we are sampling from representative populations within each individual compartment  
750 (Extended Data Fig.1B). Extracted DNA from each colony was topped off with 50ul Buffer RLT  
751 (Qiagen) and stored at -80°C.

### 752 **Laser capture microdissection**

753 Matched colonic tissue from the three 30-month-old mice used in this study was dissected and  
754 snap frozen at the time of bone marrow harvest. Colon tissue sectioning and laser capture  
755 microdissection (LCM) was performed as previously described<sup>80</sup>. Briefly, previously snap frozen  
756 colon tissue was fixed in PAXgene FIX (Qiagen) at room temperature for 24 hours and  
757 subsequently transferred into PAXgene Stabilizer for storage until further processing at -20°C. The  
758 fixed tissue was then paraffin-embedded, cut into 10 µm sections, and mounted on PEN-membrane  
759 slides. Staining of histology sections was done using haematoxylin and eosin as previously  
760 described<sup>23</sup>, with scans of each section captured thereafter. Individual colonic crypts were  
761 identified, demarcated, and isolated by LCM using a Leica Microsystems LMD 7000 microscope  
762 (Extended Data Fig.1D) followed by lysis using the Arcturus Picopure DNA Extraction kit  
763 (ThermoFisher).

### 764 **Whole genome sequencing**

765 For low DNA input whole genome sequencing of haematopoietic colonies (from young and aged  
766 mice) and colonic crypts (from aged mice), enzymatic fragmentation-based library preparation was  
767 performed on 1-10 ng of colony DNA, as previously described<sup>80</sup>. Whole genome sequencing (2x150  
768 bp) was performed at a median sequencing depth of 14X for haematopoietic colonies and 17X for  
769 colonic crypts on the NovaSeq platform. Reads were aligned to the GRCm38 mouse reference  
770 genome using bwa-mem. For whole genome single molecule (nanorate) sequencing, we used  
771 matched whole blood genomic DNA collected from the three aged mice during tissue harvest.

772 Nanorate sequencing library preparation was performed as previously described<sup>12</sup>, followed by  
773 sequencing to 146-153X coverage on the Illumina Novaseq platform.

#### 774 **Somatic mutation identification and quality control in haematopoietic colonies**

775 Single nucleotide variants (SNVs) in each colony were identified using CaVEMan<sup>81</sup>, including an  
776 unmatched normal mouse control sample that had previously undergone whole genome  
777 sequencing (MDGRCm38is). Insertions and deletions were identified using cpGPindel<sup>82</sup>. Filters  
778 specific to low-input sequencing artefacts were applied<sup>80</sup>. As variant calling utilised an unmatched  
779 control, both somatic and germline variants were initially called. Germline variants and recurrent  
780 sequencing artefacts were then identified using pooled information across mouse-specific colonies  
781 and filters as follows: i) *Homopolymer run filter*. To reduce artefacts due to mapping errors or  
782 introduced by polymerase slippage, SNVs and indels adjacent to a single nucleotide repeat of  
783 length 5 or more were excluded. ii) *Strand bias filter*. Variants supported by reads only in positive  
784 or negative directions are likely artefacts. For SNVs, a two-sided binomial test was used to assess  
785 if the proportion of forward reads among mutant allele-supporting reads differed from 0.5. Any  
786 variant with significantly uneven mutant read support (cutoff of  $p < 0.001$ ) and with over 80% of  
787 unidirectional mutant reads were excluded. For each indel, if the Pindel call in the originally  
788 supporting colony lacked bidirectional support, the indel was excluded. iii) *Beta binomial filter*.  
789 Variants were filtered based on a beta-binomial distribution across all colonies, as previously  
790 described<sup>23</sup>. The beta-binomial distribution assesses the variance in mutant read support at all  
791 colonies for a given mutation. True somatic variants are expected to be present at high VAF (~0.50)  
792 in some colonies and absent in others, yielding a high beta-binomial overdispersion parameter ( $\rho$ ).  
793 In contrast, artefactual calls are likely to be present at low VAF across many colonies, which  
794 corresponds to low overdispersion. The maximum likelihood estimate of the overdispersion  
795 parameter  $\rho$  was calculated for each loci. For samples with greater than 25 colonies, SNVs with  $\rho$   
796  $< 0.1$  and indels with  $\rho < 0.15$  were discarded. For samples with fewer than 25 colonies, SNVs and  
797 indels with  $\rho < 0.20$  were discarded. iv) *VAF filters*. Variants with VAF significantly lower than the  
798 expected VAF for clonal samples across all mutant genotyped colonies, as assessed with a  
799 binomial test with  $p$  threshold  $< 0.001$ , were discarded. Additionally, variants with VAF less than half  
800 the median VAF of variants that pass the beta-binomial filter were discarded. v) *Germline filter*. All  
801 sites at which the aggregate VAF is not significantly less than 0.45 are assumed to be germline  
802 and discarded. The aggregate VAF is derived from the mutant read count across all colonies for a

803 sample. The binomial test with a confidence threshold  $<0.001$  was used to assess departure from  
804 germline VAF. vi) *Indel proximity filter*. SNVs were discarded if they occurred within 10 base pairs  
805 (bp) of a neighbouring indel. vii) *Missing site filter*. Loci at which genotype information is unavailable  
806 due to poor sequencing coverage will interfere with accurate phylogeny construction. Variants  
807 which have no genotype or coverage less than 6X in over one-third of samples were discarded.  
808 viii) *Clustered site filter*. SNVs and indels within 10 bp of a neighbouring SNV or indel, respectively,  
809 were filtered. ix) *Non-variable site filter*. Sites genotyped as mutant or wildtype in all colonies do  
810 not inform phylogeny relationships and are likely recurrent artefacts or germline variants, thus were  
811 discarded.

812 Some colonies were excluded based on low coverage or evidence of non-clonality or  
813 contamination. Visual inspection of filtered variant VAF distributions per colony was used to identify  
814 colonies with mean variant allele fraction (VAF) $<0.4$  or with evidence of non-clonality (Extended  
815 Data Fig.1C).

#### 816 **Mutation burden estimation**

817 Total SNV burden from WGS of individual colonies was corrected for differing depths of sequencing  
818 using a per-sample asymptomatic regression fit<sup>80</sup> (Extended Data Fig.1E). A linear mixed effect  
819 model was used to estimate the rate of mutation acquisition with age, taking into account individual  
820 animals as a random effect as follows:  $burden \sim age + (0 + age | sampleID)$ .

821 We filtered nanorate sequencing calls as previously described<sup>12</sup>, with the following modifications:  
822 we excluded variants (i) mapped to the mitochondrial genome, (ii) located within 15 bp of  
823 sequencing read ends, or (iii) observed in all duplex consensus reads as these are likely germline  
824 events. Matched colony whole genome sequencing data was used as a normal control. Mutation  
825 burdens were normalised to diploid genome size to determine the global SNV burdens.

#### 826 **Phylogeny construction and quality control**

827 Phylogenetic trees were constructed based on shared mutations between colonies for each mouse,  
828 as extensively described previously<sup>9,34</sup>. The steps, in brief, were as follows: (i) *Create genotype*  
829 *matrix*. Every colony has high sequencing coverage (median 14X) distributed evenly across the  
830 genome, allowing the determination of a genotype for nearly every mutated site observed across  
831 colonies. Each locus was annotated as Present, Absent, or Unknown in a read depth-specific  
832 manner. The number of unattributable sites was low, allowing precise inferences of colony  
833 interrelatedness. (ii) *Infer phylogenetic tree from genotype matrix*. We applied the maximum

834 parsimony algorithm MPBoot to construct phylogenetic trees from the genotype matrix. Only SNVs  
835 were used to infer tree topology, but both SNVs and indels (if any) were assigned to inferred  
836 branches using *treemut*. Loci with unknown genotypes in at least one-third of colonies were  
837 annotated as missing sites and not used in phylogeny inference. (iii) *Normalise branch lengths for*  
838 *differing sequencing depth and sensitivity*. Branch lengths at this stage are defined by the number  
839 of mutations supporting each branch (molecular time). However, each colony has slightly different  
840 sequencing coverage, which correlates with differences in mutation detection sensitivity. Thus, we  
841 normalised branch lengths based on genome coverage to correct for sensitivity differences across  
842 colonies with varying depth, as described in ref. <sup>34</sup> (Extended Data Fig.1E). (iv) *Annotate trees with*  
843 *phenotype and genotype information*. Each terminal branch (tip) of a tree represents a specific  
844 colony. Thus, we annotated each branch of the tree with the sampled cell phenotype (HSC versus  
845 MPP).

846 Tree-level checks were used to identify any discordant branch assignments and assess the validity  
847 of tree topology. Any branches supported by variants with mean VAF <0.4 likely contained  
848 contamination by non-clonal variants and suggested the filtering strategy (see above) was  
849 insufficient. Similarly, the branch-level VAF distributions of every colony (tip) in the tree were  
850 manually inspected to confirm supporting variants were not present in unrelated portions of the tree  
851 (topology discordance). Finally, the trinucleotide spectra of individual somatic mutations were  
852 compared between those mutations located on shared branches (that is, mutations supported by  
853 >2 colonies) and mutations only observed once, and thus present on terminal branches. Mutation  
854 spectra were highly similar, indicating that mutations not shared by more than one colony were not  
855 populated by a relative excess of artefacts (Extended Data Fig.1F).

## 856 **Population size trajectories**

857 We use the *phylodyn* package, which uses the density of coalescent events (bifurcations) in a  
858 phylogenetic tree to estimate the trajectory of  $N(t)/\lambda(t)$  over time<sup>9,10</sup>. Ultrametric lifespan-scaled  
859 trees were used to infer chronological timing. Under a neutral model of population dynamics, the  
860 phylogeny of a sample is a realisation of the coalescent process. In the coalescent process, the  
861 rate of coalescent events at time  $t$  is proportional to the ratio of population size,  $N(t)$ , to the birth  
862 rate,  $\lambda(t)$  (which in the context of stem cell dynamics is the symmetric cell division rate). The  
863 sequence of inter-coalescent intervals across any time interval  $[t_1, t_2]$  is informative about the  
864 value of the parameter ratio  $N(t)/\lambda(t)$  across the same time interval. We note that only with a

865 constant cell division rate  $\lambda$  over time can the trajectory parameter be interpreted as a scalar  
866 multiple of the trajectory of population size  $N(t)$ . *Phylodyn* assumes isochronous sampling and a  
867 neutrally evolving population. We overlaid separate population size trajectories for HSCs and MPPs  
868 in Figure 3A.

### 869 **Approximate Bayesian computation**

870 We used inference from phylodynamic trajectories to inform the development of an HSC population  
871 dynamics model. Population size trajectories from *phylodyn* indicated two successive ‘epochs’ of  
872 exponential growth, with some variation in growth rate between epochs, and a steady increase in  
873 population size over time (Fig.3A). Given the constraint of tissue volumes, it may be implausible  
874 that the HSC population grows constantly. We reconcile this discrepancy by noting that there are  
875 very few late-in-life coalescences in our phylogenies, and, as a consequence, the estimated  
876 *phylodyn* trajectory in late adulthood is associated with very wide credible intervals. We employed  
877 a population growth model based on a linear birth-death process<sup>83</sup> (in which a population tends to  
878 grow exponentially, subject to stochastic fluctuations), together with a fixed upper bound  $N$  on  
879 population size. The model assumed a constant birth rate  $\lambda$  and constant death rate  $\nu$ , with the  
880 population trajectory growing at a rate  $\lambda-\nu$  within an epoch. The shape of the trajectory of  $N(t)/\lambda(t)$   
881 depends on the cell division rate parameter  $\lambda$ , not only through the denominator in the ratio  
882  $N(t)/\lambda(t)$ , but also on  $\lambda-\nu$ , through the tendency of the population size  $N(t)$  to grow exponentially  
883 at a rate  $\lambda-\nu$ . In particular, if we increase the fixed upper limit  $N$ , and at the same time increase the  
884 cell division rate  $\lambda$ , so that their ratio remains constant, the shape of the trajectory of  $N(t)/\lambda(t)$  will  
885 change as a consequence of the changes in the value of the parameter  $\lambda$ . This suggests that the  
886 parameters  $\lambda$ ,  $\nu$ , (in each epoch), and  $N$ , are all identifiable, and so can be estimated separately.  
887 The identifiability of  $\lambda$ ,  $\nu$ , and  $N$  are expanded upon in Supplementary Note 4.

888 We applied Bayesian inference procedures<sup>41</sup> to estimate the parameters ( $\lambda$ ,  $\nu$ , and  $N$ ) of the  
889 bounded birth-death process. We used Approximate Bayesian Computation (ABC). This method  
890 first generates simulations of population trajectories and (sample) phylogenetic trees across a  
891 lifespan. Each population simulation is run with specific values for the population dynamic  
892 parameters drawn from a prior distribution over biologically plausible ranges of parameter values.  
893 The ABC method includes a rejection step that retains only those parameter values which  
894 generated simulated phylogenies resembling the observed phylogeny (as measured by an  
895 appropriate Euclidean distance). The accepted simulations constitute a sample from the

896 (approximate) posterior distribution. Population trajectories and sample phylogenies were  
897 simulated using the *rsimpop* R package. Approximate posterior distributions were computed using  
898 the R package, *abc*. We specified uniform joint prior densities for  $\lambda$ ,  $\nu$ , and N which encompassed  
899 published estimates for N (population size) and  $\lambda$  (symmetric division rate)<sup>31,73–75,84</sup>. N ranged from  
900  $10^2$  to  $10^5$  cells,  $\lambda$  ranged from 0.01 to 0.15 cell division per day, and  $\nu$  ranged from 0 to  $\lambda$ , such  
901 that the growth rate ( $\lambda - \nu$ ) is always positive (as observed in the *phylodyn* trajectories).

902 Our population dynamics model was a birth-death process incorporating two separate growth  
903 epochs. The first (early) epoch lasted until 10 weeks post-conception, and the second (later) epoch  
904 lasted from 10 weeks onwards and corresponded to murine adulthood. Inferences were weak for  
905 the early epoch; thus, the later epoch was used for parameter inferences. Posterior densities from  
906 the three older mice were computed using the ‘rejection’ method (Extended Data Fig.6) and pooled  
907 to yield parameter estimates and credible intervals.

### 908 **Early life polytomy analysis**

909 The polytomies were used to estimate lower and upper bounds for the mutation rate per symmetric  
910 division during embryogenesis. The method detailed in Lee-Six *et al.*<sup>10</sup> was used, whereby the  
911 number of edges with zero mutation counts at the top of the tree (up to the first 12 mutations) is  
912 inferred from the number and degree of polytomies assuming an underlying tree with binary  
913 bifurcations. The mutations per division are assumed to be Poisson distributed. A maximum  
914 likelihood range is then calculated in two steps: first, using the 95% confidence interval of the  
915 proportion of zero length edges, with this next leading to a maximum likelihood estimate for the  
916 Poisson rate. Sample M7183 lacked sufficient early life diversity (<10 unique lineages within 12  
917 mutations molecular time) and thus was excluded.

### 918 **Shared variants between blood and colonic crypts**

919 Mutation genotype matrices (described above) were generated for colonic crypt samples at loci  
920 observed in truncal (shared) branches in the matched HSC tree. Every variant was annotated as  
921 present or absent for each colonic crypt. We applied two stages to crypt annotation. First, a crypt  
922 sample was marked positive if the given variant exceeded a per-sample minimum VAF threshold.  
923 The minimum VAF threshold was defined as half the median VAF for all pass-filter colonic crypt  
924 variants (as described above). Next, for each variant represented in at least one crypt, any  
925 remaining crypt with >2 mutant allele read support was marked positive. This tiered definition

926 allowed for shared variant capture despite differences in coverage among crypt samples. The  
927 proportion of a shared variant present among crypt samples was illustrated as a pie chart and  
928 annotated to the respective branch of the matched HSC tree (Extended Data Fig.3).

### 929 **Mutational signature analysis**

930 We used the Hierarchical Dirichlet Process (HDP) algorithm to extract mutation signatures across  
931 aged and young HSC and MPP colony samples, following the process detailed in ref <sup>85</sup>. Prior work  
932 in humans has applied mutation signature extraction to SNVs found only on terminal branches of  
933 phylogenetic trees – such terminal branches displayed mutation burdens in excess of 1000  
934 mutations, depending on the organ. Given the low mutation burden in mouse hematopoietic  
935 colonies (terminal branch lengths spanning 30-150 mutations), and thus reduced mutational  
936 information, we utilised all branches with length  $\geq 30$  mutations as input. To circumvent any bias  
937 against shared variants, branches with less than 30 SNVs were collapsed to a single 'shared  
938 branch' sample. We generated mutation count matrices for each branch, using the 96 possible  
939 trinucleotide mutational contexts as input to the R package *hdp*. HDP was run i) without priors (de  
940 novo), ii) with the reference catalogue of all 79 signatures derived from the PanCancer Analysis of  
941 Whole Genomes study (COSMIC version 3.3.1) as priors, or iii) with the signatures previously  
942 defined as active in mouse colon<sup>23</sup>, SBS1, SBS5, SBS18, as priors. Trinucleotide signature  
943 definitions were adjusted to mouse genome mutation opportunities before usage as priors, and all  
944 prior signatures were weighted equally. Signature extraction parameters i) and ii) produced profiles  
945 that did not resemble any existing signatures (cosine similarity  $< 0.9$ ), likely due to relatively limited  
946 SNV burden in mouse colony data. Usage of mouse colon signatures as prior information (iii)  
947 yielded four signature components. Two signature components demonstrated high similarity to  
948 SBS1 and SBS5 (cosine  $> 0.9$ ). The remaining two unknown components were deconvoluted to  
949 reattribute their composition to known signatures using the *fit\_signatures* function from *sigfit*. This  
950 yielded three components with a reconstruction cosine similarity metric exceeding 0.99 for similarity  
951 to SBS1, SBS5, and SBS18, indicating these three signatures explain the majority of our data  
952 (Extended Data Fig.4A). We surmise the final reattribution step was necessary because of the log-  
953 fold lower SNV burdens in mouse blood colonies (30-200 mutations) relative to other tissues  
954 examined in previous work ( $> 1000$  mutations).

## 955 **Branch signatures assignment and analyses**

956 For each mouse, we pooled the assigned SNVs into a “private” or “shared” category depending on  
 957 whether the variant maps to a shared branch or not. Signature attribution to signatures SBS1,  
 958 SBS5, and SBS18 was then carried out for each of these per mouse category using  
 959 *sigfit:fit\_to\_signature* with the default “multinomial” model. The per-branch attributions were then  
 960 carried out by 1) assigning a per-mutation signature membership probability and then 2) summing  
 961 these signature membership probabilities over all SNVs assigned to a branch to obtain a branch-  
 962 level signature attribution proportion. The per-mutation signature probability was calculated using:

$$963 \quad P(\text{mutation} \in \text{Sig}) = \frac{P(\text{mutation} \in \text{Sig})P_0(\text{mutation} \in \text{Sig})}{\sum_{\text{Sig}' \in \{\text{SBS1}, \text{SBS5}, \text{SBS18}\}} P(\text{mutation} \in \text{Sig}')P_0(\text{mutation} \in \text{Sig}')}$$

964  
 965 Where the prior probability,  $P_0(\text{mutation} \in \text{Sig})$ , is given by the mean Sigfit attribution probability  
 966 of the specified signature, *Sig*, for the category that the mutation belongs to.

967 A linear mixed effect model was used to assess the relationship between age and the signature-  
 968 specific substitution burden for each colony while accounting for repeated measures. The  
 969 signature-specific burdens per colony were estimated using a linear mixed model (R package *lme4*)  
 970 with age as a random effect and mouse ID as grouping variable:

$$971 \quad \text{burden}_{\text{signature}} \sim \text{age} + (0 + \text{age} | \text{mouseID}).$$

## 972 **Hidden Markov tree approach**

973 *Modelling the ancestral unobserved MPP and HSC states with a hidden Markov tree:* We defined  
 974 three unobservable (“hidden”) ancestral states, embryonic precursor cell (EMB), HSC and MPP,  
 975 and used the observed outcomes (HSC or MPP tip states) to infer the transition probabilities  
 976 between these identities and the most likely sequence of cell identity transitions during life. The  
 977 transitions between states are modelled by a discrete time Markov chain with one step in time  
 978 representing one mutation in molecular time. We require the root of the tree, presumably the  
 979 zygote, to start in the “EMB” state and to stay in that state until 10 mutations in molecular time.  
 980 After 10 mutations the cell then has a non-zero probability of transitioning to another state given by  
 981 the transition transition matrix **M**:

$$982 \quad \mathbf{M} = \begin{pmatrix} 1 - p_{\text{HSC} \rightarrow \text{MPP}} & p_{\text{HSC} \rightarrow \text{MPP}} & 0 \\ p_{\text{MPP} \rightarrow \text{HSC}} & 1 - p_{\text{MPP} \rightarrow \text{HSC}} & 0 \\ p_{\text{EMB} \rightarrow \text{HSC}} & p_{\text{EMB} \rightarrow \text{MPP}} & 1 - p_{\text{EMB} \rightarrow \text{HSC}} - p_{\text{EMB} \rightarrow \text{MPP}} \end{pmatrix}$$

983 This then implies the following transition probabilities for branch  $u$ , having length  $l(u)$  (excluding  
984 any overlap with molecular time less than 10 mutations), starting in state  $i$  and ending in state  $k$  :

985 
$$P_{i,k}(u) = (\mathbf{M}^{l(u)})_{i,k}$$

986 Now for a node that is in a specified state, the probability of descendent states is independent of  
987 the rest of the tree. This conditional independence property facilitates recursive calculation of a  
988 best path (“Viterbi path”), the likelihood of the Viterbi path, and the full likelihood of the observed  
989 phenotypes given the model. The approach is essentially an inhomogeneous special case of the  
990 approach previously described<sup>86</sup>.

991 *Upward algorithm for determining likelihood of the observed states given  $\mathbf{M}$  and a prior probability*  
992 *of root state  $\boldsymbol{\pi}$ :* The probability of the observed data descendant from a node  $u$  whose end of branch  
993 state is  $i$  is given by:

994 
$$P_u(D_u|i) = \prod_{v \in \text{children}(u)} \left( \sum_{k=1}^S P_{i,k}(v) P_v(D_v|k) \right)$$

995 where  $S$  is the number of hidden states ( $S = 3$  in our usage), and  $D_u$  denotes the observed data  
996 descendant of  $u$ , that is, the observed tip phenotypes of the clade defined by  $u$ .

997 *Initialisation of terminal branches:* The probability of observing a matching phenotype is assumed  
998 to be:

999 
$$P_u(\text{Observed Phenotype of } u = i|i) = 1 - \epsilon$$

1000 The probability of observing a mismatching phenotype,  $j \neq i$ , is:

1001 
$$P_u(\text{Observed Phenotype of } u = j|i) = 0.5\epsilon$$

1002 The root probability  $P_{root}(D_{root}|i)$  is calculated recursively from the above and the model likelihood  
1003 is given by:

1004 
$$P = \sum_{i=1}^S \pi_i P_{root}(D_{root}|i)$$

1005 Given the two-stage cell sorting approach described above, we assume nearly error-free  
1006 phenotyping and set  $\epsilon = 10^{-12}$ .

1007 *Determining the most likely sequence of hidden end-of-branch states:* This Viterbi-like algorithm  
1008 can be run in conjunction with the upward algorithm. Here, instead of summing over all possible  
1009 states, we keep track of the most likely descendant states for each possible state of the current  
1010 node  $u$ .

1011 The quantity  $\delta_u(i)$  is the probability of the most likely sequence of descendant states given that  
1012 node  $u$  ends in state  $i$ :

$$1013 \quad \delta_u(i) = \prod_{v \in \text{children}(u)} \left( \max_k \{ \delta_v(k) P_{i,k}(v) \} \right)$$

1014 Additionally, for each node we store the most probable child states given that  $u$  is in state  $i$ :

$$1015 \quad \Psi_{u,v}(i) = \operatorname{argmax}_k \{ \delta_v(k) P_{i,k}(v) \}$$

1016 The tip deltas are initialised using the emission probabilities:

$$1017 \quad \delta_u(i) = P_u(\text{Observed Phenotype of } u|i)$$

1018 The above provides a recipe for recursively finding  $\delta_{root}(i)$  and is combined with prior root  
1019 probability  $\pi$  to give the most likely root state,  $\max_k \{ \delta_{root}(i) \}$ , in our case we set the prior probability  
1020 of “EMB” to unity - so EMB is the starting state. The child node states are then directly populated  
1021 using  $\Psi$ .

## 1022 Targeted duplex-consensus sequencing

1023 Genomic DNA from freshly collected peripheral blood was purified using the Zymo Quick-DNA  
1024 Miniprep Plus kit according to the manufacturer's instructions. 1650 ng of high-molecular-weight  
1025 DNA was ultrasonically sheared to an average 300 bp fragment size using a Covaris M220 and  
1026 ligated to duplex identifier sequencing adapters<sup>87</sup> using the Twinstrand Biosciences DuplexSeq  
1027 library prep kit. A large input of gDNA was used to ensure that the number of input genomic  
1028 equivalents (about 275,000-330,000 genomes) did not limit the achievable duplex sensitivity. A  
1029 custom baitset of biotinylated probes was used to enrich sequences targeting mouse orthologues  
1030 of common human CH driver genes over two overnight hybridisation reactions. Our target panel  
1031 spanned 61.8 kb and captured homologous regions from the entire coding region of the following  
1032 genes: *Dnmt3a*, *Tet2*, *Asxl1*, *Trp53*, *Rad21*, *Cux1*, *Runx1*, *Bcor*, and *Bcorl1*, and specific exons  
1033 with hotspot mutations (as observed in COSMIC) for the following genes: *Ppm1d*, *Sf3b1*, *Srsf2*,  
1034 *U2af1*, *Zrsr2*, *Idh1*, *Idh2*, *Gnas*, *Gnb1*, *Cbl*, *Jak2*, *Ptpn11*, *Brcc3*, *Nras*, and *Kras*. Targeted loci  
1035 encompass >95% of human CH events<sup>43</sup> and are described in Supplementary File 2. Libraries were  
1036 sequenced on the NovaSeq platform to a raw depth between 1-3 million reads, corresponding to  
1037 duplex-consensus depths between 30,000-50,000X that vary across targeted exons  
1038 (Supplementary Note 3). Quality control of duplex sequencing is discussed in Supplementary Note  
1039 3.

## 1040 **Variant identification in targeted gene duplex-consensus sequencing**

1041 Duplex-consensus and single-strand consensus reads were generated using the *fgbio* suite of tools  
1042 according the *fgbio* Best Practices FASTQ to Consensus Pipeline Guidelines  
1043 (<https://github.com/fulcrumgenomics/fgbio/blob/main/docs/best-practice-consensus-pipeline.md>).  
1044 To build a duplex-consensus read, we required at least 3 reads in each supporting read family (i.e.,  
1045 at least 3 sequenced PCR duplicates of matched top and bottom strands from an original DNA  
1046 molecule). The ‘DuplexSeq Fastq to VCF’ (version 3.19.1) workflow hosted on DNANexus was also  
1047 used to generate duplex-consensus reads. Next, VarDict<sup>88</sup> was used to identify all putative variants,  
1048 followed by functional annotation using Ensembl Variant Effect Predictor<sup>89</sup>. Finally, numerous post-  
1049 processing filters were applied to remove false positives and artefactual variants: (i) *Quality flag*  
1050 *filter*. VarDict annotates all variants using a series of quality flags that assess mapping and read-  
1051 level fidelity<sup>88</sup>. Any variant with a quality flag other than “PASS” was discarded. (ii) *Read support*  
1052 *filter*. Duplex sequencing enables detection of somatic variants even from a single read<sup>87</sup>; however,  
1053 variants supported by a consensus read (singletons) were found to be highly enriched for spurious  
1054 calls. Thus, any variant supported by a single read was discarded. (iii) *Mismatches per read filter*.  
1055 Variants were excluded if the mean number of mismatches per supporting read exceeded 3.0. (iv)  
1056 *End Repair & A-tailing artefact filter*. Library preparation enzymatic steps may introduce false  
1057 positive SNVs near read ends due to misincorporation of adenine bases during A-tailing or  
1058 mistemplating during blunting of fragmented 3’ ends. The *fgbio* FilterSomaticVcf tool was used to  
1059 assess the probability that any variant within 20 bp of read ends was due to such enzymatic errors;  
1060 probable end-repair artefacts were discarded. (v) *Read position filter*. Variants in positions  $\leq 15$  bp  
1061 from the 5’ or 3’ end of a consensus read were observed to be enriched for spurious variants based  
1062 on trinucleotide signature and were discarded. (vi) *Oxidative damage filter*. Mechanical  
1063 fragmentation (prior to duplex adapter attachment) creates oxidative DNA damage, often in the  
1064 form of 8-oxoguanine<sup>90,91</sup>, which mis-pairs with thymine and is fixed after PCR amplification. Any  
1065 variant fitting the previously described oxidative artefact signature (SBS45) were discarded. (vii)  
1066 *Sequencing coverage filter*. Variants at loci with duplex depth of  $\leq 20,000X$  were considered under-  
1067 sequenced and discarded. (viii) *Strand bias filter*. We employed a Fisher’s exact test to assess for  
1068 forward or reverse strand bias between wildtype and mutant reads. Any variant enriched for  
1069 unidirectional read support was discarded. (ix) *Recurrent variant filter*. Variants present in  $\geq 5\%$  of  
1070 samples per duplex-sequencing batch or in  $\geq 5$  independent samples were discarded. (x) *Indel*  
1071 *length filter*. Long insertions or deletions could be attributed to poor mapping, erroneous fragment

1072 ligation, or false positive calls by VarDict. Any indels  $\geq 15$  bp were excluded. (xi) *High VAF filter*.  
1073 Germline variants display a VAF of 0.5 or 1.0. Any variants with  $VAF \geq 0.4$  were excluded as  
1074 putative germline variants. (xii) *Impact filter*. CH is driven by functional coding sequence changes  
1075 in driver genes. Thus, synonymous mutations were excluded during generation of the dot-plots in  
1076 Figures 4-5. This filter was not utilised for analyses that require synonymous variant information  
1077 (dN/dS, fitness effect estimation). (xiii) *Homologous position filter*. Residues conserved with  
1078 humans are likely to be functional in mice. Variants at loci without a matching reference allele at  
1079 homologous position in humans were discarded. This filter primarily eliminated intronic variants and  
1080 was not utilised for analyses incorporating synonymous variant information. Variants identified are  
1081 detailed in Supplementary File 2.

### 1082 **Murine perturbation experiments**

1083 Perturbation experiments were initiated in aged (21-month) male and female mice unless otherwise  
1084 described. Mice were randomly allocated to control or experimental groups. Investigators were not  
1085 blinded to the group assignment during experiments. For *Mycobacterium avium* infection, mice  
1086 were infected with  $2 \times 10^6$  colony-forming units of *M. avium* delivered intravenously as previously  
1087 described<sup>92</sup>. Mice were infected once every 8 weeks (twice in total) to ensure chronic infection. For  
1088 cisplatin exposure, mice were exposed to 3 mg/kg cisplatin delivered intraperitoneally every four  
1089 weeks, as indicated. Dose spacing was selected to allow for sufficient recovery following  
1090 myeloablation and blood counts were not altered in cisplatin-treated mice (Fig.4C), indicating  
1091 recovery of haematopoiesis. For 5-Fluorouracil exposure, 150 mg/kg 5-FU was delivered  
1092 intraperitoneally every four weeks two times; this 5-FU dose has previously been shown to drive  
1093 temporary activation of HSCs in mice<sup>93,94</sup>. Exposure to a normalised microbial experience (NME)  
1094 of murine transmissible pathogens was performed as previously described<sup>51</sup>. Briefly, immune-  
1095 experienced “pet store” mice were purchased from pet stores around Minneapolis, MN. Aged (24-  
1096 month) C57BL/6J:FVB/NJ laboratory mice were either directly cohoused with pet store mice or on  
1097 soiled (fomite) bedding collected from cages of pet store mice. Mice were exposed to continuous  
1098 fomite bedding for 1 month, followed by 5 months recovery on SPF bedding before tissue collection.  
1099 All NME work was performed in the Dirty Mouse Colony Core Facility at the University of Minnesota,  
1100 a BSL-3 facility. Age-matched C57BL/6J:FVB/NJ F1 laboratory mice maintained in specific  
1101 pathogen free (SPF) conditions were used as controls. For monitoring, peripheral blood (~50uL)  
1102 was collected in EDTA-coated tubes and analysed on an OX-360 automated hemocytometer (Balio

1103 Diagnostics). For all aforementioned mouse cohorts, peripheral blood genomic DNA was purified  
1104 and converted to duplex sequencing libraries as described above.

1105 Differences in clone burden between control and treated cohorts was quantified using a Mann-  
1106 Whitney test on cumulative VAFs per sample. Gene-level enrichment was measured using a  
1107 Fisher's exact test on the number and mutant and wildtype reads, normalised for coverage  
1108 differences between samples. Gene-level dN/dS estimates were generated as described below.

### 1109 **dN/dS analysis**

1110 The ratio of nonsynonymous to synonymous mutation rates (dN/dS) can be used to assess for  
1111 selection within somatic mutations by comparing the observed dN/dS to that expected under  
1112 neutral selection. We use the R package *dNdScv*<sup>95</sup> to estimate dN/dS ratios of somatic mutations  
1113 derived from whole genome and targeted gene duplex-consensus sequencing. To incorporate  
1114 mouse-specific differences in trinucleotide context composition and background mutation rates, we  
1115 generated a murine reference CDS dataset using the *buildref* function and genome annotations in  
1116 Ensembl (version 102). For the phylogenetic trees, we input all tree variants to the *dndscv*  
1117 function. dN/dS output and all coding variants detected in trees are listed in Supplementary File 1. To  
1118 examine dN/dS in targeted duplex-consensus sequencing data, we pooled all variants observed in  
1119 cross-sectionally sampled mice across ages (Fig.3A) and ran *dndscv* limited to exons only included  
1120 on our targeted panel (Supplementary File 2).

### 1121 **Targeted capture of tree variants**

1122 We designed a custom targeted DNA baitset (Agilent SureSelect) targeting mutations on the  
1123 phylogenetic trees of the aged mice, and then queried genomic DNA purified from matched  
1124 peripheral blood for tree-specific mutations using high-depth targeted sequencing. The baitset was  
1125 designed to capture mutations on the phylogenetic trees of all 3 aged mice (MD7180, MD7181,  
1126 and MD7182), and to cover mutations found in HSCs, MPPs and LSKs. The baitset was designed  
1127 as follows: (i) All variants on shared branches that pass the SureDesign tool's "moderately stringent  
1128 filters". (ii) All variants on a random subset of private branches that pass SureDesign's "most  
1129 stringent filters". Approximately 25% of the private branches of each mouse were randomly  
1130 selected. (iii) The exons and 3' and 5' UTRs for all CH driver genes used in our duplex sequencing  
1131 panel (listed above). Target-enriched libraries were generated according to the manufacturer's  
1132 protocol and sequenced using the Illumina Novaseq platform. Baits were sequenced to median  
1133 depths of 2616X, 2549X and 2628X for MD1780, MD7181 and MD7182 respectively.

1134 To quantify the degree of HSC and MPP contribution to peripheral blood, we estimated the posterior  
1135 distribution of true VAF for every mutation captured with our targeted baitset. This was done using  
1136 the Gibbs sampling method previously developed<sup>96</sup>. Then, for each molecular time  $t$ , and for each  
1137 branch that overlaps  $t$ , we estimate the VAF of a hypothetical mutation at time  $t$ . This is done by  
1138 arranging our baitset variants in descending estimated VAF order at equally spaced intervals down  
1139 the branch and then linearly interpolating the VAF at time  $t$  based on the estimated VAF of the  
1140 neighbouring mutations. The aggregate VAF at time  $t$  for a tree or lineage is then calculated as the  
1141 sum of the estimated VAFs of the overlapping branches at time  $t$ .

## 1142 **Maximum likelihood estimates of fitness effects**

### 1143 *Evolutionary framework*

1144 To generate estimates of fitness effects, mutation rates, and population size, we applied an  
1145 evolutionary framework based on continuous time branching for HSCs, as previously reported<sup>53</sup>.  
1146 The framework is based on a stochastic branching model of HSC dynamics, where variants with a  
1147 variant-specific fitness effect,  $s$ , are acquired stochastically at a constant rate  $\mu$ . Synonymous and  
1148 nonsynonymous mutations detected with duplex sequencing in untreated 24-25-month-old mice  
1149 were used in the analysis. Synonymous and nonsynonymous mutations were considered  
1150 independently. Synonymous mutations are assumed to have no fitness effect and reflect behaviour  
1151 under neutral drift, while non-synonymous mutations were hypothesised to reflect behaviour under  
1152 a positive selective advantage. The density of variants declined at VAF  $5 \cdot 10^{-5}$ , so to only include  
1153 VAF ranges supported by informative variants, only variants above this threshold were included in  
1154 maximum likelihood estimations described below.

1155 How the distribution of VAFs, predicted by our evolutionary framework, changes with age ( $t$ ), the  
1156 variant's fitness effect ( $s$ ), the variant's mutation rate ( $\mu$ ), the population size of HSCs ( $N$ ) and the  
1157 time in years between successive symmetric cell differentiation divisions ( $\tau$ ) is given by the  
1158 following expression for the probability density as a function of  $l = \log(VAF)$ :

$$1159 \quad \rho(l) = \frac{\theta}{(1-2e^l)} e^{-\frac{e^l}{\varphi(1-2e^l)}} \quad \text{where } \theta = N\tau\mu \quad \text{and } \varphi = \frac{e^{st}-1}{2N\tau s}$$

1160 The value of  $\varphi = \frac{e^{st}-1}{2N\tau s}$  is the typical maximum VAF a variant can reach and this increases with  
1161 fitness effect ( $s$ ) and age ( $t$ ). To reach VAFs  $> \varphi$  requires a variant to both occur early in life and  
1162 stochastically drift to high frequencies, which is unlikely. Therefore, the density of variants falls off  
1163 exponentially for VAFs  $> \varphi$ . For neutral mutations ( $s = 0$ ),

1164 
$$\varphi = \frac{t}{2N\tau}$$
  
1165 Because the mouse age  $t$  is known and the neutral  $\varphi$  is measurable from the data, the ratio  $\varphi/t$   
1166 allows us to infer  $N\tau$  from the distribution of neutral mutation VAFs. Because the neutral  $\theta$  is  
1167 measurable from the data, and  $\theta = N\tau\mu$ , we can also infer the neutral mutation rate ( $\mu$ ).

1168 Probability density histograms, as a function of log-transformed VAFs, were generated using  
1169 Doane's method for log(VAF) bin size calculation. Densities were normalised by the product of bin  
1170 sample size and width. Estimates for  $N\tau$  and  $\mu$  were inferred using a maximum likelihood approach,  
1171 minimising the L2 norm between the cumulative log densities and the predicted densities. For  
1172 synonymous mutations, maximum likelihood estimates were optimised for  $N\tau$  and  $\mu$ . For  
1173 nonsynonymous mutations, variants with VAFs below the observed maximum synonymous VAF  
1174 ( $1.99 \cdot 10^{-4}$ ) were used – these variants are within the “neutral” range – and estimates were  
1175 optimised for with the  $N\tau$  estimated from synonymous mutations.

#### 1176 *Differential fitness effects*

1177 We estimated the distribution of fitness effects across nonsynonymous variants using our derived  
1178 estimates of  $N\tau$  and nonsynonymous  $\mu$ . We parameterised the distribution of fitness effects using  
1179 an exponential power distribution, which captures a strongly decreasing prevalence of mutations  
1180 with high fitness:

$$1181 \mu_{non-neutral}(s) \propto \exp\left[-\left(\frac{s}{d}\right)^\beta\right]$$

1182 The shape of the distribution was fixed to  $\beta = 3^{97}$ . Using the VAF density histograms from  
1183 nonsynonymous variants, we estimated the scale of the distribution and non-neutral mutation rate:  
1184  $\int_{s=0}^{\infty} \mu_{non-neutral}(s) ds$ . The maximum likelihood fit predicted a scale of about  $d=2$  and the  
1185 proportion of non-neutral nonsynonymous mutations to be about 12% (Fig.6B).

#### 1186 **Code and data availability**

1187 SNVs and indels were detected using CaVEMan (version 1.14.0,  
1188 <https://github.com/cancerit/CaVEMan>), cgppindel (version 3.9.0,  
1189 <https://github.com/cancerit/cgppindel>), and VarDict (version 1.8.3,  
1190 <https://github.com/AstraZeneca-NGS/VarDictJava>). Variants were annotated using VAGrENT  
1191 (version 3.7.0, <https://github.com/cancerit/VAGrENT>), and Ensemble VEP (release 107-110.0,

1192 <https://github.com/Ensembl/ensembl-vep>). Phylogenies were constructed using MPBoot (version  
1193 1.1.0, <https://github.com/diepthihoang/mpboot>). Variants were assigned to phylogenies using  
1194 Rtreemut (<https://github.com/nangalialab/treemut>). Population trajectories were inferred using  
1195 *phylodyn* (<https://github.com/mdkarcher/phylodyn>). Bayesian inferences utilized the packages  
1196 *rsimpop* (<https://github.com/nangalialab/rsimpop>) for simulations and *abc* (version 2.2.1,  
1197 <https://CRAN.R-project.org/package=abc>) for approximate Bayesian Computation. Mutation  
1198 signatures were inferred using the hdp (<https://github.com/nicolaroberts/hdp>) and sigfit (version  
1199 2.2.0, <https://github.com/kgori/sigfit>). Duplex consensus reads were generated using the fgbio suite  
1200 of tools (version 1.5.1-2.1.0, <http://fulcrumgenomics.github.io/fgbio/>). dN/dS ratios were calculated  
1201 using dNdScv (version 0.1.0, <https://github.com/im3sanger/dndscv>). Population genetic analyses  
1202 of clone sizes and parameter inferences were based on code available at  
1203 <https://github.com/blundelllab/ClonalHematopoiesis/>. Other analyses were carried out using  
1204 custom R scripts and will be available at <https://github.com/CDKapadia/somatic-mouse>.

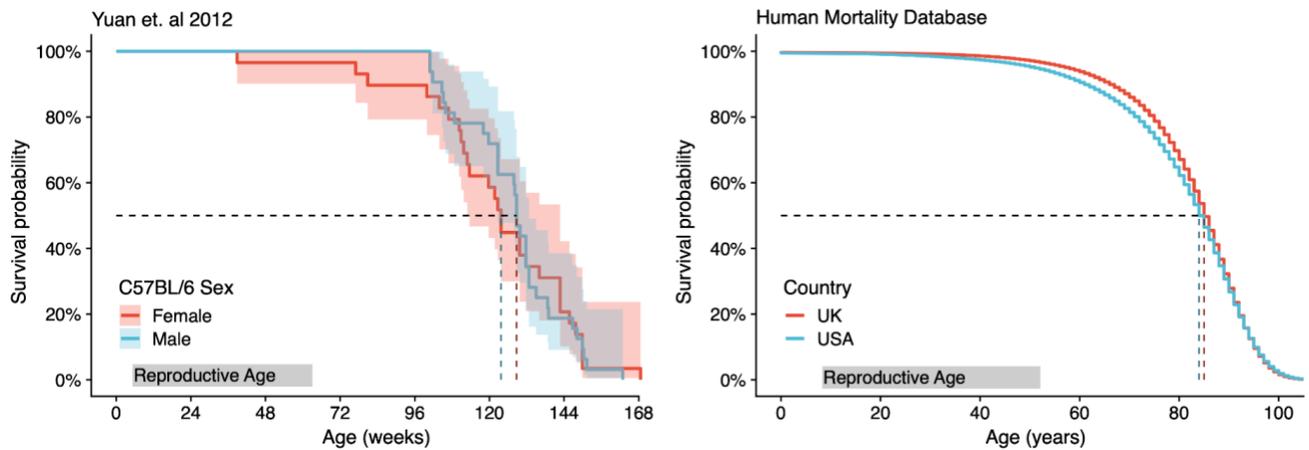
1205  
1206 Whole genome sequencing data will be deposited at the European Nucleotide Archive at accession  
1207 numbers ERP138320 and ERP144323. Targeted duplex sequencing data will be deposited at  
1208 NCBI BioProject PRJNA1033340.

1209

1210

1211 **Supplementary Note 1: Age equivalents between mouse and human**

1212  
1213 We used mouse and human survival data to estimate age equivalency between species. The  
1214 median lifespan of C57BL/6J laboratory mice is 28-months<sup>24</sup> (published data reproduced in  
1215 Supplementary Fig.S1). We retrieved 2017 life-table data from the USA and the UK compiled at  
1216 the Human Mortality Database (mortality.org). Only female data was included to match the makeup  
1217 of our aged mouse dataset. We took the average of the median lifespans in the UK (82.5) and USA  
1218 (80.7) to estimate the female mean lifespan as 81.6 years. Lastly, we normalised mouse age by  
1219 median lifespan to determine an estimated equivalent human age. The above was only performed  
1220 for the aged samples. Mice reach sexual maturity earlier in lifespan relative to humans, so age-  
1221 equivalency was determined by onset of reproductive maturity between species, as previously  
1222 described<sup>98</sup>.  
1223

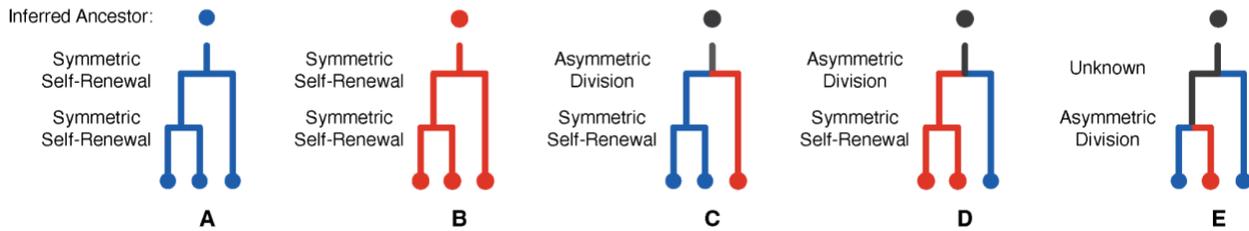


1224  
1225 **Supplementary Figure S1.** Mouse (C57BL/6J strain) survival data (left graph) by age for males (blue) and  
1226 females (red). Human survival data (right graph) by age for females in the UK (red) and USA (blue). Dashed  
1227 black lines mark the age of 50% survival probability for both species. Reproductive age is highlighted by the  
1228 grey box.  
1229  
1230

1231 **Supplementary Note 2: Ancestral cell identity inference**

1232

**No HSC-MPP hierarchy assumption**



1233

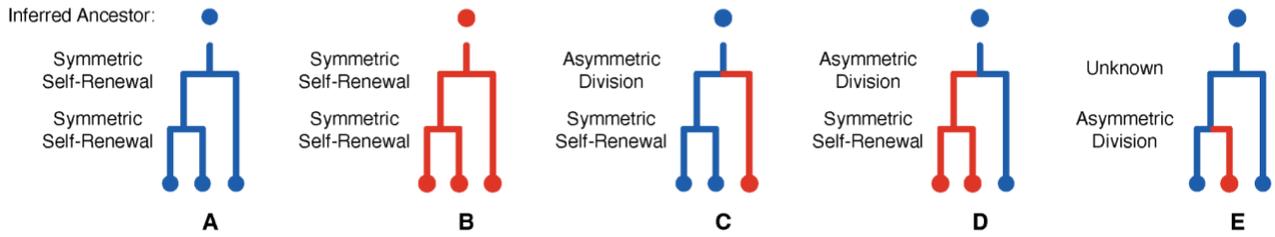
1234

1235 In our phylogenies, coalescences represent cell divisions of ancestral cells whose progeny have  
1236 been captured as observable cells (tips on the tree). Comparison of the observed cell identity  
1237 between closely related tips allows inferences of the identity of their most recent common ancestor  
1238 (MRCA) and the nature of the ancestral cell division captured as a coalescence on the phylogenetic  
1239 tree.

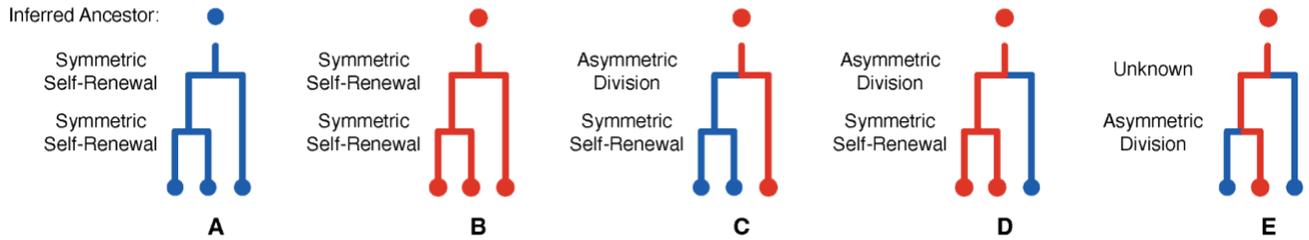
1240

1241 As an illustrative example above, if two closely related observed ('tip') cells are HSCs (scenarios A  
1242 and C), then it is inferred that their most recent common ancestor was also an HSC. This HSC  
1243 must have symmetrically divided to create two daughter HSCs, with both lines of descent also  
1244 generating HSCs that were eventually sampled as the observed cells. From this inferred cell identity  
1245 of their most recent common ancestor, if the cell state of the next closest relative is also an HSC,  
1246 then their most recent common ancestor is similarly inferred to be an HSC (scenario A). HSCs  
1247 coalescences are in blue, while MPPs are in red (scenarios B and D). Neighbouring tip states that  
1248 differ in cell type (e.g., 1 HSC and 1 MPP as in scenario E) can arise in two ways. First, there may  
1249 have been an ancestral asymmetrical cell division generating one HSC and one MPP initially, with  
1250 subsequent progeny along both lines of descent retaining these identities until sampling.  
1251 Alternatively, the same tip states could also occur via a symmetrical self-renewing division of either  
1252 MPP or HSC, followed by a later cell type change (e.g., via asymmetric cell division or direct  
1253 change) of one of the daughter cells. Either way, one cell type change from HSC to MPP (or MPP  
1254 to HSC) is required to explain these tip states; therefore we mark their ancestral coalescence as  
1255 blue/red. In these scenarios, because we cannot infer the cell identity of the MRCA, the upstream  
1256 lineage is subsequently labelled in black. These principles can be applied to all coalescences in  
1257 the observed phylogenetic trees (Fig.2A, Extended Data Fig.2, scenarios A-E above). This intuition  
1258 does not rely on any assumptions of ontogeny, such as the hierarchy of HSCs over MPPs.

### HSC-first Assumption



### MPP-first Assumption



1262 However, current models of the hematopoietic hierarchy dictate that HSCs give rise to MPPs  
1263 (HSC>MPP). With this assumption made (top row “HSC-first assumption), one can label more of  
1264 the branches and coalescences assigned as ‘black’ in the logic detailed above. For example,  
1265 assuming an ‘HSC-first’ hierarchy, the common ancestor for scenarios C-E is now inferred to be  
1266 an HSC, and the unobservable ancestral division in scenario E is inferred to be an HSC self-  
1267 renewal. In the ‘MPP-first’ assumption (lower row), the inverse is inferred. These heuristics are  
1268 applied to all coalescences in the observed phylogenetic trees.

1269

1270 We then asked how many cell state transitions are required to explain the tip states given an HSC-  
1271 first or an MPP-first model. To perform this comparison, for each tree, we subsampled the largest  
1272 category of HSC and MPPs so that there were equal numbers of MPP and HSC tips. To reduce  
1273 the risk of the downsampling being unrepresentative the subsampling was conducted 10,000 times  
1274 for each tree, and the average number of required transitions under the two unidirectional models  
1275 was calculated. We then counted the total number of transitions required to result in the observed  
1276 cell type tips. It was observed that the number of transitions required was similar for HSC-first and  
1277 MPP-first and that there was no consistent pattern of one being higher than the other (Fig.2D). It  
1278 was then natural to ask whether our cell type information was at all informative and so we randomly  
1279 permuted the tip cell types and then resolved the tree in an HSC-first fashion. This sub-sampling  
1280 and permutation was carried out 10,000 times and, as expected, the number of changes required

1281 by either the MPP-first or HSC-first models were generally far fewer than is consistent with the null  
1282 model that all balanced cell type categorisations require the same number of tree-based transitions  
1283 to explain the tip phenotype under an unidirectional model. In summary, both HSC-first and MPP-  
1284 first models are less parsimonious, i.e., requiring more cell state changes, than the model first  
1285 presented in which no assumptions are made about a hierarchy between HSCs and MPPs. The  
1286 most parsimonious model would be that of HSC and MPP lineages being derived in parallel during  
1287 similar development periods from non-overlapping common ancestors.

1288

### 1289 **A simple 3 state model for Murine Progenitor Ontogeny**

1290 To formalise the above ideas in the context of a simple model of HSC and MPP ontogeny, we  
1291 considered the state of all cells prior to 10 mutations in molecular time as being in an embryonic  
1292 precursor state (EMB), given that haematopoietic and colonic lineages remain uncommitted until  
1293 at least this time (Extended Data Fig.3). We then assumed that in each unit molecular time there  
1294 is a fixed probability of transitioning out of this embryonic state into either an HSC state,  $p_{EMB \rightarrow HSC}$ ,  
1295 or an MPP state,  $p_{EMB \rightarrow MPP}$ . Furthermore, there is a fixed probability of transitioning from an HSC  
1296 to an MPP,  $p_{HSC \rightarrow MPP}$ , and from an MPP to an HSC,  $p_{MPP \rightarrow HSC}$ . Thus, the evolution of the cells  
1297 down the tree is governed by a discrete time Markov chain process. The likelihood of the observed  
1298 tip cell types is calculated using a hidden Markov tree approach (Methods). Maximum likelihood  
1299 estimates of the model parameters are obtained by maximising the sum of the log-likelihoods  
1300 across mouse-specific phylogenetic trees. Finally, for each mouse, the most likely sequence of  
1301 unobserved states for the nodes of the phylogenetic tree is calculated using the fitted model  
1302 parameters.

1303

1304 We performed the maximum likelihood estimation using the R package “bbmle”. The maximisation  
1305 was performed on logit transformed quantities:  $p_{HSC \rightarrow MPP}$ ,  $p_{MPP \rightarrow HSC}$ ,  $\frac{p_{EMB \rightarrow HSC}}{p_{EMB \rightarrow HSC} + p_{EMB \rightarrow MPP}}$  and  
1306  $p_{EMB \rightarrow HSC} + p_{EMB \rightarrow MPP}$ . Whilst we were able to obtain parameter estimates and Hessian-based  
1307 standard errors, the profile-based estimation of confidence intervals did not work in all cases. So,  
1308 to obtain more robust estimates in the CIs of the model we implemented a Stan-based Bayesian  
1309 version of the model using the directly calculated likelihood as described above. Uniform priors on  
1310 the unit interval were assumed for  $\frac{p_{EMB \rightarrow HSC}}{p_{EMB \rightarrow HSC} + p_{EMB \rightarrow MPP}}$  and  $p_{EMB \rightarrow HSC} + p_{EMB \rightarrow MPP}$  and uniform

1311 priors on the interval (0-0.5) were assumed for both  $p_{HSC \rightarrow MPP}$  and  $p_{MPP \rightarrow HSC}$ . The model was  
1312 run with four chains, each for 10,000 iterations.

### 1313 **Separate young and old mouse cohorts provide optimal model fit**

1314 We compared fitting the model with a per-mouse, per-age, and pan cohort strata. A likelihood ratio  
1315 test analysis revealed that the best model is an age-specific model where parameters are estimated  
1316 separately in old and young mice (Supplementary Table S1).

1317 **Supplementary Table S1.** Log likelihood values and Akaike information Criterion (AIC) assessing model fit.

Model	Degrees of freedom	Log Likelihood	AIC	Likelihood Ratio Test
Pan Cohort	4	-840.2	1,688.4	
Age-Specific	8	-800.2	1,616.5	vs. Pan Cohort: P=1.78e-16
Mouse-Specific	24	-796.2	1,640.3	vs. Age Specific: P=0.945

1318 A lower AIC value indicates a better model fit.

1319

### 1320 **Young and Old mice exhibit differing patterns of differentiation**

1321 Applying our hidden Markov tree approach, we fitted an HSC-first model where,

$$1322 \frac{p_{EMB \rightarrow HSC}}{p_{EMB \rightarrow HSC} + p_{EMB \rightarrow MPP}}$$

1323 is fixed at unity. That is, all EMB must transition to an HSC before any emergence of MPPs (HSC-  
1324 first). In the context of this simple model, we can reject the HSC-first model across the combined  
1325 age group model ( $p=1.11e-18$ ) and also for the old group ( $p=1.38e-19$ ). However, we were unable  
1326 to reject the HSC-first model for the young animal group ( $p=0.397$ ).

1327

1328 Examining the types of cell-state transitions in the trees, we observed that aged animals exhibit  
1329 several independent transitions from the embryonic precursor state followed by relatively few  
1330 transitions between HSC to MPP or vice versa (Supplementary Fig.S2). In contrast, young animals  
1331 exhibit a tendency towards HSC-first followed by a relative abundance of HSC->MPP transitions  
1332 (Supplementary Fig.S2). Both HSC-specification and MPP-specification occur within the first 50  
1333 mutations molecular time (Supplementary Fig.S3). The cell identity transition rates, per unit  
1334 molecular time, are listed below and were used to generate Fig.2E.

1335

1336

	$\mathcal{P}_{EMB \rightarrow HSC}$	$\mathcal{P}_{EMB \rightarrow MPP}$	$\mathcal{P}_{HSC \rightarrow MPP}$	$\mathcal{P}_{MPP \rightarrow HSC}$
Young Donors	0.158	0.037	0.0164	0.0070
Aged Donors	0.036	0.034	0.0013	0.0006

1337

1338 The above result are fairly consistent with the Stan based results for which we show the medians  
 1339 of the marginal posterior distribution followed by the 95% credibility intervals:

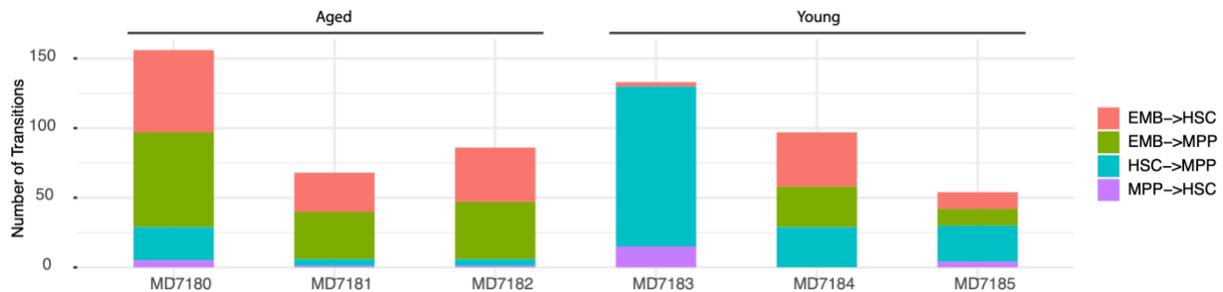
1340

	$\mathcal{P}_{EMB \rightarrow HSC}$	$\mathcal{P}_{EMB \rightarrow MPP}$	$\mathcal{P}_{HSC \rightarrow MPP}$	$\mathcal{P}_{MPP \rightarrow HSC}$
Young	0.43(0.11 - 0.9)	0.063(0.0048 - 0.27)	0.017(0.014 - 0.022)	0.0078(0.0029 - 0.016)
Aged	0.04(0.025 - 0.068)	0.038(0.025 - 0.064)	0.0014(0.00079 - 0.0022)	0.00071(0.00022 - 0.0015)

1341

1342 Of note, the mode of the marginal posterior distribution of  $\mathcal{P}_{EMB \rightarrow HSC}$  peaks at 0.19, which is  
 1343 reassuringly close to the maximum likelihood estimate of 0.158.

1344



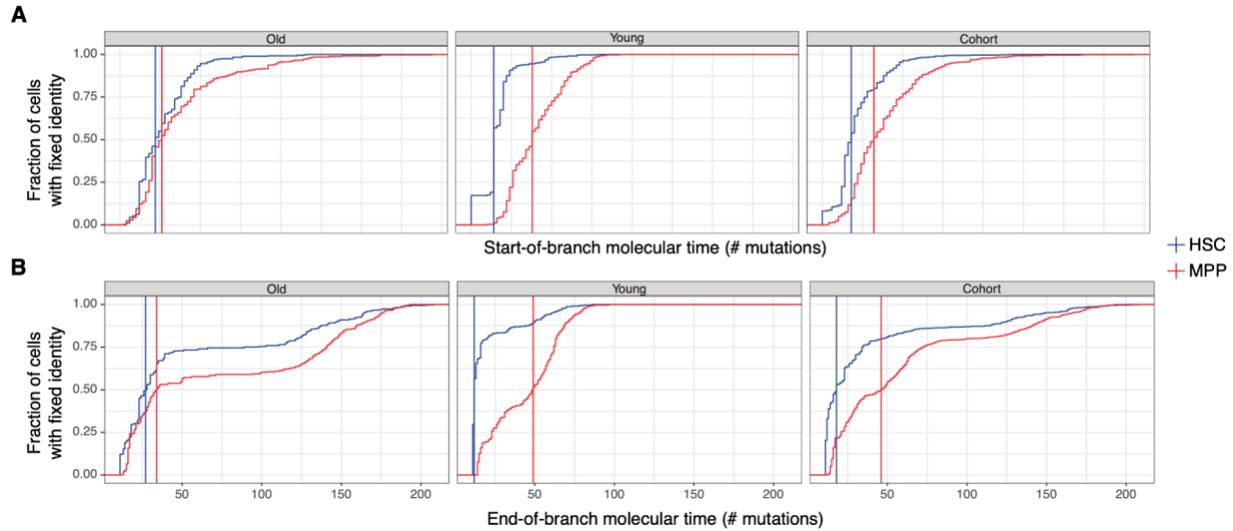
1345

1346 **Supplemental Figure S2: Transition Type Counts.** The old mice exhibit an abundance of approximately  
 1347 equally prevalent EMB->HSC and EMB->MPP transitions followed by relatively few transitions to the eventual  
 1348 observed cell types. The young mice exhibit relatively fewer EMB->HSC and EMB->MPP and then a relative  
 1349 abundance of HSC->MPP transitions.

1350

1351

1352

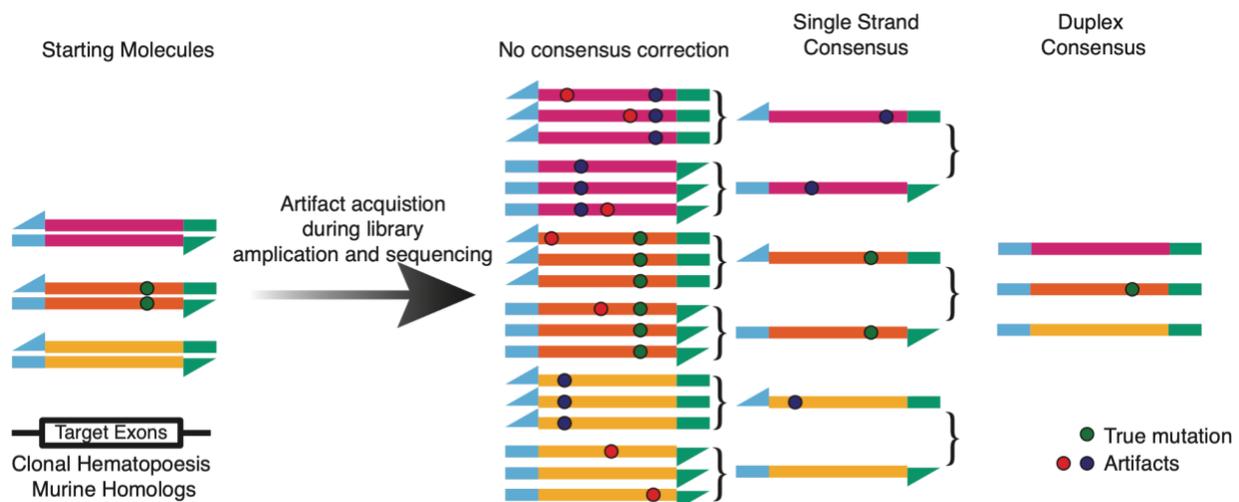


1353

1354 **Supplementary Figure S3: Cumulative Distribution of Specification Timing.** For each colony we use the  
1355 molecular time of the **A)** start of the branch, or **B)** end of the branch on which the ancestral lineage first  
1356 transitions to the observed cell type as an upper bound for the timing of its transition to its final observed  
1357 state. The panels show the cumulative distribution of these upper bounds calculated from the most likely  
1358 sequence of transitions inferred using the age-specific model and the pan cohort model. Vertical lines indicate  
1359 the time at which 50% of the sampled cells have specified identity.  
1360

1361 **Supplementary Note 3: Quality control of targeted duplex-sequencing**

1362 We expected that somatic clones in mice might be rare events at small clone sizes, thus would  
1363 require a sensitive detection assay. High-depth sequencing can be used for detection of subclonal  
1364 variants, but with increasing coverage, the error-rate intrinsic to short-read sequencing can obscure  
1365 true low variant allele fraction (VAF) variants. To circumvent this sensitivity limit, read-level error-  
1366 correction approaches are necessary. Thus, we applied duplex-consensus sequencing, which  
1367 offers among the highest sensitivity for subclonal variant detection. In duplex-sequencing, each  
1368 initial dsDNA molecule is uniquely barcoded such that reads derived from complementary 5' and  
1369 3' strands are linked, but also distinguishable. Detected variants must be present on both uniquely  
1370 barcoded strands of the initial dsDNA fragment to pass bioinformatic filtration (Supplementary  
1371 Fig.S4). By enforcing that variants are present in reads derived from both of the matched  
1372 complementary strands of DNA, one can eliminate the majority of sequencer-induced artefacts that  
1373 usually hamper sensitivity. To apply this technology to murine clonal haematopoiesis (CH), we  
1374 developed a target panel of the mouse homologs of genes most frequently mutated in human CH  
1375 (Methods, Supplemental File 2).



1376

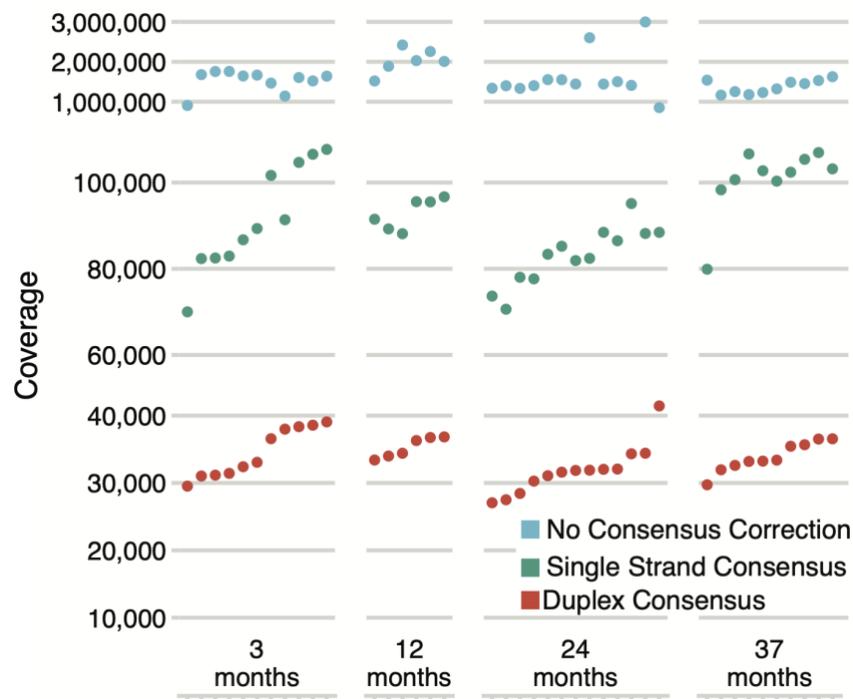
1377 **Supplementary Figure S4: Error correction strategy in targeted duplex sequencing.**

1378 PCR during library preparation and sequencing introduce low-frequency artefacts. Duplex barcodes allow  
1379 grouping of PCR duplex reads from a single DNA library molecule (read families) and single-strand read  
1380 consensus generation. Next, single strand consensus reads from complementary strands on initial dsDNA  
1381 are matched to generate a duplex consensus. To build a duplex-consensus read, we required at least 3  
1382 reads in each supporting read family (i.e., at least 3 sequenced PCR duplicates of matched top and bottom  
1383 strands from an original dsDNA molecule).

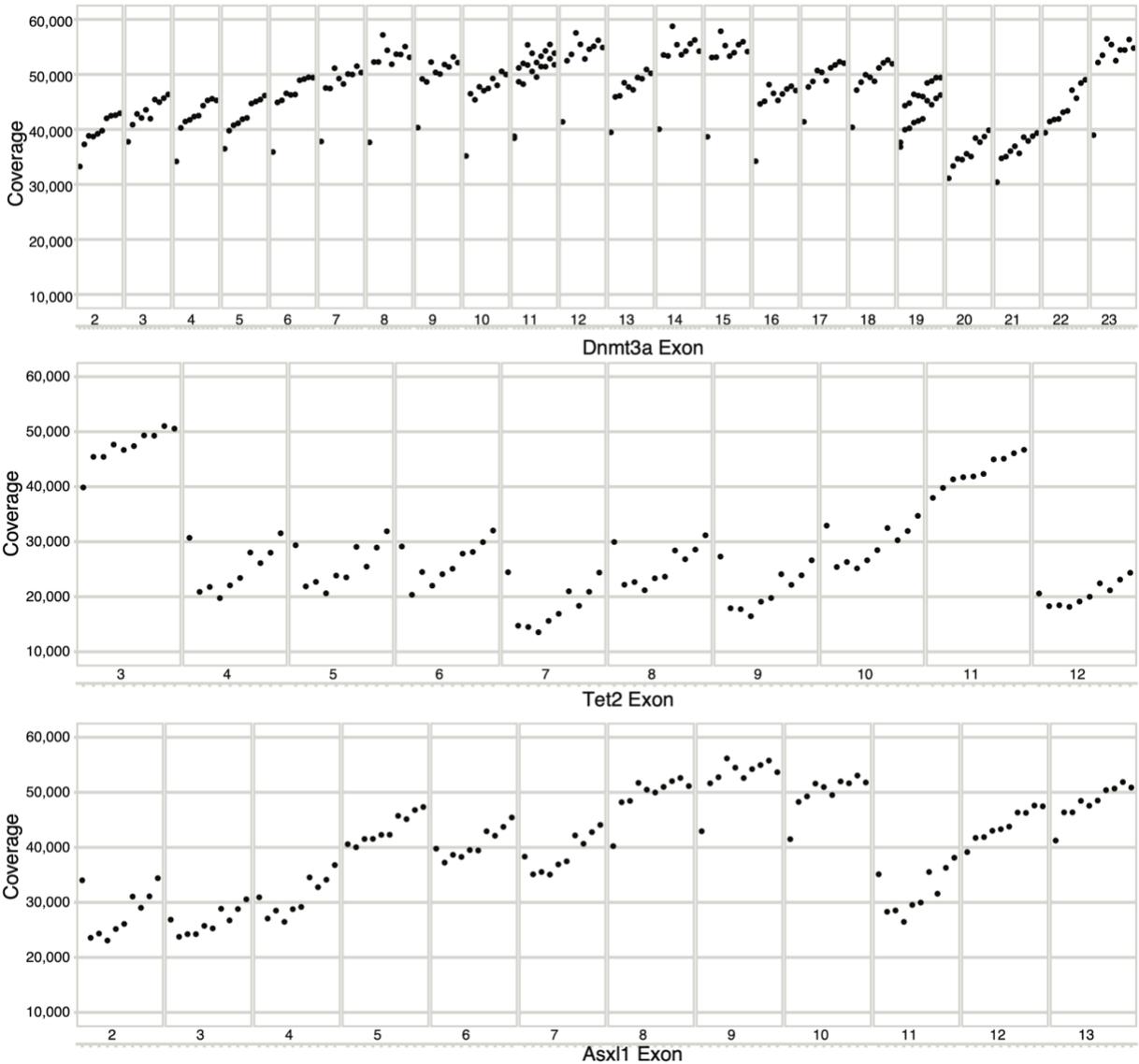
1384

1385

1386 *Coverage requirements to generate a duplex consensus:* To generate a duplex-consensus read,  
1387 an initial DNA molecule must be sequenced multiple times with reads from matched 5' and 3'  
1388 strands sufficiently represented. To ensure that clone detection sensitivity would not be limited by  
1389 input genomic DNA (*i.e.*, the libraries contained sufficient genomic complexity), we input at least  
1390 100,000 genomic equivalents (or at least 1650 ng of genomic DNA) into our library preparations.  
1391 High library complexity decreases the probability of matched 5' and 3' reads being sequenced by  
1392 chance; thus, even with a target panel enrichment, extremely high sequencing depth is required to  
1393 capture library complexity in duplex consensus reads. Median raw, non-deduplicated coverage  
1394 spanned 1,000,000X to 3,000,000X at targeted loci per sample. This correlated with a single-strand  
1395 consensus coverage spanning 60,000X-120,000X, which, after 5' and 3' linkage, further collapsed  
1396 to duplex consensus coverage spanning 30,000X-40,000X (Supplementary Fig.S5). Duplex  
1397 coverage at specific exons within targeted genes was variable between samples (Supplementary  
1398 Fig.S6)



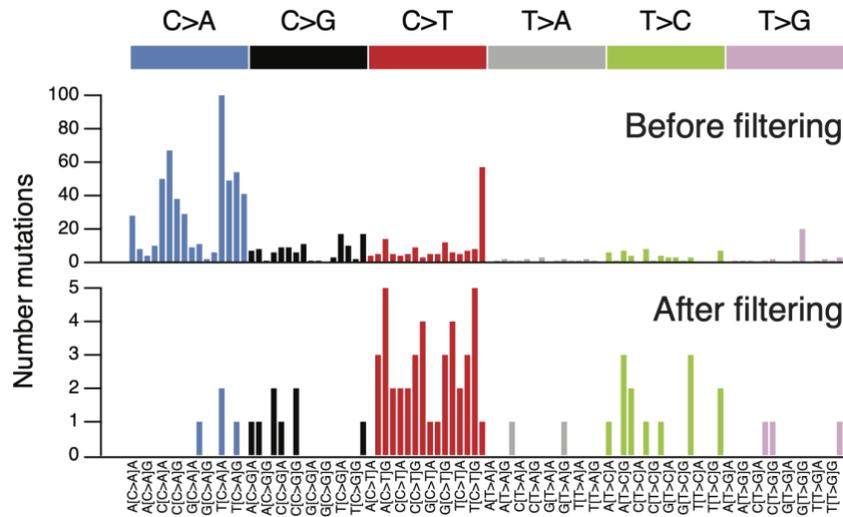
1399  
1400 **Supplementary Figure S5:** Sequencing coverage at targeted loci for all samples in Fig.4A. The relationship  
1401 between raw (not deduplicated), single strand consensus, and duplex consensus coverage is shown.



1402  
1403 **Supplementary Figure S6:** Duplex coverage at coding exons in *Dnmt3a*, *Tet2*, and *Asx1* for a series of  
1404 aged samples.  
1405

1406 *Mutation filtering:* To supplement the sensitivity afforded by duplex sequencing, stringent read- and  
1407 variant-level filters were applied to reduce the presence of false positive mutations or spurious  
1408 calls. Without filtration, we observed an enrichment of C>A mutations (Supplementary Fig.S7),  
1409 reminiscent of mutation signature SBS45, which is likely attributable to oxidative damage during  
1410 sequencing<sup>90,91</sup>. Such oxidative damage mutations likely arose after duplex barcode attachment,  
1411 were enriched at read ends, and likely caused mutations within the duplex barcode sequence. Due  
1412 to mutations in duplex barcodes, a read family derived from a single initial dsDNA molecule (a

1413 singleton) would erroneously appear as derived from an additional read family (a doublet). This  
 1414 observation led us to apply a stringent series of filters (Methods), after which the trinucleotide  
 1415 spectra of variants detected in duplex sequencing more resembled that seen with blood  
 1416 (Supplementary Fig.S7).



1417 **Supplementary Figure S7:** Trinucleotide spectra of duplex-sequencing variants before and after post-  
 1418 processing filters. See Methods for filtering strategy details.  
 1419  
 1420

1421 *Clone size calculation:* Given the differences in coverage between loci, we normalised the variant  
 1422 read counts to allow accurate clone size comparisons between samples. In general, clone size for  
 1423 a given variant is defined as:

$$Clone\ size = \frac{Mutant\ allele\ read\ count}{Total\ read\ count}$$

1424 For very small clones, there is a degree of stochasticity affecting if sufficient mutant read alleles  
 1425 will be converted to duplex consensus reads to allow detection. Duplex clones supported by very  
 1426 few mutant allele reads would have a low numerator, thus clone size estimations may be skewed.  
 1427 Given more single-strand consensus reads are generated than duplex consensus reads  
 1428 (Supplementary Fig.S6), we reasoned that mutant allele reads would be relatively more abundant  
 1429 within single-strand consensus reads – that is, mutant allele reads would be present among the  
 1430 reads ‘discarded’ due to insufficient evidence to generate a duplex consensus. To normalise clone  
 1431 size, especially in low-magnitude clones, we used de-duplicated single-strand consensus reads,  
 1432 as follows:

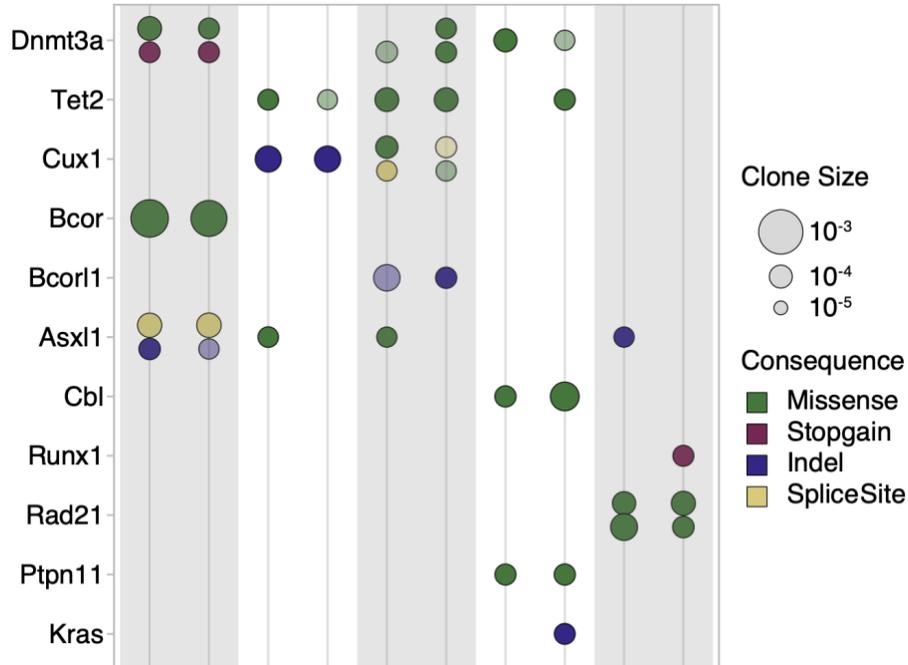
$$Clone\ size_{corrected} = \frac{Mutant\ allele\ read\ count_{single\ strand\ consensus}}{Total\ read\ count_{single\ strand\ consensus}}$$

1433 By using the de-duplicated single-strand consensus reads for clone size calculation, the numerator

1436 (variant allele count) and denominator (coverage) both increase, reducing any skewing that may  
1437 be present in clone size calculations from duplex consensus reads. All clone sizes depicted on dot  
1438 plots are calculated in this manner.

1439  
1440 *Biological replicates:* To validate reproducibility within the targeted duplex-sequencing library  
1441 preparation and variant calling pipelines, we assessed clone prevalence in biological replicate  
1442 samples. For each replicate, peripheral blood was separately collected (in different tubes) and  
1443 underwent genomic DNA extraction independently. Thus, the genomic DNA “pools”, while derived  
1444 from the same sample mouse, were purified in separate reactions. Replicate DNA samples  
1445 underwent duplex library preparation and variant calling as described in Methods. As shown in  
1446 Supplementary Figure S8, clone detection is concordant between paired replicates. Clones unique  
1447 to a single replicate were at the limit of detection for the specific locus, and thus it is likely in the  
1448 paired replicate that insufficient variant reads were sequenced to generate duplex consensus read  
1449 support. Such borderline detectable clones will likely be detectable within single-strand consensus  
1450 reads, which carry nearly double greater read depth, though at the expense of duplex sensitivity.  
1451 We examined single-strand consensus reads from the biological duplicate samples and were able  
1452 to “rescue” missing variants from the paired replicate sample, in about half of cases (Supplementary  
1453 Figure S8). This confirms that much of the missing replicate clones were lost during duplex  
1454 consensus building, for example when a clone has insufficient top or bottom strand support to  
1455 create a duplex read.

1456



1457 **Supplementary Figure S8. Native CH in biological replicate samples.** Shaded and unshaded pairs  
 1458 represent duplex libraries separately prepared from an identical initial blood sample. Clones are presented  
 1459 as described in Fig.4A. Transparency indicates a clone that was only detectable within single strand  
 1460 consensus reads but not duplex consensus reads.  
 1461  
 1462

1463 *In silico estimation of the sensitivity and specificity of duplex-sequencing results:* We next sought  
 1464 to understand if the degree of clone concordance between biological replicate samples was  
 1465 consistent with the sensitivity of our assay. We consider a simple model for SNVs of conditional  
 1466 base calling probabilities for the reference base (R), a mutant base (A) and the two other bases  
 1467 (B,C). For an individual read (or read family/bundle) the probability of observing the bases is  
 1468 modeled in the following manner:  
 1469

$$\begin{aligned}
 1472 \quad P(\text{Base is } A) = & \\
 1470 \quad & P(\text{DNA Molecule is mutant } A \text{ at site}) * \\
 1471 \quad & P(\text{Base called as mutant } A | \text{DNA Molecule is mutant } A \text{ at site}) \\
 1474 \quad & + P(\text{DNA Molecule is not mutant at site}) * \\
 1473 \quad & P(\text{Base called as mutant } A | \text{DNA Molecule is not mutant at site}).
 \end{aligned}$$

1475  
 1476 Now  $P(\text{DNA Molecule is mutant at site}) = \frac{\text{Aberrant Cell Fraction}}{\text{ploidy}} = VAF$  where for economy we now  
 1477 use the term (true) VAF to characterise the clone. Moreover we assume there is a base calling

1478 error rate (“epsilon”)  $\epsilon$ . It is assumed that this results in the one of the 3 incorrect bases to be called  
 1479 with equal probability of  $\epsilon/3$ :

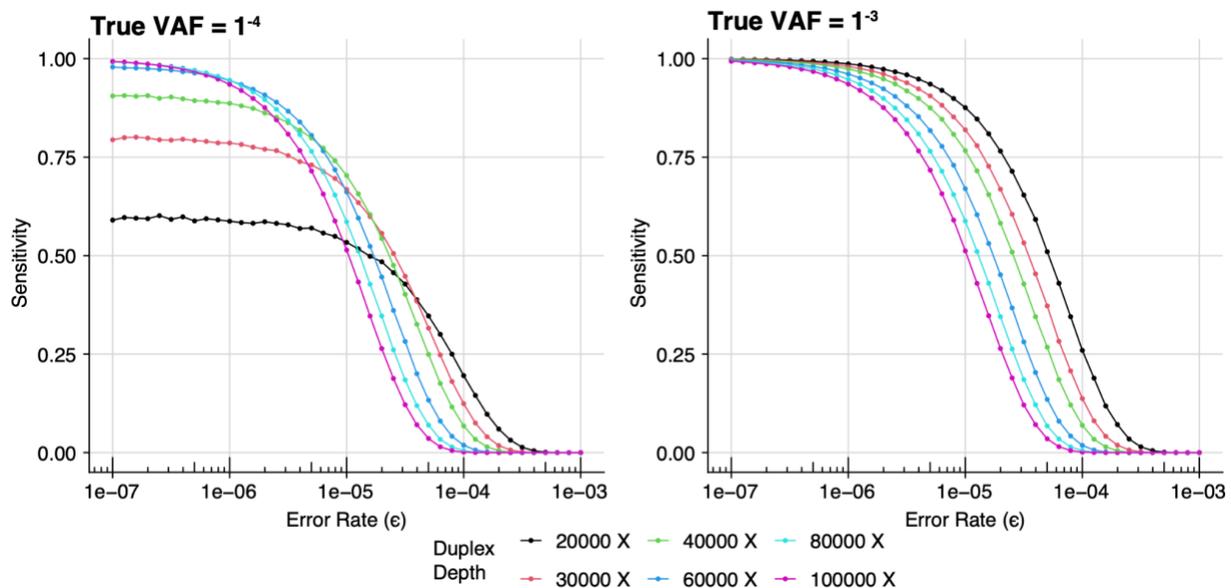
1480 
$$P(\text{Base is Reference}) = (1 - VAF)(1 - \epsilon) + VAF \frac{\epsilon}{3}$$

1481 
$$P(\text{Base is A}) = VAF(1 - \epsilon) + (1 - VAF) \frac{\epsilon}{3}$$

1482 
$$P(\text{Base is B}) = VAF \frac{\epsilon}{3} + (1 - VAF) \frac{\epsilon}{3} = \frac{\epsilon}{3}$$

1483 
$$P(\text{Base is C}) = \frac{\epsilon}{3}$$

1484 For a given bait set wide depth of sequencing, *depth*, a given site has depth that is Poisson  
 1485 distributed with mean *depth*. For a clone to be detected it is only required that at least 2 mutant  
 1486 reads are observed. We assume we have a known clone, with  $VAF=1^{-4}$  or  $VAF=1^{-3}$ , and with  
 1487 mutant allele A. The A clone is discovered if there are 2 or more mutant “A” reads, and no other  
 1488 mutant reads (“B” or “C”). With these criteria, we can plot the sensitivity for a given error rate,  $\epsilon$ ,  
 1489 shown below in Supplementary Figure S9.

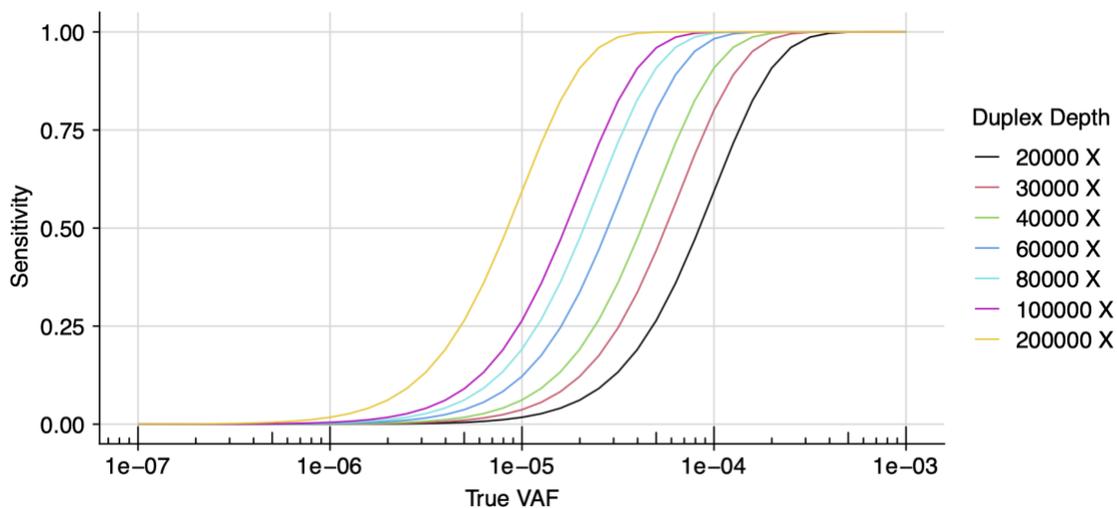


1490  
 1491 **Supplementary Figure S9: True clone discovery across error rates using multinomial modeling.**  
 1492 Estimated sensitivity of detecting a variant at a given site with true VAF 1<sup>-3</sup> (left) or 1<sup>-4</sup> (right) across  
 1493 increasing error rates. A range of duplex depth at variant sites are shown.  
 1494

1495 The above plots show that using single strand consensus sequencing with error rate of  $\sim 3^{-5}$  at  
 1496 depth 60,000x provides a sensitivity of 30% for clone sizes of  $VAF=1^{-3}$  or less. However, using

1497 duplex depth of 30,000x with an error rate of  $1^{-6}$  to  $1^{-7}$  (as described in Kennedy *et al.*<sup>87</sup>) provides  
1498 a sensitivity of >75%.

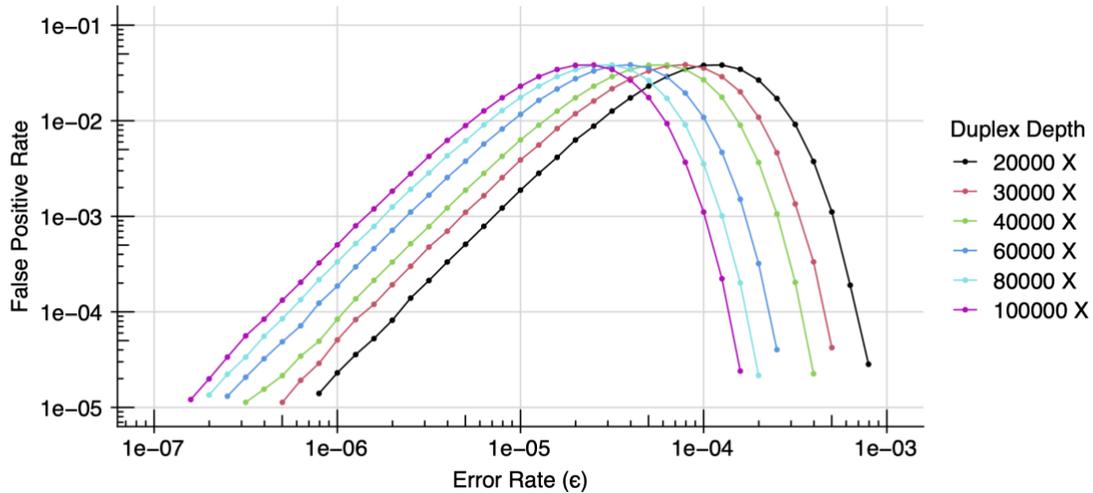
1499  
1500 If we assume the extreme (and implausible) case of error-free sequencing, then the clone detection  
1501 sensitivity is purely governed by the binomial distribution with a probability of True VAF.  
1502 Importantly, even if the sequencing was error-free, we would not expect there to be concordance  
1503 of clone detectability in different samples. In the Supplementary Figure S10 below we can see that  
1504 for error-free sequencing at depth 20,000X, we would actually have a concordance of around 60%.  
1505 This aligns with the observed duplex clone concordance seen in the biological replicate samples  
1506 shown in Supplementary Figure S8.



1507  
1508 **Supplementary Figure S10: True clone discovery with error-free sequencing.** The estimated sensitivity  
1509 for clone detection at increasing VAFs in the scenario of error-free variant detection. In the absence of an  
1510 error rate, detection sensitivity can be described with a binomial distribution. A range of duplex depth at  
1511 variant sites are shown.

1512  
1513 Finally, we can estimate the probability of false positive clone detection at a given error rate  $\epsilon$ . As  
1514 shown below, when querying a range of feasible duplex-sequencing sensitivities and duplex-  
1515 corrected sequencing depths, a false positive clone is far less likely than a false negative clone  
1516 (missing a true event). As an illustrative example, for duplex depths 20,000X to 30,000X and the  
1517 duplex error rate of  $<8e-04$  (estimated error rate of  $<1e-06$ ), the false positive rate is  $<0.01$ .  
1518 (Supplementary Figure S11).

1519



1520

1521 **Supplementary Figure S11: False positive variant detection.** The estimated incidence of incorrectly  
1522 detecting a variant at a given site is shown, using multinomial modeling of detection error rate and site-  
1523 specific duplex depth.

1524

1525

1526

*Concordance of duplex sequencing data explored through mixing mutant and wildtype reads:*

1527

The in-silico analyses described above suggest that a true variant clone may not be observed due  
1528 to insufficient duplex read support, and sensitivity increases with additional duplex depth. In this  
1529 case, an *expected* variant would likely be detectable in single-strand consensus reads  
1530 (Supplementary Figure S4), which require reduced read support to build a consensus read, and  
1531 harbour far higher coverage (Supplementary Figure 5), though at the expense of sensitivity.

1532

1533 We performed a mixing analysis using our duplex data, with the aim to evaluate 1) the concordance  
1534 of calling serially lower VAF clones in different sample, and 2) the degree missing-but-expected  
1535 clones can be found in single-strand consensus data.

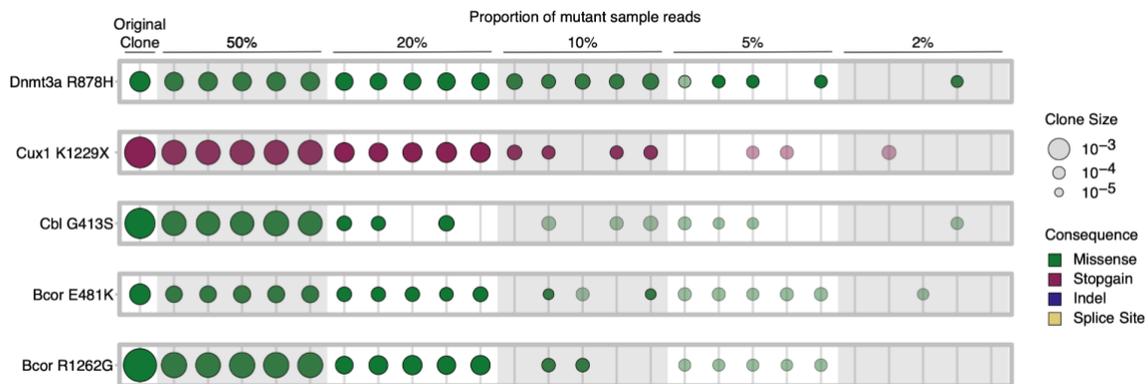
1536

1537 We selected clones with a large detectable clone size, then generated serial dilutions of input  
1538 mutant file reads with wild-type file reads to simulate diminishing read support and the subsequent  
1539 detection of an expected variant in duplex reads. Mutant file reads were diluted by the following  
1540 percentages: 50%, 20%, 10%, 5%, 2%. Five replicates of each random subsample dilution were  
1541 used as technical replicates. Read dilution was done with raw, unmodified reads; that is, before  
1542 any mapping or consensus building steps. Mutant reads were mixed with wildtype reads to the  
1543 same overall read count as the original data, then analysed using the duplex consensus building  
1544 and variant calling pipeline described herein. In cases where the expected variant was not detected

1545 in duplex consensus reads (either due to lack of read support, or failing to pass stringent filters),  
1546 we examined matched single strand consensus reads for the variant, and often were able to detect  
1547 the expected clone.

1548  
1549 As shown below in Supplementary Figure S12, we observe concordance among technical  
1550 replicates when the mutant clone is relatively less diluted from the original data, with reduced  
1551 variant detection in duplex reads as the mutant read support is diminished. The missing variants  
1552 can be rescued when examining single-strand consensus data. With increasing dilution, the variant  
1553 eventually lacks sufficient read support to build both duplex or single-strand consensus reads, and  
1554 is not detectable.

1555



1556 **Supplementary Figure S12: Dilution of mutant reads and subsequent variant call concordance.** For  
1557 five initially large clones, reads from the original input file (ie supporting the observed clone) were diluted at  
1558 the indicated proportion with wildtype reads. Five replicates for each dilution factor are grouped. The original  
1559 clone observed in these unmixed data are shown at the far left column. Clones are presented as described  
1560 in Fig.4A. Transparency indicates a clone that was not detected with standard duplex filtering, but was  
1561 detectable within single-strand consensus reads.

1562  
1563  
1564

## Supplementary Note 4: Inferring population size and division rates from cell phylogenies

### Introduction

The aim of this exercise is to understand the apparent identifiability of the model parameters, seen in Bayesian inferences about the dynamics of HSC populations in mice (and other systems with similar population dynamics), when the phylogeny of a sample of descendent cells is the only available data. These Bayesian inferences were performed using ABC (approximate Bayesian computation) methods. Here the term approximate Bayesian computation refers to a class of Monte Carlo methods for generating samples from posterior distribution, which avoid computation of the likelihood function, by relying on simulation of the model. The more descriptive term *likelihood-free* Bayesian computation is also used. These methods include rejection sampling (pioneered by Pritchard *et al.*<sup>99</sup>), and various regression methods (introduced by Beaumont *et al.*<sup>41</sup>).

The ABC results reported here were obtained by using the *rsimpop* package<sup>34</sup> to perform simulations of models of cell population dynamics, and then using the *abc* package<sup>100</sup> to compute (approximate) marginal posterior densities for modal parameters. The *rsimpop* package allows us to specify a wide range of stochastic growth models based on an underlying birth-death process.

In the case of neutral deterministic growth models, we have exact formulas for the likelihood function, where the model parameter is a sequence of *effective population sizes* (or a sequence of *drift intensities*). For these models, efficient Monte Carlo methods<sup>101</sup> are available for sampling from the exact posterior distribution of the model parameters. In the case of stochastic growth models which can be approximated by a neutral deterministic growth model, we can obtain an approximate formulas for the likelihood function in which sequence of effective population sizes is replaced by a parameter vector which includes birth rates and death rates as model parameters. We will use approximate likelihood functions obtained in this way to address the issue of identifiability of parameters for various models.

### Likelihood functions for neutral models given phylogeny data

When we have genome sequences from a sample of single cells taken from an individual donor, we can construct a phylogeny for the sample, with the mutations assigned to branches. From this phylogeny, we can obtain an ultrametric tree, in which the relative lengths of the branches can be

estimated (taking account of the number of mutations assigned to each branch). We also know the age  $t_S$  of the donor at the time point at which the sample was taken. (Here age is measured from the moment of conception.)

From the phylogeny on a sample of  $n$  cells, together with the estimated absolute branch lengths, we can label the internal nodes (coalescent event) with integers  $2, 3, \dots, n$ , where  $n$  is the label on the most recent node (closest to the time of sample collection), and where  $2$  is the label on the earliest node on the phylogeny (the root node). We let  $S$  be the sequence of node heights ( $S(n), S(n-1), \dots, S(2)$ ), where  $S(r)$  is the height (time in days or years, measured backwards from sample collection) of internal node  $r$ . These node heights are determined by the branch lengths. The same information is contained in the sequence  $T$  of inter-coalescent interval durations ( $T(n), T(n-1), \dots, T(2)$ ). The inter-coalescent interval duration  $T(r)$  is the duration (in days or years) of the time interval during which exactly  $r$  lines of descent remain.

We begin by allowing the neutral model to take a very general form, which can be viewed as a generalisation of the neutral Moran model<sup>102</sup>. For now, we measure time  $t$  forward from conception ( $t = 0$ ) when the population of cells contains a single founder cell ( $N_0 = 1$ ), which is the zygote. This time  $t$  coincides with age (measured from conception). The sequence of distinct time points (ages) at which the population size changes, together with the age  $t_C$  at the time of sample collection, is recorded as  $t = (t_1, t_2, \dots, t_C)$ , where  $0 < t_1 < t_2 < \dots < t_C$ . So we have a sequence of  $C-1$  population events which occur before sample collection. We assume that at each population event, these changes in population size occur instantaneously, so that we can define a population size  $N_k$  which persists throughout the time interval  $[t_k, t_{k+1})$  from event  $k$  to the moment immediately preceding the next event.

At each of these events ( $k = 1, 2, \dots, C-1$ ) at which the population size changes, we allow the number of *births*  $b_k$  (cell division) to be either 0 or 1, and the number of *deaths*  $d_k$  (cells which leave the stem cell population, either via cell deaths, or via cell differentiation events) to any integer value from 0 up to  $N_{k-1}$  (the size of the population when it enters event  $k$ ).

Note that in a birth-death process it is more usual to assume that each event is either a birth event (where  $b_k = 1$ , and  $d_k = 0$ ) or a death event (where  $b_k = 0$ , and  $d_k = 1$ ). However, it turns out that while the analysis outlined below is greatly complicated if we allow  $b_k$  to exceed 1, when we relax

the constraints on  $d_k$  we encounter very little additional difficulty. We have the following recursion for the population size

$$N_k = N_{k-1} + b_k - d_k, \quad (1)$$

for  $k = 1, 2, \dots, C-1$ , subject to the constraints that either  $b_k = 0$  or  $b_k = 1$ , and  $0 \leq d_k \leq N_{k-1}$ . There is one more sequence which it is useful for us to define here. This is the sequence of *drift intensities*,  $\xi = (\xi_1, \xi_2, \dots, \xi_{C-1})$ , where

$$\xi_k = \left(\frac{N_k}{2}\right)^{-1} b_k = \frac{2b_k}{N_k(N_k-1)} \quad (2)$$

for  $k = 1, 2, \dots, C-1$ . Recall that if there was no birth (cell division) at event  $k$ , then  $b_k = 0$ , and therefore  $\xi_k = 0$ . Notice that here we are using the conventional notation  $\binom{n}{k}$ , for binomial coefficients. In particular we have

$$\binom{n}{2} = \frac{n(n-1)}{2} \quad (3)$$

We can define the function

$$b(t) = \sum_{k=1}^{C-1} b_k \delta(t, t_k),$$

which represents the intensity of birth events. We can also define the function

$$\xi(t) = \sum_{k=1}^{C-1} \xi_k \delta(t, t_k), \quad (4)$$

which represents the intensity of random drift.

We can express the drift intensity function as

$$\xi(t) = \left(\frac{N(t)}{2}\right)^{-1} b(t) = \frac{2b(t)}{N(t)(N(t)-1)}, \quad (5)$$

which is in agreement with the earlier definition (Equation 4). The trajectory of the intensity of random drift, as specified by the drift intensity function  $\xi(t)$  (Equations 4 and 5), takes us a step closer to our goal of deriving an expression for the likelihood function for the sample phylogeny

data. However, in order to express the likelihood function in its most familiar and convenient form, we need to express the trajectory  $\xi(t)$  (and the related trajectories  $N(t)$ , and so on) as functions of time  $s$  measured backwards from the time point at which the sample was collected ( $s = 0$ ). The relationship between the forward time  $t$  (age from conception) and the backwards time  $s$ , is given by

$$s = t_C - t,$$

and hence  $t = t_C - s$ .

So we can represent the backwards time trajectory for population size as the function

$$\tilde{N}(s) = N(t_C - s).$$

Similarly, we can represent the backwards time trajectories for other quantities of interest as follows

$$\tilde{b}(s) = b(t_C - s)$$

and

$$\tilde{\xi}(s) = \xi(t_C - s).$$

Now that we have this definition of the (reverse time) population size function  $\tilde{N}(s)$ , we can express the (reverse time) drift intensity function as

$$\tilde{\xi}(s) = \left(\frac{\tilde{N}(s)}{2}\right)^{-1} \tilde{b}(s) = \frac{2\tilde{b}(s)}{\tilde{N}(s)(\tilde{N}(s)-1)}, \quad (6)$$

which is simply the reverse time version of Equation 5.

Recall that we defined the sequence  $t$  of distinct (forward) times (ages) at which the population size changes, together with the age  $t_C$  at the time of sample collection,  $t = (t_1, t_2, \dots, t_C)$ , where  $0 < t_1 < t_2 < \dots < t_C$ . The same sequence of time points, representing population events, which we have labelled with forward times (ages)  $t_k$ , can also be labelled with reverse times  $s_k = t_C - t_k$ , for  $k = 1, 2, \dots, C-1$ . We now define the sequence  $s$  of distinct reverse times at which the population size changes, together with the time  $s_0 (= t_C)$  at which conception occurred,  $s =$

$(s_0, s_1, s_2, \dots, s_{C-1})$ , where  $s_k = t_C - t_k$ , for each event  $k$ . Therefore, we have  $s_0 > s_1 > s_2 > \dots > s_{C-1} > 0$ . The function  $\tilde{\xi}(s)$  (and the function  $\xi(t)$ ) is completely determined by the sequence pair  $(t, \xi)$ , and also by the (equivalent) sequence pair  $(s, \xi)$ .

When the phylogeny, with (estimated) absolute branch lengths, is the only available data, the likelihood function of the model parameter given the data, is (up to a constant factor) equal to the joint probability density

$$p_n(T(n), T(n-1), \dots, T(2); s, \xi) = \prod_{r=2}^n f_r(T(r)|S(r+1); s, \xi), \quad (7)$$

where

$$S(r) = T(n) + T(n-1) + \dots + T(r),$$

and where each factor

$$f_r(w|s; s, \xi) = \binom{r}{2} \tilde{\xi}(s+w) \cdot R_r(w|s; s, \xi), \quad (8)$$

is the (marginal) probability density of the waiting time to the next coalescent event, starting from time point  $s$ , when  $r$  lines of descent remain (each of which can be traced back from the sample). The function  $\tilde{\xi}(s)$  is the drift intensity at time  $s$  (measured backwards from the time of sample collection). The function

$$R_r(w|s; s, \xi) = \exp\left[-\binom{r}{2} \int_{u=s}^{u=s+w} \tilde{\xi}(u) du\right], \quad (9)$$

gives the probability that the waiting time to the next coalescent event (starting from time point  $s$ , when  $r$  lines of descent remain) is exceeds  $w$ . We could describe  $R_r(w|s; s, \xi)$  as the *reliability* function (or survival function), and interpret  $T(r)$  as a kind of *failure time* (at which one line of descent fails to persist).

Strictly speaking, Equations 8 and 9 represent an approximation which is valid whenever the entire sample phylogeny lies within a time interval throughout which the intensity of random drift  $\tilde{\xi}(s)$  remains small (the effective population size remains large). See refs. <sup>103,104</sup> for derivation of the properties of the (reverse-time) genealogical process.

We want to draw attention to a feature of the likelihood function represented by Equations 7, 8, and 9. From the likelihood function (Equations 7 and 8) it is evident that, while segments (spanning certain time intervals) of the trajectory for the drift intensity (represented variously as a sequence pair  $s, \xi$ , or as a function of time), may constitute an identifiable parameter (when we have a phylogeny on a large enough sample), the trajectory for the population size, and the trajectory for the intensity of birth events, in the absence of additional constraints, are non-identifiable parameters. This is because the population sizes and the counts of birth events do not appear separately in the likelihood function, but only in the particular combination represented by the trajectory for the drift intensity.

In Section 3 below, parameter identifiability is defined more carefully, with some pointers to the literature. We also discuss in more detail the implications of non-identifiability for parameter estimation in our current model. In particular we will discuss how additional constraints on the population trajectory can restore identifiability of the population size, and the intensity of birth events.

### **Parameter estimation and identifiability**

We usually make some further assumptions about the possible trajectories which the population is allowed to follow through time. In the case of a deterministic growth model, we assume that the sequence pair  $(s, \xi)$  of event times and drift intensities belongs to a family of trajectories, in which the individual trajectory is completely determined by a parameter vector  $\phi$ . (Typically this parameter vector is of low dimension.) We say that the family of trajectories is parametrised by  $\phi$ . Here we have in mind models of deterministic exponential growth, where the model parameters include rates of cell division and rates of cell death.

In the case of a stochastic growth model, we assume that the sequence pair  $(s, \xi)$  is drawn from a distribution which belongs to some family of distributions. Within this family of distributions, the specific distribution is completely determined by a parameter vector  $\phi$ . We say that the family of distributions is parametrised by  $\phi$ . Here we have in mind models based on a birth death process, where the model parameters again include rates of cell division and rates of cell death.

In order to emphasise the dependence on the parameter vector  $\phi$ , it is convenient to use the notation

$$\begin{aligned}
 L(\phi|T) &= p_n(T(n), T(n-1), \dots, T(2); \phi) \\
 &= \prod_{r=2}^n f_r(T(r)|S(r+1); \phi),
 \end{aligned}
 \tag{10}$$

for the likelihood function specified by Equations 7 and 8.

We say that a parameter vector  $\phi$  is *non-identifiable* whenever there is a mapping  $\vartheta$  (to a vector of lower dimension) for which the likelihood function  $L(\phi|T)$  depends on the parameter vector  $\phi$  only through  $\theta = \vartheta(\phi)$ . In other words  $\vartheta(\phi_1) = \vartheta(\phi_2)$  implies that  $L(\phi_1|T) = L(\phi_2|T)$ . If there is no such mapping  $\vartheta$ , then we say that the parameter vector  $\phi$  is *identifiable*. When the parameter vector  $\phi$  is *identifiable*, we may also refer to the components of this vector as *identifiable* parameters. See ref. <sup>105</sup> (*non-identifiability* is introduced in Section 3.15, on page 70, and discussed further on pages 72 and 74).

If such a mapping  $\vartheta$  (to a vector of lower dimension) exists (so that  $\phi$  is non-identifiable), then this means (loosely speaking) that from the fixed data  $T$ , we can not learn anything about the unobserved parameter vector  $\phi$ , beyond what we can learn about the (lower dimensional) parameter vector  $\theta$ . We can state this more precisely. First, we can always (leaving aside technical issues and pathological cases) express the prior density  $\pi(\phi)$  for the parameter vector  $\phi$ , in the form

$$\pi(\phi) = \pi(\phi|\theta)\pi(\theta). \tag{11}$$

If  $\phi$  is non-identifiable, and  $\theta = \vartheta(\phi)$  is identifiable, then the posterior density  $\pi(\phi|T)$  of the parameter vector  $\phi$  is of the form

$$\pi(\phi|T) = \pi(\phi|\theta)\pi(\theta|T). \tag{12}$$

As a consequence, we also have

$$\pi(\phi|T, \theta) = \pi(\phi|\theta). \tag{13}$$

This means that if we knew the (lower dimensional) parameter vector  $\theta$ , then the observed data  $T$  would tell us nothing more about the (higher dimensional) parameter vector  $\phi$ .

First we consider a family of models where the population trajectory includes prolonged epochs during which birth events and death events occur equally often, so that the population size remains stable. Then we consider neutral models where the trajectory includes epochs of (deterministic) exponential population growth (Section 5). Finally, we consider birth-death processes, without an upper boundary (Section 6), and with an upper boundary (Section 7) on the population size, and how these stochastic growth models can be approximated by deterministic growth models.

### Epochs of stable effective population size

First we consider a family of models where the population trajectory includes prolonged epochs during which the population size remains stable. Suppose that across the time interval  $[a, b]$ , the population size remains constant at  $N_A$ . In order to maintain a constant population size, the birth rate  $\beta_A$  must be balanced by an equal death rate.

The observed inter-coalescent interval durations  $T(r)$ , which fall within the time interval  $[a, b]$ , contribute factors to the likelihood function which are of the form

$$f_r(T(r)|S(r+1); \phi) = \left(\frac{r}{2}\right)^{\frac{2\beta_A}{N_A}} \cdot \exp\left[-\left(\frac{r}{2}\right)^{\frac{2\beta_A}{N_A}} T(r)\right], \quad (14)$$

where  $\phi = (N_A, \beta_A)$  is the parameter vector of the model.

From the expression on the right hand side of Equation 14, it appears that the only identifiable parameter is the ratio  $\beta_A/N_A$ .

### Epochs of exponential population growth

Now we turn to neutral models where the trajectory includes epochs of exponential population growth. Suppose that the (forward time) estimated trajectory  $\hat{\xi}(t)$  of the drift intensity appears to fit an exponential growth path across the time interval  $[t_A, t_C]$ , where  $t_C$  is the time (age) at which the sample of  $n$  genome-sequenced cells was collected. The estimated trajectory  $\hat{\xi}(t)$  at time  $t$  can be interpreted as a kind of average drift intensity over some interval centred on the time point  $t$ . The (forward time) estimated trajectory is

$$\hat{\xi}(t) = \hat{k} \cdot \exp[\hat{\rho}(t - t_A)], \quad (15)$$

which is based on point estimates  $\hat{\rho}$  (for the growth rate) and  $\hat{k}$  (for the initial drift intensity). Notice that when  $\hat{\rho}$  is positive, the drift intensity declines exponentially, with increasing age  $t$ .

If we measure time backwards from sample collection, then the (reverse time) estimated trajectory  $\hat{\xi}(s)$  of the drift intensity appears to fit an exponential growth path across the time interval  $[0, s_A]$ . The (reverse time) estimated trajectory is

$$\hat{\xi}(s) = \hat{k} \cdot \exp[\hat{\rho}(s_A - s)], \quad (16)$$

where  $s_A = t_C - t_A$  is the time measured backwards from sample collection to the time point at which the epoch of exponential growth began. Notice that when  $\hat{\rho}$  is positive, the drift intensity increases exponentially, with increasing time  $s$ .

There is this one very simple model of population growth, in which births occur at a constant rate  $\lambda$ , and deaths occur at a constant rate  $\nu$ , which results in an exponential trajectory. This is an exceptionally parsimonious explanation for the observed exponential trajectory. If we can accept this parsimonious explanation, then we can set aside the general problem of making inferences about an arbitrary trajectory  $\xi(t)$  for the intensity of random drift (the reciprocal of the effective population size), and restrict our attention to the very specific problem of making inferences about the parameters of the deterministic exponential growth model, or the parameters of the birth death process.

Having observed an (approximately) exponential trajectory for the drift intensity (and its reciprocal, the effective population size), from age  $t_A$ , up to the point of sample collection (at age  $t_C$ ), we have arrived at a parsimonious explanation which we now examine in more detail. The population size has been growing at a constant growth rate  $\rho$ , while the birth rate has remained constant at a value  $\lambda$ , and the death rate has remained constant at a value  $\nu$ , which yields the constant growth rate  $\rho = \lambda - \nu$ . Now we can express the trajectory for the population size  $N(t)$ , forward in time across the epoch of exponential growth (from age  $t_A$  to age  $t_C$ ) as

$$N(t) = N_A \exp[\rho(t - t_A)], \quad (17)$$

while the forward time trajectory for the drift intensity is

$$\xi(t) = \frac{2\lambda}{N_A} \cdot \exp[-\rho(t - t_A)], \quad (18)$$

where  $N_A$  is the size of the ancestral population at age  $t_A$  (when the epoch of exponential growth begins).

We now return to time measured backwards from sample collection. The reverse time trajectory for the population size is

$$\tilde{N}(s) = N_A \exp[\rho(s_A - s)], \quad (19)$$

where  $s_A = t_c - t_A$  is the time measured backwards from sample collection to the time point at which the epoch of exponential growth began. The reverse time trajectory for the drift intensity is

$$\tilde{\xi}(s) = \frac{2\lambda}{N_A} \cdot \exp[-\rho(s_A - s)]. \quad (20)$$

The (marginal) probability density  $f_r(w|s; \phi)$  of the waiting time to the next coalescent event (starting from time point  $s$ , when  $r$  lines of descent remain), is in this case

$$f_r(w|s; \phi) = \left(\frac{r}{2}\right) \frac{2\lambda}{N_A} \cdot \exp[\rho(w + s - s_A)] \cdot R_r(w|s; \phi), \quad (21)$$

where  $\phi = (\lambda, \nu, N_A)$  is the parameter vector of this model, and where

$$R_r(w|s; \phi) = \exp\left[-\left(\frac{r}{2}\right) \frac{2\lambda}{N_A} \cdot \frac{1}{\rho} \exp[\rho(s - s_A)](e^{\rho w} - 1)\right], \quad (22)$$

is the reliability function.

The observed inter-coalescent interval durations  $T(r)$ , which fall within the time interval  $[0, s_A]$  (the epoch of exponential growth), contribute factors to the likelihood function which are of the form

$$\begin{aligned} & f_r(T(r)|S(r+1); \phi) \\ &= \left(\frac{r}{2}\right) \frac{2\lambda}{N_A} \cdot e^{-\rho(U(r)-T(r))} \cdot \exp\left[-\left(\frac{r}{2}\right) \frac{2\lambda}{N_A} \cdot \frac{1}{\rho} e^{-\rho U(r)}(e^{\rho T(r)} - 1)\right], \end{aligned} \quad (23)$$

where  $U(r) = s_A - S(r+1)$ .

The parameter vector of this model is  $\phi = (\lambda, \nu, N_A)$ , where  $\lambda$  is the birth rate,  $\nu$  is the death rate, and  $N_A$  is the size of the ancestral population at the start of the epoch of exponential growth. (This occurs at age  $t_A$ , which precedes sample collection by time interval of duration  $s_0 = t_C - t_A$ .) From the formula for this factor of the likelihood function, it appears that the parameter vector  $\phi$  is *non-identifiable*, while the parameter vector  $\theta = (N_A/\lambda, \rho)$  is *identifiable*. The components of the parameter vector  $\theta$  are the ratio  $N_A/\lambda$ , and the difference  $\rho = \lambda - \nu$  (the population growth rate).

In the special case where the epoch of exponential growth (at constant growth rate  $\rho$ ) extends all the way back to the founding individual (zygote cell), we know  $N_A = 1$ , and we know that (reverse) time  $s_A = s_C$  (age  $t_A = 0$ ) corresponds to the moment of conception. In this special case, the unobserved parameters  $\lambda$  and  $\nu$ , are identifiable. More generally, if the population size at the beginning of the epoch of exponential growth  $N_A$  is known with certainty, then the parameter vector  $\theta = (\lambda, \nu)$  is *identifiable*.

In the case of a sample of single cell genome sequences obtained from blood-derived colonies, from a mouse (or any species with similar HSC dynamics), the parameter  $N_A$  is the size of the ancestral population of HSCs at age  $t_A$  (when the epoch of exponential growth begins); or if the time  $t_A$  is even earlier, then  $N_A$  is the size of the population of embryonic cells existing at this time which are ancestral to the HSCs. Unfortunately we do not have direct observations of the ancestral HSC population size  $N_A$  (at the age  $t_A$  when the epoch of exponential growth begins).

However, we can place some bounds on the value of  $N_A$ . First of all there is an upper bound  $M_A$ , on  $N_A$ , which can be obtained from embryological observations. We know the approximate number of cells in the embryo at age  $t_A$ . If some differentiation has already occurred, we may be able to exclude some cell types as HSC ancestors, and thus perhaps obtain an upper bound  $M_A$  which is somewhat lower than the average total number of cells in a mouse embryo at age  $t_A$ . Secondly, we have a lower bound on  $N_A$ , which we can obtain directly from the phylogeny. This is the number of lines of descent  $n_A$  present on the tree at time  $t_A$ .

### **The linear birth-death process**

A linear birth-death process is a simple stochastic growth model in which birth events and death events occur at constant rates (birth rate  $\lambda$  and death rate  $\nu$ ) per individual (cell) per unit of time (day or year). Therefore the total rate of birth (respectively death) events in the population at each

time point is proportional to the total number of individuals in the population at that time point (hence a *linear* birth-death process. The total size  $N(t)$  of the population at each time point is determined by the (stochastic) sequence of events (births and deaths) up to that time point. For the properties of the linear birth-death process, see ref. <sup>83</sup>, and ref <sup>106</sup>, pages 174–177.

Whenever the population size is not too small, and the growth rate is not too close to zero, the linear birth-death process behaves much like deterministic exponential growth. The trajectory for the population size  $N(t)$  is well approximated by Equation 17, with growth rate  $\rho = \lambda - \nu$ , provided that the birth rate  $\lambda$  exceeds the death rate  $\nu$ , so that  $\rho$  is positive.

In the case of an epoch of stochastic growth (under a linear birth-death process) it is important to bear in mind that the formula for the factors of the likelihood function in Equation 21, is an approximation, which can break-down. A conclusive argument about the identifiability of the model parameters should be based on an exact formula for the likelihood function for the linear birth-death process, when the phylogeny is the only available data.

### **The birth-death process with an upper boundary on population size**

If a mouse lives long enough, we would expect that the propensity of the mouse HSC population to grow exponentially will eventually be checked by the physical constraints on the space available to accommodate the HSC cells within the bone marrow.

In the case of a model where the population undergoes deterministic exponential growth until an upper boundary  $N_B$  on population size is reached, the phylogeny may contain additional information about the time  $T_B$  at which the population first hits the upper boundary  $N_B$ . Such information can be present only if the sample of cells has been collected from the population at a time point after the time  $T_B$ . In this case, the hitting time parameter  $T_B$  occurs in the likelihood function.

In the case of a model where the population undergoes deterministic exponential growth until an upper boundary  $N_B$  population size is reached. The hitting time  $T_B$  is determined by model parameters ( $N_A/N_B$  and  $\rho = \lambda - \nu$ ). Using Equation 17, we can obtain

$$\frac{N_B}{N_A} = \exp[\rho(T_B - t_A)], \quad (24)$$

and therefore

$$T_B = t_A + \frac{1}{\rho} \ln \left( \frac{N_B}{N_A} \right). \quad (25)$$

When the population reaches the upper boundary on population size, the marginal birth rate and the marginal death rate must be equal ( $\delta_B = \beta_B$ ). The parameter vector of the model is now  $\phi = (\lambda, \nu, N_A, N_B, \beta_B)$ .

As usual we inspect the formula for the likelihood function in order to discover which parameters may be identifiable, and which are clearly non-identifiable. The factors of the likelihood function representing the epoch of exponential growth are of the form given in Equation 21, in which the parameter combinations  $\lambda/N_A$  and  $\rho$  appear. The factors of the likelihood function representing the epoch of stable population size are of the form given in Equation 14, in which the parameter combination  $\beta_B/N_B$  appears. We have also seen from Equation 24 that the ratio  $N_B/N_A$  is determined by the parameter  $\rho$  and the hitting time  $T_B$ . The hitting time  $T_B$  is a change point, which also appears in the likelihood function. Therefore, from the formulas for the factors of the likelihood function, it appears that the parameter vector  $\theta = (\rho, \lambda/N_A, \beta_B/N_B, N_B/N_A)$  is identifiable. Notice also that by combining the last three components of  $\theta$ , we obtain

$$\frac{N_B}{N_A} \cdot \frac{\xi_B}{\xi_A} = \frac{N_B}{N_A} \cdot \frac{\beta_B}{N_B} \cdot \frac{N_A}{\lambda} = \frac{\beta_B}{\lambda}.$$

So the ratio  $\beta_B/\lambda$  is also identifiable.

In the special case where  $N_A$  is known for certain, the parameter vector  $\theta = (\lambda, \nu, N_B, \beta_B)$  is identifiable. As already discussed in Section 5, when the epoch of exponential growth (at constant growth rate  $\rho$ ) extends all the way back to the founding individual (zygote cell), we know  $N_A = 1$ . So, in this case, the parameters  $\lambda$ ,  $\nu$ ,  $N_B$ , and  $\beta_B$ , are all identifiable, and amenable to estimation from the phylogeny of a sample.

## REFERENCES

1. Sender, R. & Milo, R. The distribution of cellular turnover in the human body. *Nat Med* **27**, 45–48 (2021).
2. Patel, S. H. *et al.* Lifelong multilineage contribution by embryonic-born blood progenitors. *Nature* **606**, 747–753 (2022).
3. Sun, J. *et al.* Clonal dynamics of native haematopoiesis. *Nature* **514**, 322–327 (2014).
4. Kucinski, I. *et al.* A time- and single-cell-resolved model of murine bone marrow hematopoiesis. *Cell Stem Cell* **31**, 244-259.e10 (2024).
5. Takizawa, H., Regoes, R. R., Boddupalli, C. S., Bonhoeffer, S. & Manz, M. G. Dynamic variation in cycling of hematopoietic stem cells in steady state and inflammation. *J Exp Med* **208**, 273–284 (2011).
6. Munz, C. M. *et al.* Regeneration after blood loss and acute inflammation proceeds without contribution of primitive HSCs. *Blood* **141**, 2483–2492 (2023).
7. Fanti, A.-K. *et al.* Flt3- and Tie2-Cre tracing identifies regeneration in sepsis from multipotent progenitors but not hematopoietic stem cells. *Cell Stem Cell* **30**, 207-218.e7 (2023).
8. Trumpp, A., Essers, M. & Wilson, A. Awakening dormant haematopoietic stem cells. *Nature Reviews Immunology* **10**, 201–209 (2010).
9. Mitchell, E. *et al.* Clonal dynamics of haematopoiesis across the human lifespan. *Nature* **606**, 343–350 (2022).
10. Lee-Six, H. *et al.* Population dynamics of normal human blood inferred from somatic mutations. *Nature* **561**, 473–478 (2018).

11. Osorio, F. G. *et al.* Somatic Mutations Reveal Lineage Relationships and Age-Related Mutagenesis in Human Hematopoiesis. *Cell Reports* **25**, 2308-2316.e4 (2018).
12. Abascal, F. *et al.* Somatic mutation landscapes at single-molecule resolution. *Nature* **593**, 405–410 (2021).
13. Spencer Chapman, M. *et al.* Lineage tracing of human development through somatic mutations. *Nature* **595**, 85–90 (2021).
14. Jaiswal, S. Clonal hematopoiesis and nonhematologic disorders. *Blood* **136**, 1606–1614 (2020).
15. Xie, M. *et al.* Age-related mutations associated with clonal hematopoietic expansion and malignancies. *Nature Medicine* **20**, 1472–1478 (2014).
16. Jaiswal, S. & Ebert, B. L. Clonal hematopoiesis in human aging and disease. *Science* **366**, (2019).
17. Kapadia, C. D. & Goodell, M. A. Tissue mosaicism following stem cell aging: blood as an exemplar. *Nat Aging* **4**, 295–308 (2024).
18. Moore, L. *et al.* The mutational landscape of normal human endometrial epithelium. *Nature* **580**, 640–646 (2020).
19. Martincorena, I. *et al.* High burden and pervasive positive selection of somatic mutations in normal human skin. *Science* **348**, 880–886 (2015).
20. Lawson, A. R. J. *et al.* Extensive heterogeneity in somatic mutation and selection in the human bladder. *Science* **370**, 75–82 (2020).
21. Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus with age.

- Science* **362**, 911–917 (2018).
22. Ng, S. W. K. *et al.* Convergent somatic mutations in metabolism genes in chronic liver disease. *Nature* **598**, 473–478 (2021).
  23. Cagan, A. *et al.* Somatic mutation rates scale with lifespan across mammals. *Nature* **604**, 517–524 (2022).
  24. Yuan, R. *et al.* Genetic coregulation of age of female sexual maturation and lifespan through circulating IGF1 among inbred mouse strains. *Proc Natl Acad Sci U S A* **109**, 8224–8229 (2012).
  25. Chin, D. W. L. *et al.* Aged healthy mice acquire clonal hematopoiesis mutations. *Blood* **139**, 629–634 (2022).
  26. Osawa, M., Hanada, K., Hamada, H. & Nakauchi, H. Long-Term Lymphohematopoietic Reconstitution by a Single CD34-Low/Negative Hematopoietic Stem Cell. *Science* **273**, 242–245 (1996).
  27. Adolfsson, J. *et al.* Upregulation of Flt3 Expression within the Bone Marrow Lin<sup>-</sup>Sca1<sup>+</sup>c-kit<sup>+</sup> Stem Cell Compartment Is Accompanied by Loss of Self-Renewal Capacity. *Immunity* **15**, 659–669 (2001).
  28. Kiel, M. J. *et al.* SLAM family receptors distinguish hematopoietic stem and progenitor cells and reveal endothelial niches for stem cells. *Cell* **121**, 1109–1121 (2005).
  29. Challen, G. A., Boles, N. C., Chambers, S. M. & Goodell, M. A. Distinct Hematopoietic Stem Cell Subtypes Are Differentially Regulated by TGFβ1. *Cell Stem Cell* **6**, 265–278 (2010).
  30. Cabezas-Wallscheid, N. *et al.* Identification of regulatory networks in HSCs and their

immediate progeny via integrated proteome, transcriptome, and DNA methylome analysis.

*Cell Stem Cell* **15**, 507–522 (2014).

31. Pietras, E. M. *et al.* Functionally Distinct Subsets of Lineage-Biased Multipotent Progenitors Control Blood Production in Normal and Regenerative Conditions. *Cell Stem Cell* **17**, 35–46 (2015).
32. Sun, J. *et al.* Clonal dynamics of native haematopoiesis. *Nature* **514**, 322–327 (2014).
33. Busch, K. *et al.* Fundamental properties of unperturbed haematopoiesis from stem cells in vivo. *Nature* **518**, 542–546 (2015).
34. Williams, N. *et al.* Life histories of myeloproliferative neoplasms inferred from phylogenies. *Nature* **602**, 162–168 (2022).
35. Machado, H. E. *et al.* Diverse mutational landscapes in human lymphocytes. *Nature* **608**, 724–732 (2022).
36. Coorens, T. H. H. *et al.* Inherent mosaicism and extensive mutation of human placentas. *Nature* **592**, 80–85 (2021).
37. Bryder, D., Rossi, D. J. & Weissman, I. L. Hematopoietic stem cells: the paradigmatic tissue-specific stem cell. *Am J Pathol* **169**, 338–346 (2006).
38. Qin, P. *et al.* Integrated decoding hematopoiesis and leukemogenesis using single-cell sequencing and its medical implication. *Cell Discov* **7**, 2 (2021).
39. de Haan, G. & Van Zant, G. Dynamic Changes in Mouse Hematopoietic Stem Cell Numbers During Aging. *Blood* **93**, 3294–3301 (1999).
40. Morrison, S. J., Wandycz, A. M., Akashi, K., Globerson, A. & Weissman, I. L. The aging of

- hematopoietic stem cells. *Nat Med* **2**, 1011–1016 (1996).
41. Beaumont, M. A., Zhang, W. & Balding, D. J. Approximate Bayesian computation in population genetics. *Genetics* **162**, 2025–2035 (2002).
  42. Koonin, E. V. Splendor and misery of adaptation, or the importance of neutral null for understanding evolution. *BMC Biology* **14**, 114 (2016).
  43. Challen, G. A. & Goodell, M. A. Clonal hematopoiesis: mechanisms driving dominance of stem cell clones. *Blood* **136**, 1590–1598 (2020).
  44. Yoshizato, T. *et al.* Somatic Mutations and Clonal Hematopoiesis in Aplastic Anemia. *New England Journal of Medicine* **373**, 35–47 (2015).
  45. King, K. Y., Huang, Y., Nakada, D. & Goodell, M. A. Environmental influences on clonal hematopoiesis. *Experimental Hematology* **83**, 66–73 (2020).
  46. Florez, M. A. *et al.* Clonal hematopoiesis: Mutation-specific adaptation to environmental change. *Cell Stem Cell* **29**, 882–904 (2022).
  47. Coombs, C. C. *et al.* Therapy-Related Clonal Hematopoiesis in Patients with Non-hematologic Cancers Is Common and Associated with Adverse Clinical Outcomes. *Cell Stem Cell* **21**, 374-382.e4 (2017).
  48. Bolton, K. L. *et al.* Cancer therapy shapes the fitness landscape of clonal hematopoiesis. *Nat Genet* **52**, 1219–1226 (2020).
  49. Wong, T. N. *et al.* Role of TP53 mutations in the origin and evolution of therapy-related acute myeloid leukaemia. *Nature* **518**, 552–555 (2015).
  50. Beura, L. K. *et al.* Normalizing the environment recapitulates adult human immune traits

- in laboratory mice. *Nature* **532**, 512–516 (2016).
51. Camell, C. D. *et al.* Senolytics reduce coronavirus-related mortality in old mice. *Science* **373**, (2021).
  52. Matatall, K. A. *et al.* Chronic Infection Depletes Hematopoietic Stem Cells through Stress-Induced Terminal Differentiation. *Cell Rep* **17**, 2584–2595 (2016).
  53. Watson, C. J. *et al.* The evolutionary dynamics and fitness landscape of clonal hematopoiesis. *Science* **367**, 1449–1454 (2020).
  54. Kimura, M. & Ohta, T. The Average Number of Generations until Fixation of a Mutant Gene in a Finite Population. *Genetics* **61**, 763–771 (1969).
  55. Challen, G. A., Pietras, E. M., Wallscheid, N. C. & Signer, R. A. J. Simplified murine multipotent progenitor isolation scheme: Establishing a consensus approach for multipotent progenitor identification. *Experimental Hematology* **104**, 55–63 (2021).
  56. Sheikh, B. N. *et al.* MOZ (KAT6A) is essential for the maintenance of classically defined adult hematopoietic stem cells. *Blood* **128**, 2307–2318 (2016).
  57. Kobayashi, M. *et al.* HSC-independent definitive hematopoiesis persists into adult life. *Cell Reports* **42**, 112239 (2023).
  58. Abkowitz, J. L., Catlin, S. N., McCallie, M. T. & Gutter, P. Evidence that the number of hematopoietic stem cells per animal is conserved in mammals. *Blood* **100**, 2665–2667 (2002).
  59. Dykstra, B. *et al.* Long-term propagation of distinct hematopoietic differentiation programs in vivo. *Cell Stem Cell* **1**, 218–229 (2007).
  60. Gros, P. & Casanova, J.-L. Reconciling Mouse and Human Immunology at the Altar of

Genetics. *Annu Rev Immunol* **41**, 39–71 (2023).

61. Lindsay, S. J., Rahbari, R., Kaplanis, J., Keane, T. & Hurles, M. E. Similarities and differences in patterns of germline mutation between mice and humans. *Nat Commun* **10**, 4053 (2019).
62. Bergeron, L. A. *et al.* Evolution of the germline mutation rate across vertebrates. *Nature* **615**, 285–291 (2023).
63. Gould, S. J., Lewontin, R. C., Maynard Smith, J. & Holliday, R. The spandrels of San Marco and the Panglossian paradigm: a critique of the adaptationist programme. *Proceedings of the Royal Society of London. Series B. Biological Sciences* **205**, 581–598 (1997).
64. Brayton, C. F., Treuting, P. M. & Ward, J. M. Pathobiology of aging mice and GEM: background strains and experimental design. *Vet Pathol* **49**, 85–105 (2012).
65. Shepherd, B. E. *et al.* Hematopoietic stem-cell behavior in nonhuman primates. *Blood* **110**, 1806–1813 (2007).
66. Koelle, S. J. *et al.* Quantitative stability of hematopoietic stem and progenitor cell clonal output in rhesus macaques receiving transplants. *Blood* **129**, 1448–1457 (2017).
67. Shin, T.-H. *et al.* A macaque clonal hematopoiesis model demonstrates expansion of TET2-disrupted clones and utility for testing interventions. *Blood* **140**, 1774–1789 (2022).
68. Yu, K.-R. *et al.* The impact of aging on primate hematopoiesis as interrogated by clonal tracking. *Blood* **131**, 1195–1205 (2018).
69. Hsu, J. I. *et al.* PPM1D Mutations Drive Clonal Hematopoiesis in Response to Cytotoxic Chemotherapy. *Cell Stem Cell* **23**, 700-713.e6 (2018).

70. Heyde, A. *et al.* Increased stem cell proliferation in atherosclerosis accelerates clonal hematopoiesis. *Cell* **184**, 1348-1361.e22 (2021).
71. Meisel, M. *et al.* Microbial signals drive pre-leukaemic myeloproliferation in a Tet2-deficient host. *Nature* **557**, 580–584 (2018).
72. Hormaechea-Agulla, D. *et al.* Chronic infection drives Dnmt3a-loss-of-function clonal hematopoiesis via IFN $\gamma$  signaling. *Cell Stem Cell* **28**, 1428-1442.e6 (2021).
73. Cheshier, S. H., Morrison, S. J., Liao, X. & Weissman, I. L. In vivo proliferation and cell cycle kinetics of long-term self-renewing hematopoietic stem cells. *Proceedings of the National Academy of Sciences* **96**, 3120–3125 (1999).
74. Wilson, A. *et al.* Hematopoietic Stem Cells Reversibly Switch from Dormancy to Self-Renewal during Homeostasis and Repair. *Cell* **135**, 1118–1129 (2008).
75. Bernitz, J. M., Kim, H. S., MacArthur, B., Sieburg, H. & Moore, K. Hematopoietic Stem Cells Count and Remember Self-Renewal Divisions. *Cell* **167**, 1296-1309.e10 (2016).
76. Chen, J., Astle, C. M. & Harrison, D. E. Genetic regulation of primitive hematopoietic stem cell senescence. *Exp Hematol* **28**, 442–450 (2000).
77. Ahmad, A. *et al.* ERCC1-XPF endonuclease facilitates DNA double-strand break repair. *Mol Cell Biol* **28**, 5082–5092 (2008).
78. Nadon, N. L., Strong, R., Miller, R. A. & Harrison, D. E. NIA Interventions Testing Program: Investigating Putative Aging Intervention Agents in a Genetically Heterogeneous Mouse Model. *EBioMedicine* **21**, 3–4 (2017).
79. Bowie, M. B. *et al.* Hematopoietic stem cells proliferate until after birth and show a

- reversible phase-specific engraftment defect. *J Clin Invest* **116**, 2808–2816 (2006).
80. Ellis, P. *et al.* Reliable detection of somatic mutations in solid tissues by laser-capture microdissection and low-input DNA sequencing. *Nat Protoc* **16**, 841–871 (2021).
81. Jones, D. *et al.* cgpCaVEManWrapper: Simple Execution of CaVEMan in Order to Detect Somatic Single Nucleotide Variants in NGS Data. *Curr Protoc Bioinformatics* **56**, 15.10.1-15.10.18 (2016).
82. Ye, K., Schulz, M. H., Long, Q., Apweiler, R. & Ning, Z. Pindel: a pattern growth approach to detect break points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics* **25**, 2865–2871 (2009).
83. Kendall, D. G. Stochastic Processes and Population Growth. *Journal of the Royal Statistical Society. Series B (Methodological)* **11**, 230–282 (1949).
84. Challen, G. A., Boles, N., Lin, K. K.-Y. & Goodell, M. A. Mouse hematopoietic stem cell identification and analysis. *Cytometry. Part A : the journal of the International Society for Analytical Cytology* **75**, 14–24 (2009).
85. Lee-Six, H. *et al.* The landscape of somatic mutation in normal colorectal epithelial cells. *Nature* **574**, 532–537 (2019).
86. Durand, J.-B., Goncalves, P. & Guedon, Y. Computational methods for hidden Markov tree models-an application to wavelet trees. *IEEE Transactions on Signal Processing* **52**, 2551–2560 (2004).
87. Kennedy, S. R. *et al.* Detecting ultralow-frequency mutations by Duplex Sequencing. *Nature protocols* **9**, 2586–606 (2014).

88. Lai, Z. *et al.* VarDict: a novel and versatile variant caller for next-generation sequencing in cancer research. *Nucleic Acids Res* **44**, e108 (2016).
89. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biology* **17**, 122 (2016).
90. Costello, M. *et al.* Discovery and characterization of artifactual mutations in deep coverage targeted capture sequencing data due to oxidative DNA damage during sample preparation. *Nucleic Acids Res* **41**, e67 (2013).
91. Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer. *Nature* **578**, 94–101 (2020).
92. Feng, C. G., Weksberg, D. C., Taylor, G. A., Sher, A. & Goodell, M. A. The p47 GTPase Lrg-47 (*Irgm1*) links host defense and hematopoietic stem cell proliferation. *Cell Stem Cell* **2**, 83–89 (2008).
93. Lerner, C. & Harrison, D. E. 5-Fluorouracil spares hemopoietic stem cells responsible for long-term repopulation. *Exp Hematol* **18**, 114–118 (1990).
94. Dong, S. *et al.* Chaperone-mediated autophagy sustains haematopoietic stem-cell function. *Nature* **591**, 117–123 (2021).
95. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell* **171**, 1029-1041.e21 (2017).
96. Campbell, P. *et al.* Clonal dynamics after allogeneic haematopoietic cell transplantation using genome-wide somatic mutations. Preprint at <https://doi.org/10.21203/rs.3.rs-2868644/v1> (2023).

97. Poon, G. Y. P., Watson, C. J., Fisher, D. S. & Blundell, J. R. Synonymous mutations reveal genome-wide levels of positive selection in healthy tissues. *Nat Genet* **53**, 1597–1605 (2021).
98. Flurkey, K., M. Curren, J. & Harrison, D. E. Mouse Models in Aging Research. in *The Mouse in Biomedical Research (Second Edition)* (eds. Fox, J. G. et al.) 637–672 (Academic Press, Burlington, 2007). doi:10.1016/B978-012369454-6/50074-1.
99. Pritchard, J. K., Seielstad, M. T., Perez-Lezaun, A. & Feldman, M. W. Population growth of human Y chromosomes: a study of Y chromosome microsatellites. *Molecular Biology and Evolution* **16**, 1791–1798 (1999).
100. Csilléry, K., François, O. & Blum, M. G. abc: an R package for approximate Bayesian computation (ABC). *Methods in ecology and evolution* **3**, 475–479 (2012).
101. Lan, S., Palacios, J. A., Karcher, M., Minin, V. N. & Shahbaba, B. An efficient Bayesian inference framework for coalescent-based nonparametric phylodynamics. *Bioinformatics* **31**, 3282–3289 (2015).
102. Moran, P. A. P. Random processes in genetics. *Proceedings of the Cambridge Philosophical Society* **54**, 60–71 (1958).
103. Kingman, J. F. C. On the Genealogy of Large Populations. *J. Appl. Probab.* **19A**, 27–43 (1982).
104. Griffiths, R. C. & Tavaré, S. Sampling theory for neutral alleles in a varying environment. *Philosophical Transactions of the Royal Society, London, Series B* **344**, 403–410 (1994).
105. O’Hagan, A. & Forster, J. *Bayesian Inference*. vol. 2B (Arnold, London, UK, 2004).
106. Moran, P. A. P. *An Introduction to Probability Theory*. (Oxford University Press, Oxford,

UK, 1968).