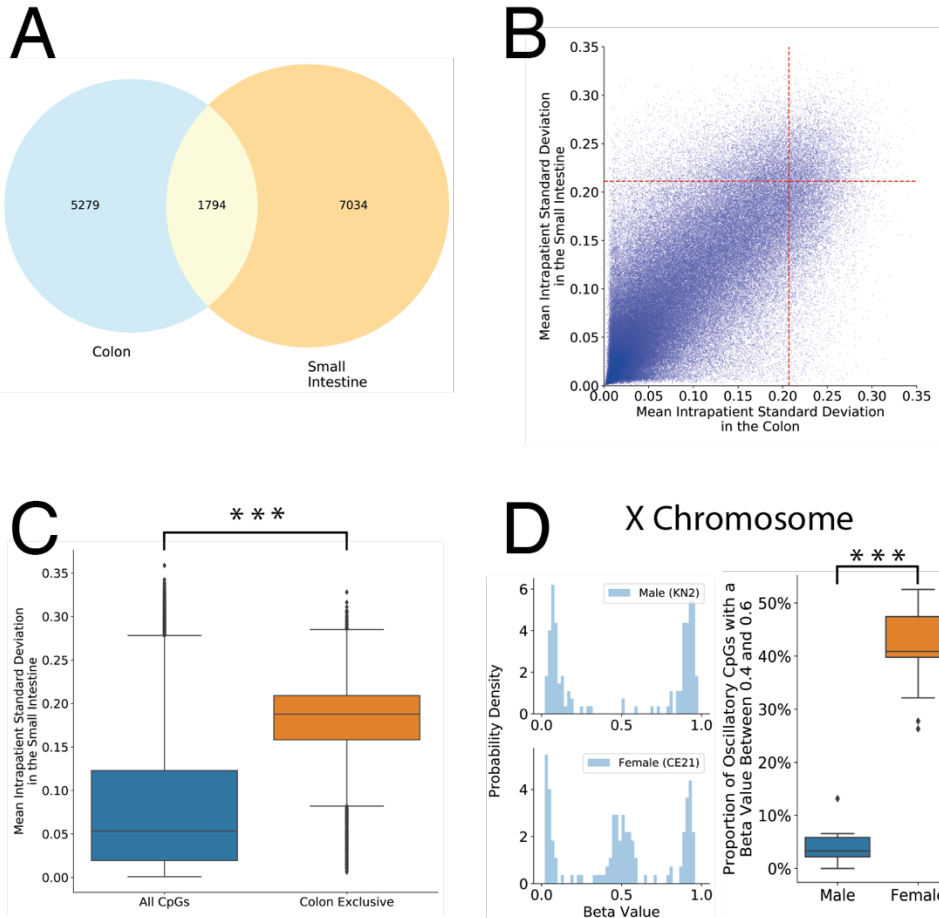**Supplementary information**

# Fluctuating methylation clocks for cell lineage tracing at high temporal resolution in human tissues

In the format provided by the authors and unedited

# 1  Supplementary Information

2

# 3  Supplementary Figures



**Supplementary Figure 1: Additional analysis of fCpG identification process**

**A:** Venn diagram showing the overlap of CpG loci identified as fCpGs in the colon and the small intestine. **B:** Scatter density plot (with the density plotted on the log-scale) of the heterogeneity metric (mean intra-patient standard deviation) of CpGs in the colon and small intestine, with the cutoff of the top 5% most heterogenous loci indicated in red. **C:** Comparison of the heterogeneity metric of the colon exclusive fCpG loci (i.e. those identified in the colon but not the small intestine) to all type II CpGs, within the small intestine samples (center line, median; box limits, upper and lower quartiles; whiskers, 1.5 IQR). The colon exclusive CpG loci are significantly more variable ($p < 2 \times 10^{-16}$, two-sided Mann Whitney U test). **D:** An extension of the fCpG identification process to CpG loci located on the X chromosome. (left) We present example methylation distributions for these X-chromosome fCpG loci for a male and a female crypt, confirming the predictions from theory that the male crypts lack the peak at 50% as they contain only a single copy. (right) To test whether this relationship holds in general, we compare the proportion of fCpG's on the X chromosome with an intermediate beta value ($0.4 \leq \beta \leq 0.6$) between all colon crypts from males and females (center line, median; box limits, upper and lower quartiles; whiskers, 1.5 IQR), confirming that males have a significantly lower probability mass near 50% ($p = 3.9 \times 10^{-6}$, two-sided Mann Whitney U test).

4

## Evidence for fluctuating human CpG sites

6

7   Oscillatory DNA methylation has been previously documented in cell lines systems where
8   changes in DNA methylation have been directly observed and correlated with changes in
9   gene expression and developmental cell fates[1–4]. These active changes in DNA methylation
10  are modulated by Tet and Dnmt family enzymes, and occur over hours or a few days. By
11  contrast, the erratic fluctuations at fCpG loci outlined in this paper occur over years and in
12  adult human tissues. For these fCpG loci, the proposed mechanism is more likely stochastic
13  replication errors, where DNA methylation and demethylation are equally likely.

14

15  Such stochastic, neutral changes in DNA methylation are more likely to occur in non-genic
16  or non-regulatory CpG sites rather than in expressed genes that could alter cell phenotypes.
17  The identified fCpG loci were relatively enriched for non-genic CpG sites (Fig. 2B) and were
18  about three times less likely to be in promoter regions that largely control gene expression.
19  Moreover, genes with promoter fCpG sites had significantly lower expression than genes
20  without promoter fCpG loci ($p < 0.001$ Welch's t-test performed upon the log-transformed
21  data, Fig. S2A). We note that none of the genes with promoter fCpG sites had an expression
22  greater than 10 TPM, a typical cutoff for "intermediate" expression. Hence, fCpG methylation
23  variation is unlikely to be associated with cell differentiation or actively confer significant
24  changes in phenotype.

25

26  Although direct serial observations of the methylation and demethylation of fCpG loci in the
27  same human colon crypt are impractical, the proposed stochastic mechanism has testable
28  predictions. We illustrate some of these predictions by comparing methylation extremes,
29  defined as "0" if the methylation is less than 0.2 and "1" if the methylation is greater than 0.8
30  (Fig. S2B). The clonal crypts have hundreds of fCpG sites with values of either "0" or "1",
31  which should effectively serve as lineage barcodes over short time intervals. (CpG sites with
32  methylation between 0.2 and 0.8 are not compared.) For example, two identical samples
33  should have the same barcode. With increasing time, the fluctuations should randomize the
34  fCpG methylation, and the probability of matching should be about 0.5, or the same as
35  flipping a coin.

36

37  We compared FMCs between the top and bottoms of colon crypts (Fig. S2C). Crypt stem
38  cells are located in the bottoms of crypts and cell differentiation occurs at crypt bottoms.
39  Differentiated cells migrate upwards and the tops of crypts are differentiated cells that are
40  lost within a week. As expected for a relatively short interval of a week, there was minimal
41  switching between 0 to 1 or 1 to 0 between the top and bottom of the same crypt (Fig. S2D).
42  This data also indicates the fluctuations in methylation pattern are not associated with cell
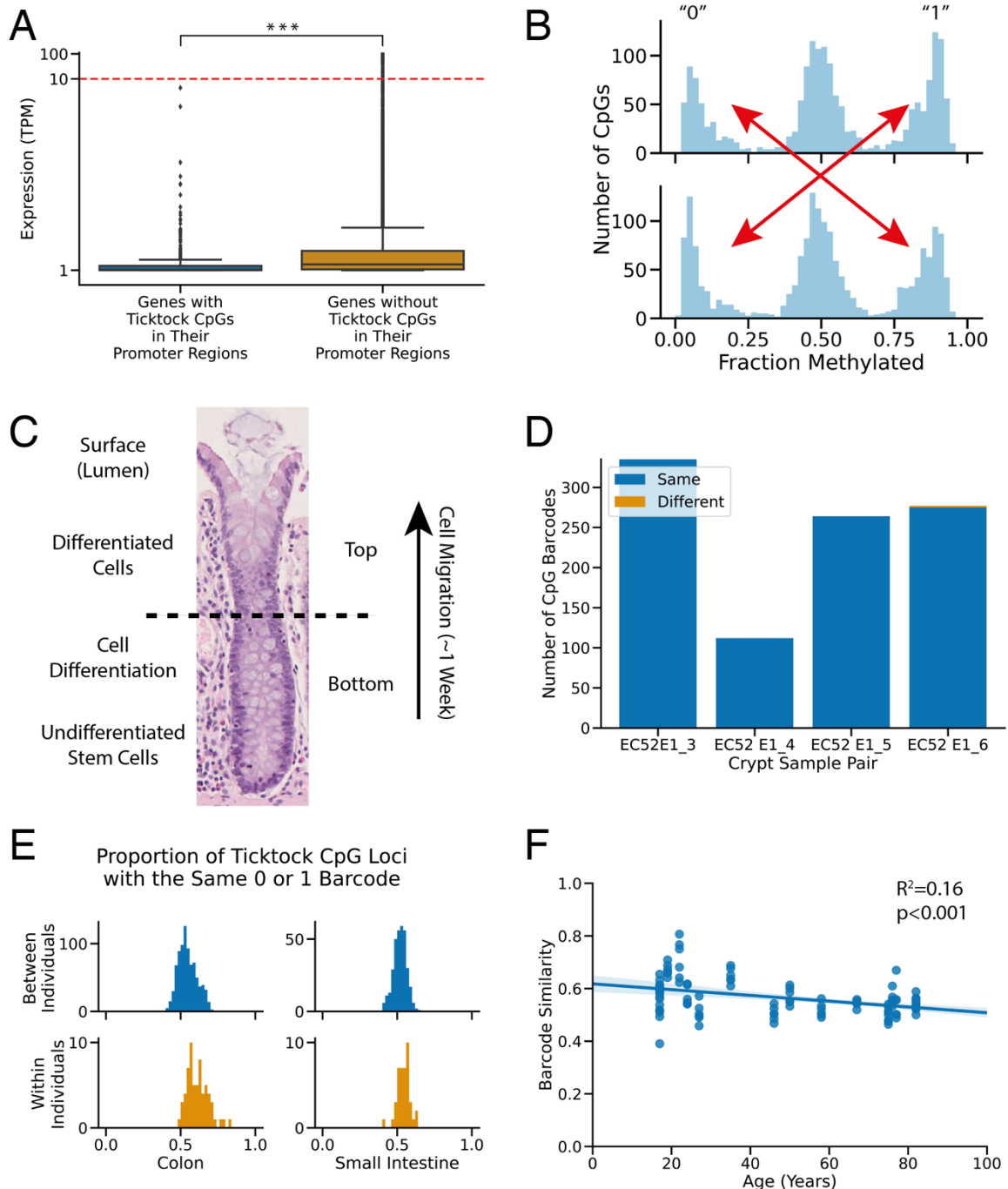43  differentiation.

44

45  By contrast, essentially a random distribution was found when FMCs were compared
46  between adult crypts from different individuals or within an individual, with an average
47  identity of barcodes of about 0.5 (Fig. S2E). The youngest individual was 17 years old,
48  providing many years for methylation fluctuations to randomize between crypts that last
49  share a common ancestor before birth. Consistent with a requirement for time for
50  randomization, there was a trend for increasingly more random barcodes with age, with
51  younger individuals tending to have more similar fCpG barcodes (Fig. S2F).

52

53    In summary, although we are unable to directly observe changes in the fCpG loci
54    methylation, these observations are consistent with fluctuations in methylation at CpG sites
55    with change occurring over many years. The fCpG sites do not switch during cell
56    differentiation occurring over a week, and are rarely found in promoter regions or in highly
57    expressed genes. Consistent with stochastic fluctuations that occur independently in
58    different cells, FMC barcodes are essentially randomized between adult crypts from different
59    individuals, and become increasingly different with aging within an individual. These findings
60    are more consistent with dynamic methylation fluctuations that randomly change between 0,
61    0.5 and 1.0, rather than active regulatory oscillations or static random states. The FMCs
62    distribution reflect the population structure of their cells. In stable polyclonal populations, the
63    fCpG sites are asynchronous between cells and average methylation will be around 0.5, with
64    narrow distributions as seen in normal whole blood. However, fCpG sites become
65    synchronized in clonal populations and exhibit "W" shaped distributions, as seen in normal
66    human crypts, endometrial glands and clonal hematopoiesis.

67

**Supplementary Figure 2: Additional evidence supporting the existence of fluctuating CpG loci**

**A:** RNA expression is on average 3 times lower ($p = 3.0 \times 10^{-28}$ two-sided Welch's t-test performed upon the log-transformed data) for genes with fCpG sites in their promoter regions compared to genes without promotor fCpG loci in the colon (RNAseq single cell data from GSE132257 - center line, median; box limits, upper and lower quartiles; whiskers, 1.5 IQR). **B:** FMC methylation reversibly fluctuates between 0, 50 and 100% methylated. Two samples can be compared by counting the number of fCpG loci that switch from 0 to 1 or vice versa. **C:** Colon crypts have small numbers of undifferentiated stem cells at their bases/bottoms. Cell differentiation occurs near crypt bottoms and fully differentiated cells migrate upwards and are lost within a week or two. **D:** Comparisons

between four top and bottom halves of the same crypts reveals FMC barcode methylation (0 or 1) is nearly always identical (986/988 barcode comparisons). **E:** Comparisons of colon and SI crypt FMC barcode methylation between and within individuals. The probability of another crypt having the same 0 or 1 barcode is approximately 0.5, or essentially random. **F:** Performing pairwise comparisons of the proportion of fCpG loci with the same 0 or 1 barcode, there was a slight trend towards increasingly random barcodes with age. Confidence band was calculated via bootstrapping and represents 95% confidence intervals.

68

## Derivation of model describing methylation within the stem cell niche

70

Consider a single fCpG locus within a fixed population of $S$ stem cells. Within each stem cell, the locus is assumed to be diploid, so each stem cell contains 2 alleles at this locus. In this way, there are 3 possible "states" for a given stem cell, (i) neither allele methylated, (ii) both alleles methylated, (iii) or one allele methylated whilst the other is unmethylated. We are interested in the population methylation level, so we assume that the population is well-mixed, which allows us to characterize the system using just 2 state variables: $k$ – the number of stem cells containing a single methylated allele, and $m$ – the number of stem cells containing 2 methylated alleles. The number of stem cells containing 0 methylated alleles is then given by $S - m - k$.

80

The states are constrained such that $0 \leq k, m \leq S$ and $k + m \leq S$, which allows us to calculate the total number of possible states by considering all possible combinations of $k$ and $m$. If we first consider the case when $m = 0$, then $k$ can take any value between $0$ and $S$ giving a total of $S + 1$ states. If we next consider the $m = 1$ case, then $k$ can take any value between $0$ and $S - 1$, a total of $S$ states. We can continue in this fashion for each of the $S + 1$ possible states for $m$, such that the total number of states is

$$\sum_{m=0}^{S} S + 1 - m = \frac{1}{2}(S + 1)(S + 2) \tag{1}$$

87

We assume that there are three possible processes that can change the population methylation level $(k, m) \rightarrow (k', m')$:

(1) an unmethylated allele spontaneously becoming methylated (which, for a single unmethylated CpG locus, occurs at a rate $\mu$ per allele per stem cell per year)
(2) a methylated allele spontaneously becoming unmethylated (which occurs at a rate $\gamma$ per allele per stem cell per year)
(3) one stem cell replacing one of the other $S - 1$ stem cells (which occurs at a rate $\lambda$ per stem cell per year).

To formulate a system of differential equations that characterize the rates at which the population methylation changes, we first consider the probability the system in state $(k, m)$ at time $t$ transitions to state $(k', m')$ within the time $t + \delta t$ (where we assume $\delta t$ is small enough that the probability of a "double-jump" is negligible).

100

If we are in state $(k, m)$, then the probability that one of the $k$ heterozygous methylated stem cells becomes unmethylated (via process (2)) in a time period $\delta t$ is:

$$\mathrm{P}\big((k, m) \rightarrow (k - 1, m)\big) = k\gamma\delta t \tag{2}$$

103

104 And the probability that the one of the $m$ homozygous methylated stem cells (representing
105 2m methylated alleles) undergoes process (2) is:

$$P\big((k,m) \to (k+1, m-1)\big) = 2m\gamma\delta t \qquad\qquad 3$$

106

107 Similarly, considering methylation (process (1)), there are a total of $2S - k - 2m$
108 unmethylated alleles where the process could occur. The probability that one of the
109 homozygous S-k-m unmethylated stem cells becomes heterozygous is:

$$P\big((k,m) \to (k+1, m)\big) = 2(S - k - m)\mu\delta t \qquad\qquad 4$$

110

111 And the probability that one of the heterozygous methylated stem cells becomes
112 homozygous methylated is:

$$P\big((k,m) \to (k-1, m+1)\big) = k\mu\delta t \qquad\qquad 5$$

113

114 Let us now consider the replacement process. In a time period $\delta t$ the probability that a
115 replacement occurs is $S\lambda\delta t$. There are $S(S-1)$ possible replacements: $S$ possible cells that
116 can expand, which must replace any of the $S-1$ other cells. To go from state $(k,m)$ to a
117 different state $(k',m')$, we require the expanding stem cell to replace a cell with a different
118 methylation status. Therefore, the probability of the transition $(k,m) \to (k',m')$ is equal to the
119 probability that any of the cells replaces another, $S\lambda\delta t$, multiplied by the number of ways that
120 particular transition could occur, and normalized by the total possible number of transitions.

121

122 To give a concrete example, consider the stem the cell niche illustrated in Figure 1C, which
123 contains 5 stem cells and is initially in the state $(k = 3, m = 1)$. There are a total of $5 * 4 =$
124 20 possible replacements. Clearly, if one of the heterozygous stem cells replaces another of
125 the heterozygous stem cells, the population methylation level will not change. To jump to the
126 state $(k = 3, m = 2)$ as illustrated in Figure 1C, only one replacement (the homozygous
127 methylated stem cell replacing the homozygous unmethylated stem cell) allows the specified
128 jump, hence the probability of the jump $(3,1) \to (3,2)$ in the time $\delta t$ is $\frac{1}{5*4}5\lambda\delta t = \frac{1}{4}\lambda\delta t$. To
129 generalise this, the fraction of possible transitions that give rise to the particular jump
130 $(k,m) \to (k',m')$ is equal to the multiplicity of the expanding cell multiplied by the multiplicity
131 of the replaced cell, divided by $S(S-1)$.

132

133 Applying the same logic, we can derive the probability of all six possible state transitions via
134 replacement:

$$P\big((k,m) \to (k, m+1)\big) = \frac{m(S - m - k)\lambda\delta t}{S - 1} \qquad\qquad 6$$

$$P\big((k,m) \to (k+1, m)\big) = \frac{k(S - m - k)\lambda\delta t}{S - 1} \qquad\qquad 7$$

$$P\big((k,m) \to (k-1, m+1)\big) = \frac{km\lambda\delta t}{S - 1} \qquad\qquad 8$$

$$P\big((k,m) \to (k+1, m-1)\big) = \frac{km\lambda\delta t}{S - 1} \qquad\qquad 9$$

$$P\big((k,m) \to (k,m-1)\big) = \frac{m(S-m-k)\lambda\delta t}{S-1} \qquad\qquad 10$$

$$P\big((k,m) \to (k-1,m)\big) = \frac{k(S-m-k)\lambda\delta t}{S-1} \qquad\qquad 11$$

135

136 The methylation switching and replacement processes that we have considered separately
137 above are independent, allowing us to simply add the probabilities together (again,
138 assuming that $\delta t$ is small enough that the probability of two processes occurring in $\delta t$ is
139 negligible) to find the total probability that a given transition would occur:

$$P\big((k,m) \to (k,m+1)\big) = \frac{m(S-m-k)\lambda\delta t}{S-1} \qquad\qquad 12$$

$$P\big((k,m) \to (k+1,m)\big) = \frac{k(S-m-k)\lambda\delta t}{S-1} + 2(S-m-k)\mu\delta t \qquad\qquad 13$$

$$P\big((k,m) \to (k-1,m+1)\big) = \frac{km\lambda\delta t}{S-1} + k\mu\delta t \qquad\qquad 14$$

$$P\big((k,m) \to (k+1,m-1)\big) = \frac{km\lambda\delta t}{S-1} + 2m\gamma\delta t \qquad\qquad 15$$

$$P\big((k,m) \to (k,m-1)\big) = \frac{m(S-m-k)\lambda\delta t}{S-1} \qquad\qquad 16$$

$$P\big((k,m) \to (k-1,m)\big) = \frac{k(S-m-k)\lambda\delta t}{S-1} + k\gamma\delta t \qquad\qquad 17$$

140

141 We have considered above the transitions leading "out" of the state $(k,m)$ into adjacent
142 states $(k',m')$. However, we can also consider the jumps "into" the state $(k,m)$ from the
143 adjacent states $(k',m')$:

$$P\big((k,m-1) \to (k,m)\big) = \frac{(m-1)(S-(m-1)-k)\lambda\delta t}{S-1} \qquad\qquad 18$$

$$P\big((k-1,m) \to (k,m)\big) = \frac{(k-1)(S-m-(k-1))\lambda\delta t}{S-1} + 2(S-m-(k-1))\mu\delta t \qquad\qquad 19$$

$$P\big((k+1,m-1) \to (k,m)\big) = \frac{(k+1)(m-1)\lambda\delta t}{S-1} + (k+1)\mu\delta t \qquad\qquad 20$$

$$P\big((k-1,m+1) \to (k,m)\big) = \frac{(k-1)(m+1)\lambda\delta t}{S-1} + 2(m+1)\gamma\delta t \qquad\qquad 21$$

$$P\big((k,m+1) \to (k,m)\big) = \frac{m(S-(m+1)-k)\lambda\delta t}{S-1} \qquad\qquad 22$$

$$P\big((k+1,m) \to (k,m)\big) = \frac{(k+1)(S-m-(k+1))\lambda\delta t}{S-1} + (k+1)\gamma\delta t \qquad\qquad 23$$

144 So far, we have considered the probability that the system changes from state $(k,m)$ to state
145 $(k',m')$ within time $\delta t$. However, we primarily want to know the probability of the system
146 being in state $(k,m)$ at time $t$, $P(k,m;t)$, and how this changes over time. For the system to
147 be in state $(k,m)$ at time $t+\delta t$, either (i) the system must have been in state $(k,m)$ at time $t$
148 and has not transitioned out of the state (which is equal to 1 minus the probability of
149 transitioning to an adjacent state, defined by equations 12-17), (ii) or the system was in a

150 different (adjacent) state $(k', m')$ at time $t$ and has transitioned into the state $(k, m)$ in time $\delta t$
151 (defined by equations 18-23):

$$
\begin{aligned}
\text{P}(k, m; t + \delta t) = \text{P}(k, m; t)\left(1 - \sum_{k', m'} \text{P}\big((k, m) \to (k', m')\big)\right) \\
+ \sum_{k', m'} \text{P}(k', m'; t)\text{P}\big((k', m') \to (k, m)\big)
\end{aligned}
\tag{24}
$$

152

153 We can rearrange equation 24, factoring out the common factor of $\delta t$ in the $\text{P}\big((k', m') \to$
154 $(k, m)\big)$ terms and take the limit $\delta t \to 0$:

$$
\begin{aligned}
\frac{d\text{P}(k, m; t)}{dt} &= \lim_{\delta t \to 0}\left(\frac{\text{P}(k, m; t + \delta t) - \text{P}(k, m; t)}{\delta t}\right) \\
&= \sum_{k', m'} \text{P}(k', m'; t)\frac{\text{P}\big((k', m') \to (k, m)\big)}{\delta t} \\
&\quad - \text{P}(k, m; t)\frac{\text{P}\big((k, m) \to (k', m')\big)}{\delta t}
\end{aligned}
\tag{25}
$$

155

156 The sum over equation 12-17 in the final term evaluates as:

$$
\begin{aligned}
\sum_{k', m'} \frac{\text{P}\big((k, m) \to (k', m')\big)}{\delta t} &= \big(k(S - k) + m(S - k - m)\big)\frac{2\lambda}{S - 1} \\
&\quad + \big(2S - (k + 2m)\big)\mu + (k + 2m)\gamma
\end{aligned}
\tag{26}
$$

157 Due to the constraints on $k$ and $m$, we consider the differential equations for $(k = 0, m = 0)$,
158 $(k = S, m = 0)$ and $(k = 0, m = S)$ separately. Combining equations 25, 26 and 18-23, we
159 derive the following set of differential equations:

$$
\frac{d\text{P}(0,0|\lambda, \mu, \gamma; t)}{dt} = (\lambda + \gamma)\text{P}(1,0|\lambda, \mu, \gamma; t) + \lambda\text{P}(0,1|\lambda, \mu, \gamma; t) - S\mu\text{P}(0,0|\lambda, \mu, \gamma; t)
\tag{27}
$$

$$
\begin{aligned}
\frac{d\text{P}(S, 0|\lambda, \mu, \gamma; t)}{dt} &= (\lambda + 2\mu)\text{P}(S - 1,0|\lambda, \mu, \gamma; t) \\
&\quad + (\lambda + 2\gamma)\text{P}(S - 1,1|\lambda, \mu, \gamma; t) \\
&\quad - S(\mu + \gamma)\text{P}(S, 0|\lambda, \mu, \gamma; t)
\end{aligned}
\tag{28}
$$

$$
\begin{aligned}
\frac{d\text{P}(0, S|\lambda, \mu, \gamma; t)}{dt} &= (\lambda + \mu)\text{P}(1, S - 1|\lambda, \mu, \gamma; t) \\
&\quad + \lambda\text{P}(0, S - 1|\lambda, \mu, \gamma; t) \\
&\quad - S\gamma\text{P}(0, S|\lambda, \mu, \gamma; t)
\end{aligned}
\tag{29}
$$

160

161 Otherwise:

$$\frac{d\mathrm{P}(k,m|\lambda,\mu,\gamma;t)}{dt}$$

$$= (S-m-(k-1))\left((k-1)\frac{\lambda}{S-1}+2\mu\right)\mathrm{P}(k-1,m|\lambda,\mu,\gamma;t)$$

$$+ (m-1)(S-(m-1)-k)\frac{\lambda}{S-1}\mathrm{P}(k,m-1|\lambda,\mu,\gamma;t)$$

$$+ (k+1)\left((m-1)\frac{\lambda}{S-1}+\mu\right)\mathrm{P}(k+1,m-1|\lambda,\mu,\gamma;t)$$

$$+ (k+1)\left((S-m-(k+1))\frac{\lambda}{S-1}+\gamma\right)\mathrm{P}(k+1,m|\lambda,\mu,\gamma;t) \qquad \textit{30}$$

$$+ (m+1)(S-(m+1)-k)\frac{\lambda}{S-1}\mathrm{P}(k,m+1|\lambda,\mu,\gamma;t)$$

$$+ (m+1)\left((k-1)\frac{\lambda}{S-1}+2\gamma\right)\mathrm{P}(k-1,m+1|\lambda,\mu,\gamma;t)$$

$$- (2\big(k(S-k)+m(S-k-m)\big)\frac{\lambda}{S-1}+\big(2S-(k+2m)\big)\mu$$
$$+ (k+2m)\gamma)\mathrm{P}(k,m|\lambda,\mu,\gamma;t)$$

162  This master equation determines how the methylation statues of a single CpG locus within
163  the stem cell niche evolves over time. The replacement, methylation and demethylation rate
164  are assumed to be constant, hence this process is Markovian and we are able to solve this
165  using standard matrix exponentiation.

166

167  Bayesian analysis of the effect of tissue location and disease state on stem cell
168  dynamics

169

170  The Bayesian pipeline described in the main body of the text allowed the posterior
171  distribution of the parameters defining the stem cell dynamics (namely, the effective number
172  of stem cells, $S$, and the replacement rate per stem cell, $\lambda$) of each individual crypt to be
173  inferred. To interrogate the effect of age, sex, tissue location (colon, small intestine and
174  endometrium) and the disease state of colonic crypts (AFAP/FAP) on stem cell dynamics,
175  we take the posterior mean of $S$ and $\lambda$ as representative of the inferred distribution for each
176  crypt.

177  As a matter of notation, let there be $K$ patients subscripted with $k = [1 .. K]$ and $N$ crypts
178  subscripted with $i = [1 .. N]$. The age of the $k^{th}$ patient is $t_k$, which we normalise to be
179  between 0 and 1 by dividing each patient's age by the maximum age in the patient cohort.
180  Similarly, the sex the $k^{th}$ patient is encoded as a dummy variable, which equals 0 for female
181  patients and 1 for male patients. The location/disease state of each crypt is encoded with the
182  dummy variables $x_{i,j}$ for $j \in \{Colon, Small\ Intestine, FAP, AFAP, Endometrium\}$.

183  We fit the parameters determining stem cell dynamics $y = \{S, \lambda\}$ using a generalised linear
184  model with a gamma-distributed dependent variable (this accounts for the fact $S$ and $\lambda$ are
185  strictly positive). Let $y_{i,k}$ be the dependent variable with expectation $\hat{y}_{i,k}$, then we employ the
186  natural log as a link function $g(\hat{y}_{i,k}) = \ln(\hat{y}_{i,k})$. $y_{i,k}$ is then gamma distributed with mean $\hat{y}_{i,k}$
187  and a tissue/disease-specific standard deviation $\phi_j$.

188  We use the parameterization of the gamma distribution in terms of its shape ($\psi$) and rate
189  ($\omega$):

190
$$\text{Gamma}(y|\psi, \omega) = \frac{\omega^{\psi}}{\Gamma(\psi)} y^{\psi-1} e^{-\omega y}$$

191 The mean of this distribution is $\frac{\psi}{\omega}$ and the variance is $\frac{\psi}{\omega^2}$. Hence, to parameterize the gamma
192 distribution in terms of its mean ($\hat{y}$) and standard deviation ($\phi$), we apply the transformation
193 $\psi = \frac{\hat{y}^2}{\phi^2}, \omega = \frac{\hat{y}}{\phi^2}$.

194 Our dataset contains multiple samples from the same patient, so we assume the offset in the
195 linear predictor is drawn for each patient from a hierarchical normal distribution with mean $\nu$
196 and variation $\sigma$ (hence accounting for random inter-patient variability, not attributable to the
197 factors we are explicitly modelling). Similarly, to maximize the information that can be drawn
198 from the data, we allowed the tissue/disease-specific intrapatient standard deviation, $\phi_j$, to
199 be drawn from a lognormal distribution, with a population mean $\rho$ and standard deviation $\zeta$.

200 Priors:

201
$$a_k \sim \text{normal}(\nu, \sigma)$$

202
$$\ln(\phi_j) \sim \text{normal}(\rho, \zeta)$$

203 Model:

204
$$\ln(\hat{y}_{i,k}) = a_k + b_j x_i^{\ j} + c t_k + d s_k$$

205
$$y_{i,k} \sim \text{gamma}\left(\frac{\hat{y}_{i,k}^{\ 2}}{\phi_j^{\ 2}}, \frac{\hat{y}_{i,k}}{\phi_j^{\ 2}}\right)$$

206 The hierarchical Bayesian model was fit to the data using pystan, a python implementation
207 of Stan[5], a probabilistic programming language that allows for rapid MCMC sampling.

208

209 Because a log-link function was used to ensure the positivity of $\hat{y}_{i,k}$, the coefficients of the
210 regression, $b_j$, encode the difference between each tissue-type or disease-state, and colon
211 on the log scale. We take the exponential transform of each of these regression coefficients
212 to derive the posterior for the relative stem cell number and replacement rate of each tissue-
213 type and disease-state relative to colon. We take a hypothesis testing by parameter
214 inference approach, where the effect of a particular tissue/disease on the dependent
215 variable is termed significant when the 95% equal-tailed credible interval does not overlap 0.
216 The hierarchical Bayesian model that we have developed naturally penalizes increasing
217 numbers of parameters, hence there is no need to apply a multiple test correction [6].

218

219 ## Investigating the well-mixed assumption
220

221 One of the major assumptions taken in the development of the mathematical model
222 describing the FMC distribution was that the stem cells within the niche are well-mixed; that
223 is, each stem cell can replace any other stem cell with equal probability. This assumption
224 was made to minimize the mathematical and computational complexity of the model, as it
225 allowed the state of the stem cell system to be fully characterized with just 2 state variables.
226 However, previous work in mouse suggests that stem cells within the crypt are organized
227 into a ring-like structure, where each stem cell can only replace the 2 cells directly
228 neighboring it (Fig. S3A). It is worth noting that the replacement rate in mouse is on the
229 order of months[7], whereas in human the replacement rate is on the order of years[8]. Hence
230 the potential of stem cells to randomly swap position, as identified by Ritsma et al.[9], raises
231 the possibility that within human crypts, the most accurate model of stem cell replacement is
232 neither perfectly well-mixed nor perfectly organized into a ring.

233

234 To investigate the effect of this well-mixed assumption, a Gillespie simulation was developed
235 to model the stem cell replacement process for the well-mixed and ring geometry (code
236 accessible at https://github.com/CalumGabbutt/flipflop.git[10], see gillespie_crypt.py). Briefly,
237 we generated a 3-dimensional binary array of size $[S, 2, n]$ to model the stem cell niche. The
238 time until the next replacement event was drawn from a Poisson distribution by generating a
239 random uniform number $r \sim \mathrm{uniform}(0, 1)$ as follows:

240
$$\Delta t = \frac{1}{\lambda S} \log\left(\frac{1}{r}\right)$$

241

242 We accounted for the (de)methylation of the individual CpG loci by recognizing that each
243 individual CpG locus on a particular DNA strand was effectively a 2-state system with
244 forward rate $\mu$ and reverse rate $\gamma$. Given that a given CpG locus is methylated at time $t$, the
245 probability that the CpG locus is still methylated at time $t + \Delta t$ is:

246
$$P(on|on) = \frac{\mu}{\mu + \gamma} + \left(1 - \frac{\mu}{\mu + \gamma}\right) e^{-(\mu+\gamma)\Delta t}$$

247 Whilst the probability that the same locus is not methylated is:

248
$$P(off|on) = 1 - P(on|on)$$

249 Similarly, the probability that a given CpG locus that is demethylated at time $t$ is methylated
250 at $t + \Delta t$ is:

251
$$P(on|off) = \frac{\mu}{\mu + \gamma}\left(1 - e^{-(\mu+\gamma)\Delta t}\right)$$

252 And the probability that that CpG locus is still demethylated at time $t + \Delta t$ is:

253
$$P(off|off) = 1 - P(on|off)$$

254

255 Hence, once the time until the next replacement had been drawn, we could update the
256 methylation states of the individual CpG sites by drawing new methylation states with the
257 aforementioned probabilities. The replacement could then be handled by choosing one cell
258 to clonally expand and another to recede. Depending on the geometry of the crypt that the
259 simulation was intended to model, this could be done by selecting one cell at random and
260 then selecting one of that cell's neighbors with equal probability (ring model), or by selecting
261 two cells at randomly without replacement (well-mixed model). To ensure that that our
262 analysis was probing the effect of differing geometry, in each case 100 synthetic crypts were
263 generated and the average methylation probability distribution determined.

264

265 For crypts containing $S = 5$ stem cells which replace each other at a rate $\lambda = 1.0$
266 replacements/stem cell/year. The average FMC distribution of the ring simulations is not
267 significantly different from predictions of the well-mixed model (Fig. S3B, $p > 0.05$
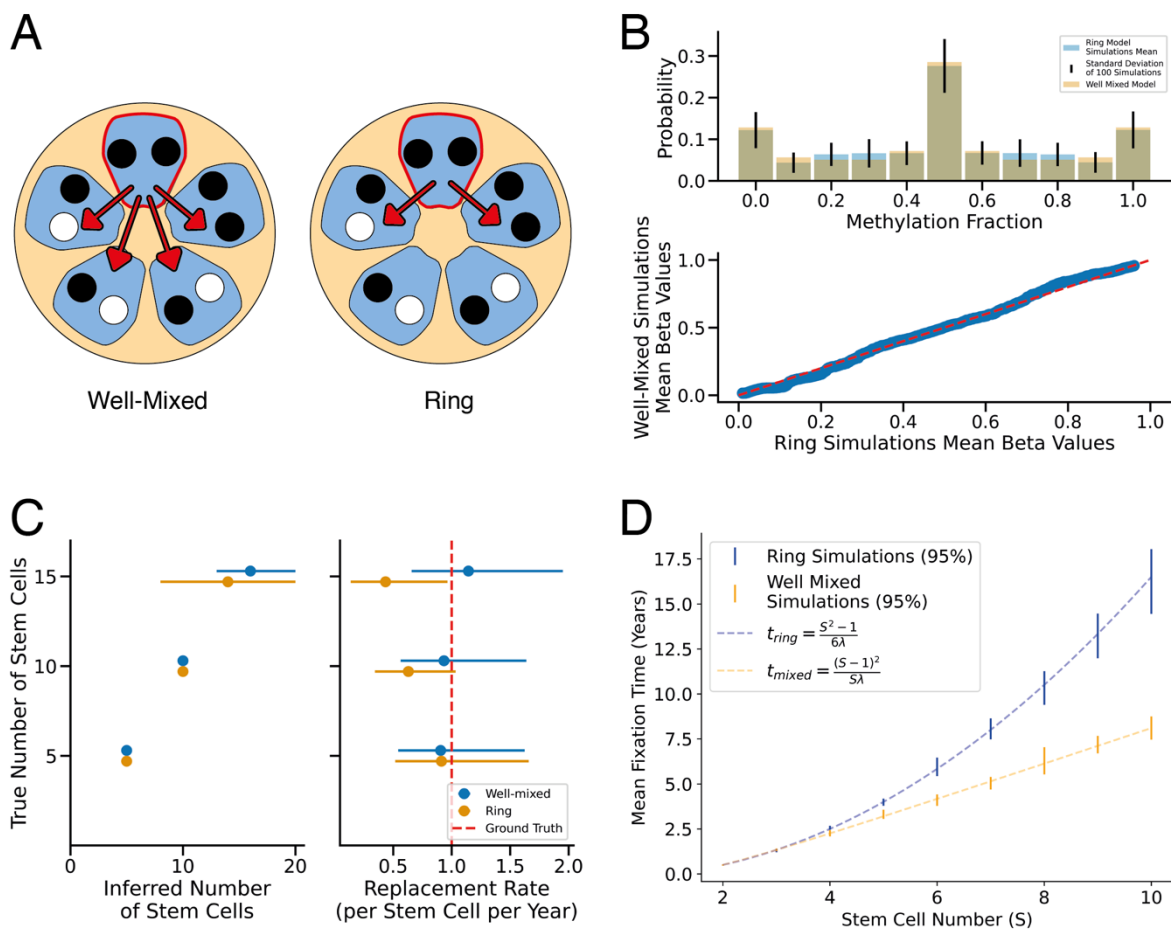268 Kolmogorov–Smirnov test).

269

270 The impact of the geometry with increasing numbers of stem cells on the inference
271 performed upon the methylation distribution was also explored. The above simulations were
272 replicated for $S = 10$ and $S = 15$ for both ring and well-mixed geometries, and the Bayesian
273 inference framework was run upon each of the average methylation distributions. At low
274 values of $S$, the inference model was able to accurately recover the known stem cell
275 dynamics parameters; however, for higher values of $S$ the inferred replacement rate of the

276  ring simulations was lower than the ground truth value (Fig. S3C). The majority of samples
277  analyzed in this study had an inferred stem cell number lower than 10, hence, the impact of
278  the well-mixed assumption was likely to be minimal.

279

280  This result intuitively aligns with our understanding of the mean fixation time (the average
281  time for a mutation that occurs at $t = 0$ to sweep through the population and become clonal,
282  conditioned upon survival), at low stem cell number the two geometries have similar fixation
283  times; however, the fixation time for the ring model scales as $\sim S^2$, whereas the fixation time
284  for the well-mixed model only scales $\sim S$ (Fig. S3D). In the case of a ring geometry with a
285  large number of stem cells, the mean fixation time is larger than for the well-mixed case, and
286  hence there are a greater number of subclonal fCpG loci, which the inference model
287  accounts for by proposing a lower inferred replacement rate.

288



**Supplementary Figure 3: Well-mixed vs. ring stem cell geometry**

**A:** To test the effect of the stem cell geometry on the resulting FMC distribution and the inference process, a simulation of the crypt stem cell dynamics was developed with either a well-mixed (each cell can replace any other) or ring (each cell can only replace its neighbors) geometry. **B:** (top) A histogram displaying the discrete methylation fraction distribution (i.e. before sampling/technical noise has been added) for a set of 100 simulations of $S = 5$ stem cells arranged in a ring geometry, with the well-mixed model predictions overlaid. Error bars denote 1 standard deviation. (bottom) A Q-Q plot comparing the average methylation fraction distribution of the set of 100 ring simulations against that of 100 well-mixed distributions. The two distributions are not statistically different ($p > 0.05$ two-sided KS test). **C:** The 95% credible interval of the posterior for the

number of stem cells (left) and replacement rate (right) compared to the ground truth. **D:** A modified version of the simulation was designed to track the fixation time of a mutation introduced at $t = 0$ to fix within the population. The resulting mean fixation time (10000 total simulations per stem cell number), conditioned on survival, is plotted against the number of stem cells, with the analytic predictions of each model overlaid ($t_{ring} = \frac{S^2 - 1}{6\lambda}$, $t_{mixed} = \frac{(S-1)^2}{S\lambda}$). Bars denote 95% confidence intervals.

289

## Non-fluctuating CpG loci

291

292 In the development of our inference framework, one of the implicit assumptions was that
293 every CpG loci that we had identified as an fCpG was actually behaving in a fluctuating
294 manner. However, whilst every precaution was taken to filter out CpGs which were likely to
295 be under active selection/regulation (see main text), there remained the possibility that a
296 fraction of the identified CpGs do not truly fluctuate. CpG loci that do not behave in a
297 clocklike manner will not track the clonal dynamics of the stem cells, and therefore will dilute
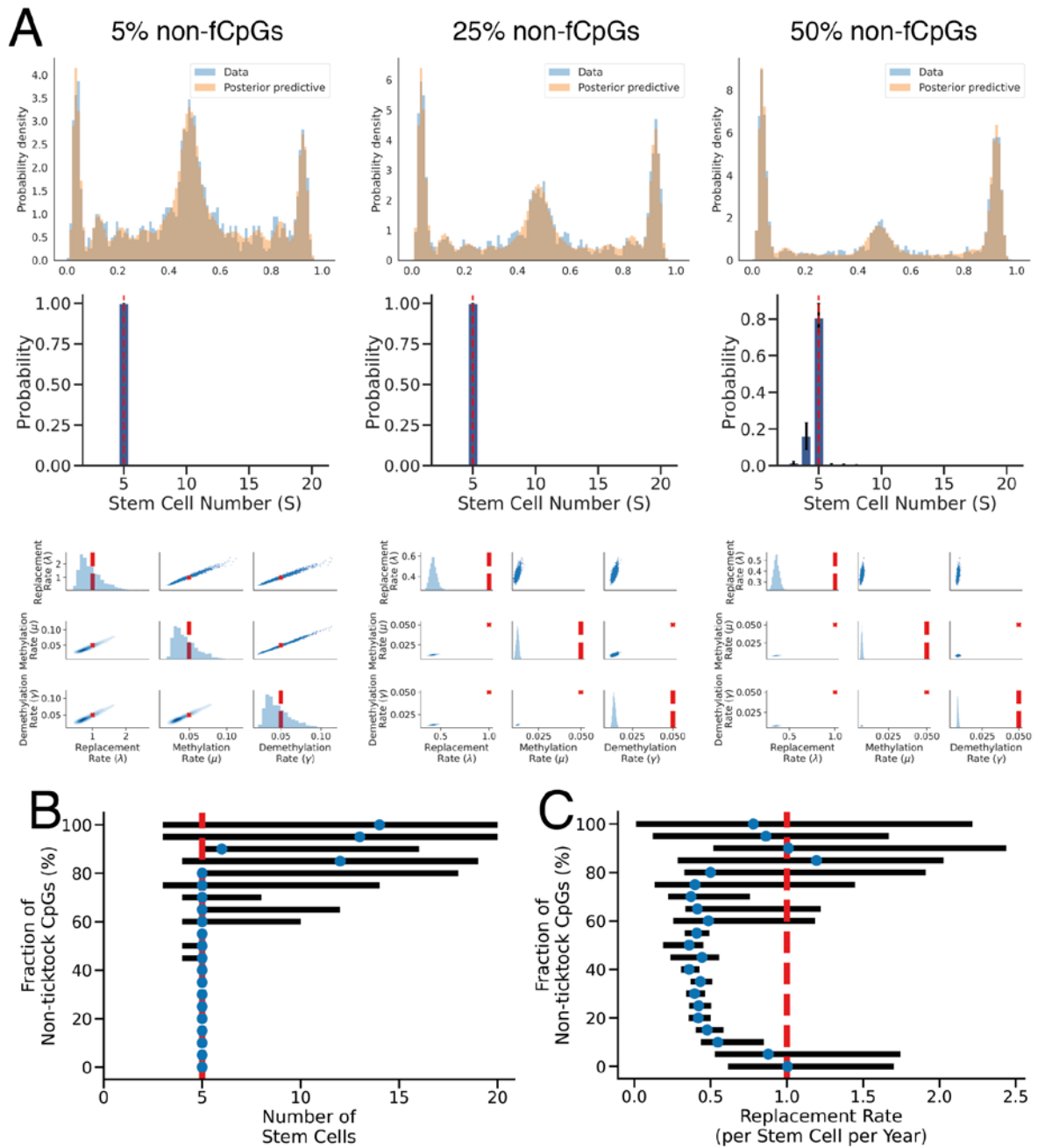298 the timing signal of the fCpGs.

299

300 To investigate the effect that this would have upon our inference framework, we first
301 generated a synthetic crypt containing $S = 5$ stem cells, with a replacement rate $\lambda = 1.0$
302 replacements/ stem cell/year and a (de)methylation rate $\mu = \gamma = 0.05$ per allele/stem
303 cell/year. The noise due to sampling was then simulated so that each of the resulting
304 samples were identical with respect to the fCpG loci using the transforms specified in the
305 main text ($\delta = 0.04, \epsilon = 0.92, \kappa = 200$). Then we replaced a fraction of the fCpGs, $\omega =$
306 $\{0, 0.05, 0.1, ..., 0.9, 0.95, 1\}$ with CpGs that were randomly assigned a beta value of either 0 or
307 1 with equal probability (with noise added such that the non-fluctuating CpG loci had the
308 same noise profile as the fCpG sites).

309

310 The inference framework pipeline was run upon each of the 11 resulting synthetic crypts
311 (example posteriors for $\omega = 005, 0.25, 0.5$ are presented in Fig. S4A). The effect of non-
312 fluctuating CpG loci was to effectively reassign clonal heterozygous and subclonal CpGs to
313 clonal homozygous states. For $\omega \leq 0.5$, there are still sufficient subclonal mutations to
314 accurately infer the number of stem cells (Fig. S4B). Initially, the effect of the relabeling on
315 the inferred parameters was to decrease both the (de)methylation rate and the replacement
316 rate, effectively assuming that the methylation distribution has not relaxed as far from the
317 initial conditions (Fig. S4C). However, once $\omega > 0.5$, the inference struggles to infer the
318 number of stem cells, and the uncertainty over $S$ is propagated into the posterior for $\lambda$ due to
319 the correlations in the posterior between $\lambda$ and $S$.

320

321 This modelling applies equally well to mistakenly identifying fCpG loci that do not fluctuate
322 and are static over time, and the possibility that a fraction of fCpG loci are actively regulated,
323 dynamically setting the methylation status of all the stem cells in the crypt at that locus to be
324 the same. Note that our analysis of fCpG loci located on chromosome X (Fig. S1D) and our
325 finding that fCpG methylation status is preserved along the crypt (Fig. 2C) suggests the
326 influence of cell-type specific methylation is unlikely to be major.

**Supplementary Figure 4: Non-fluctuating CpG loci**

To investigate the effect of a fraction of the CpG sites that we have identified as FMCs in fact being under active regulation, we generated synthetic crypts with an increasing fraction of non-fluctuating CpG loci and ran the inference framework upon them.

**A:** Example posterior predictive distributions and posterior distributions for 5, 25 and 50% of non-fCpGs. (middle) Error bars were calculated from the estimated error (1 standard deviation) on the log-evidence. **B-C:** 95% credible intervals for the inferred replacement rate and stem cell number respectively for an increasing fraction of non-fluctuating CpG sites.

327

328 <span style="color:blue">Mean fixation time</span>

329

330  For small populations of cells in a process of neutral competition, one cell stochastically
331  clonally expanding until it dominates the niche is an inevitability; however, the time-scale
332  over which this process occurs varies depending on both the replacement rate per stem cell
333  and the number of stem cells within the niche. The time it takes for a mutation appearing
334  within a single cell to undergo monoclonal convergence, conditioned upon that mutation not
335  going extinct, follows a positive skewed distribution (Fig. S5A); however, this distribution is
336  often summarized using a single statistic – the mean fixation time.
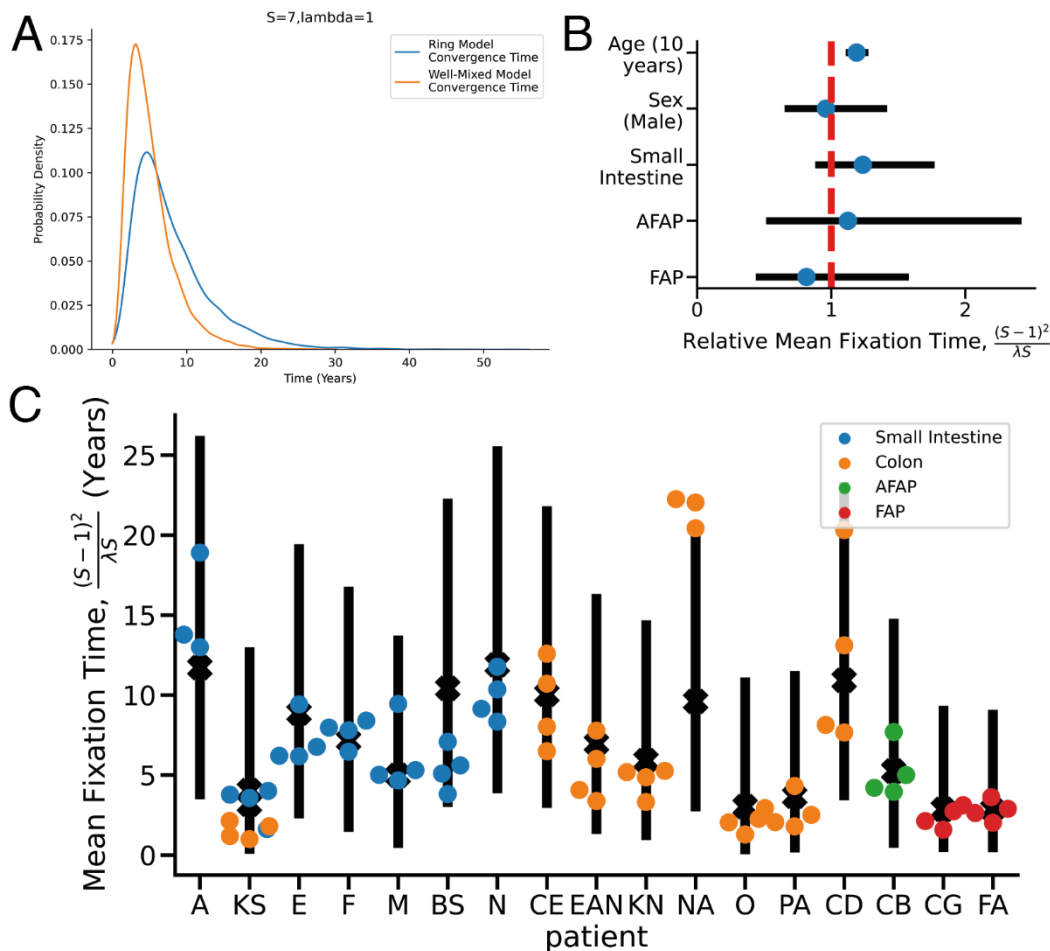
337

338  For a ring geometry, where stem cells can only replace their neighbors, the mean fixation
339  time scales $\sim \frac{S^2}{\lambda}$, whereas for a well-mixed geometry the mean fixation time scales $\sim \frac{S}{\lambda}$ (Fig.
340  S3D). Our inference framework relies upon a well-mixed geometry, and therefore for
341  reasons of self-consistency we shall assume the formula for the mean fixation time is

342  $t_{mixed} = \frac{(S-1)^2}{\lambda S}.$

343

344  For each crypt sample, for every point in the posterior we calculated the mean fixation time
345  to generate a corresponding posterior for the mean fixation time. We followed the same
346  hierarchical Bayesian generalized linear model as in the main text. The only significant factor
347  was the age of the patient (Fig. S5B), suggesting that the rate at which a stem cell colonizes
348  a crypt slows down over the course of a patient's lifetime. We present the mean fixation time
349  for all of the intestinal crypts in Fig. S5C.
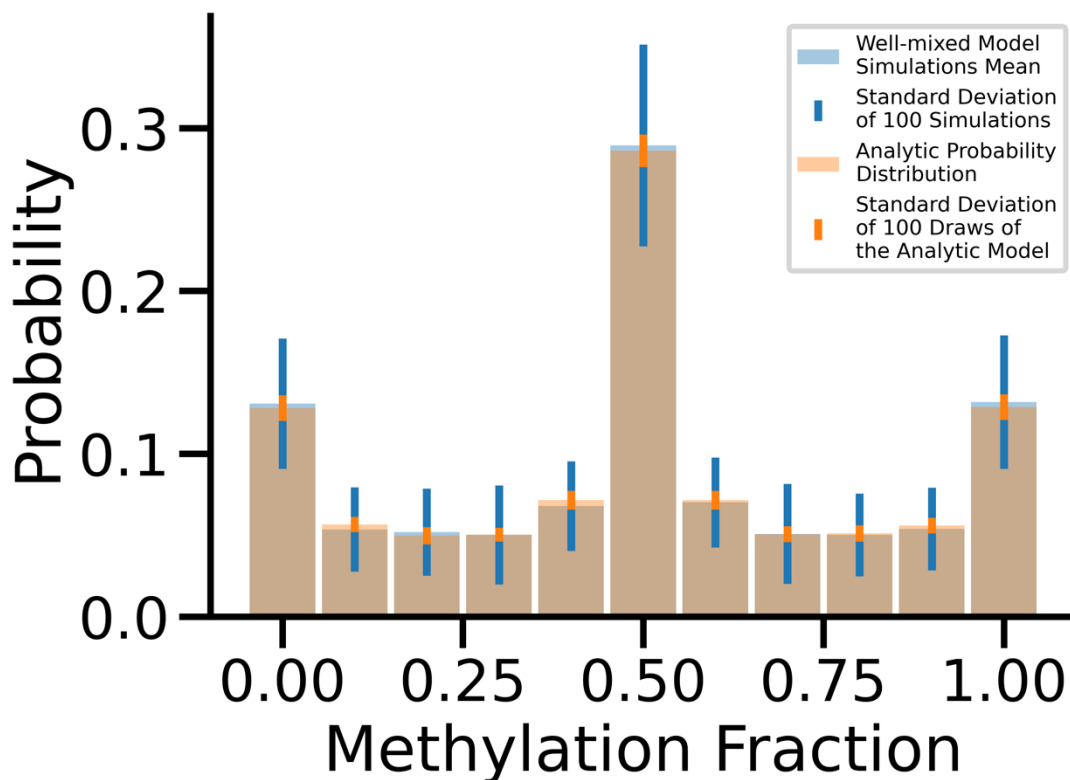
350

**Supplementary Figure 5: Mean fixation time**

**A:** The time for a mutation introduced at $t = 0$ to "fix" within the population, conditioned upon that mutation not going extinct, was simulated using the same simulations outlined in "Supplementary Materials - Investigating the Well-Mixed Assumption". The fixation time follows a skewed distribution, with different shapes depending upon the geometry of the stem cells population. **B:** posterior distributions (95% credible intervals) for the effect of patient age (per decade), sex (with female encoded as reference), tissue type and disease state on the relative mean fixation time compared to normal colon. **C:** individual crypt and posterior mean per patient for the mean fixation time, with the 95% credible range of the generalized linear model (GLM) expectation, accounting for age, sex, tissue, disease state and intra- and inter-patient heterogeneity.

351

## Linkage between CpG Loci

353

354 The mathematical model describing how the distribution of methylation patterns evolves over
355 time presented in the main text treats each fCpG locus independently; however, the
356 replacement process will couple the methylation status of individual CpG loci (because a cell
357 contains a set of genetically "linked" CpGs). Unlike traditional unidirectional lineage tracing
358 markers (e.g. SNVs), the relabeling of individual CpG loci will cause these correlations to
359 naturally erode over time. To investigate the effect that these correlations between individual
360 CpG sites will have upon the FMC distribution, we employed the well-mixed Gillespie
361 simulations described above (Investigating the Well-Mixed Assumption) to generate 100
362 synthetic crypts ($S = 5, \lambda = 1.0, \mu = \gamma = 0.05$) which will naturally include intra- CpG
363 correlations. The resulting probability mass functions (PMF) were compared to 100 draws
364 from the analytic model. The mean PMF of the synthetic crypts exactly matched that of the
365 analytic model (Fig. S6), but the crypt simulations exhibited a wider degree of variability than
366 would be expected from the analytic model alone.

**Supplementary Figure 6: Correlations Between Different fCpG Loci**

The analytic model derived in the text assumes that individual CpG loci behave independently; however, the replacement process will correlate the fates of CpG loci located on separate chromosomes. To investigate this, 100 simulations of well-mixed crypt were performed, and the mean and the standard deviation of the methylation recorded for each possible methylation state. Similarly, the analytic probability distribution was sampled 100 times. Error bars denote 1 standard deviation. The effect of these correlations is that the mean methylation distribution of the simulated crypts exactly matches that of the analytic model, but that the simulated crypts had a higher degree of intra-crypt variability than the analytic model predicts.

367

368  Identifiability of rate parameters

369

370  The posterior of the inference presented in the main text exhibits a degree of collinearity
371  between $\mu$, $\gamma$ and $\lambda$ (Fig. 3A, 4A and 5B). To demonstrate that the system is sensitive to the
372  absolute values of the rate parameters, rather than just their relative values, the methylation
373  distribution was generated for an identical crypt to the "just right" crypt generated in Fig. 3,
374  but with all the rate parameters 10 times smaller (i.e. $\lambda = 0.1$ replacements/stem cell/year,
375  $\mu = \gamma = 0.005$ per CpG locus/year). The resulting methylation was markedly different (Fig.
376  S7A), lacking the clonal heterozygous peak at 0.5 as the system has not yet decayed far
377  from the initial conditions.

378

379

380  The intuition behind this result, and the reason why the absolute rate parameters are
381  obtainable, is because of the slow speed at which the model converges to the steady state.
382  The matrix exponentiation step discussed above can be re-written in terms of the
383  eigenvalues ($u_i$) and eigenvectors ($\vec{v}_i$) of the transition matrix.

384
$$\vec{P}(t) = e^{tT}\vec{P}(t=0) = \sum c_i e^{u_i t} \vec{v}_i$$

385

386  The system decays from the initial state according to the magnitude of the non-zero
387  eigenvalues (which are all negative) towards the steady state (the eigenvector
388  corresponding to the eigenvalue with 0 magnitude). Therefore, the smallest magnitude non-
389  zero eigenvalue determines how rapidly the system decays to the steady state.

390

391  The eigenvalues can be calculated numerically for a given set of parameters $\{\lambda, \mu, \gamma, S\}$. If we
392  use the same set of parameters as in the "just right" simulated data (Fig. 3, $\{\lambda = 1.0, \mu =$
393  $0.05, \gamma = 0.05, S = 5\}$), then the smallest magnitude non-zero eigenvalue has a value of $-0.1$
394  /year. After 30 years, the slowest eigenvalue will have decayed to $e^{-3} \approx 0.05$ of its initial
395  value – certainly low, but sufficient for the inference model to accurately infer the
396  replacement rates in real units, as exhibited in Fig. 3. When we scale each of the rate
397  parameters to be 10 times smaller, as presented in Fig. S7A, then although the ratio of the
398  replacement:methylation:demethylation rates are unchanged, the system has decayed much
399  less towards the steady state, which the inference framework is able to detect and quantify.
400  This reinforces the importance of selecting fCpG loci with a (de)methylation rate that is "just
401  right". If sites (de)methylate too quickly, the system is indistinguishable from the steady state
402  distribution and the information on the absolute values of the rate parameters is lost (as in
403  the "too fast" synthetic crypt in Fig. 3). If sites (de)methylate too slowly, the methylation state
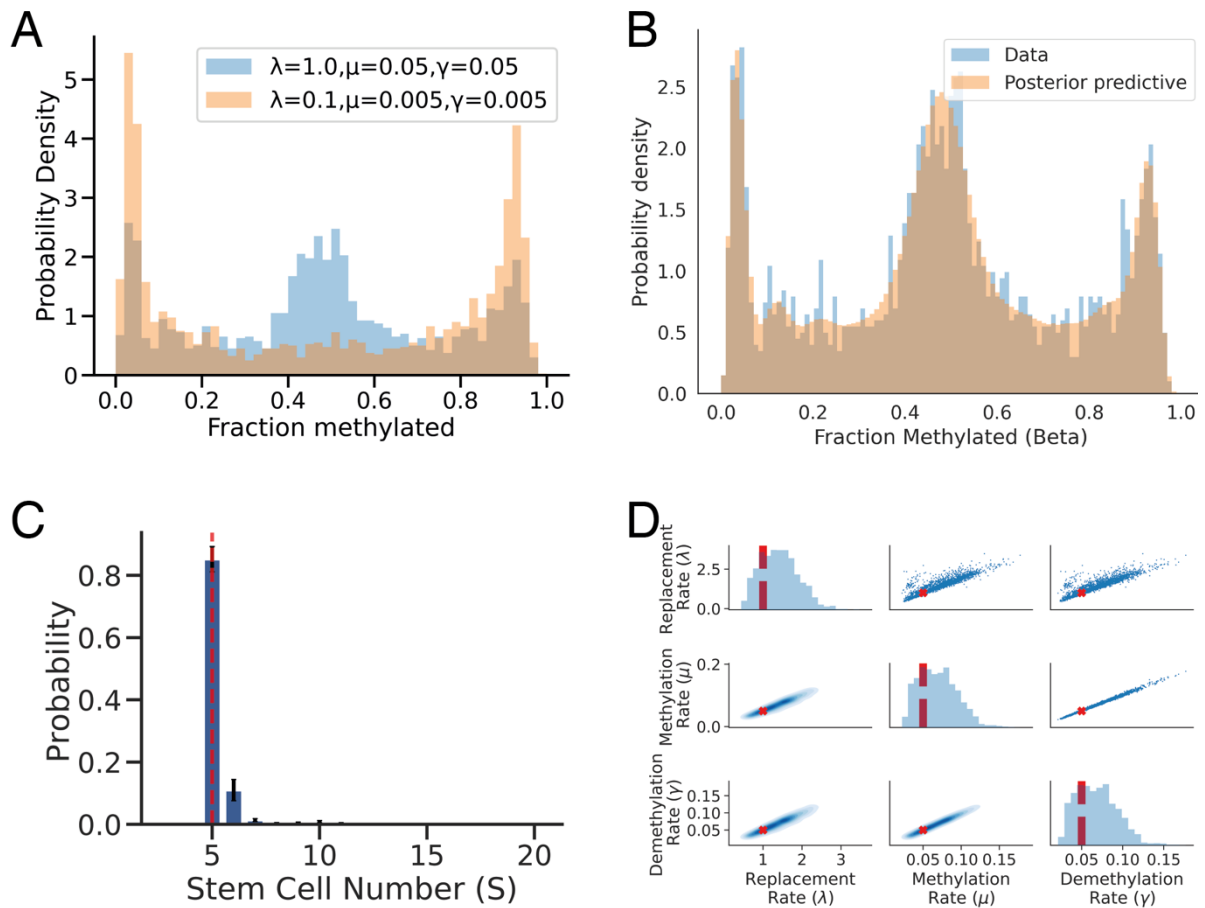404  of the system does not evolve sufficiently over the observed time period.

405

406  We note that in the real patient data, the mean inferred methylation rate is 0.027 across the
407  sample cohort, approximately half that of the assumed value in the simulated data, and so
408  the decay to steady state in patient is slower than presented in the simulations in Fig 3.

409

410  To ensure that the parameters were separately identifiable, the inference on the "just right"
411  simulated data (Fig. 3) was repeated with a wider prior on the (de)methylation rate
412  ($\mu, \gamma \sim \text{halfnormal}(0.5)$). The 95% credible interval still contained the ground truth parameters
413  values and the posterior predictive distribution well-matched the data (Fig. S7B-D), hence
414  the parameters were separately identifiable.

415

**Supplementary Figure 7: Identifiability of parameters**

**A:** Example methylation distributions for identical crypts, except with the rate parameters of one crypt 10 times smaller than the other. The resulting methylation distributions are strikingly different. **B-C** To demonstrate that the parameters are separately identifiable, the inference upon the "just right" synthetic crypts (Fig. 3) was repeated with 10 times wider priors upon the (de)methylation rates. The resulting posterior is still able to accurately recover the ground truth values. Error bars were calculated from the estimated error (1 standard deviation) on the log-evidence.

416

## Whole blood simulations

417
418
419 Whole blood was simulated in Java using the HAL framework[11] as a non-spatial agent-
420 based model using 27,634 fCpG sites as measured in the experimental data. Parameters
421 (Supplemental Table 1) for normal hematopoiesis are numbers of hematopoietic stem cells
422 (N, HSCs), number of possible division events (T), CpG error rates (S, methylation and
423 demethylation) for the fCpG sites, and HSC replacement dynamics ($\lambda$). To model clonal
424 expansion, a single cell was selected to grow upon induction, and added parameters are its
425 expansion rate (E) and its final blood frequency of the clonal expansion ($\omega$). These clonal
426 expansions resulted in the overall population size to grow until the appropriate final blood
427 frequency was reached. The output of the simulations provided the beta values at the fCpG
428 sites and the overall distribution variance over time.

429 The number of HSCs was set at a lower value of 1000 initiating cells. This was much lower
430 than the 30,000 based on the large number of HSC inferred by DNA sequencing studies[12,13];
431 however, the results shown here are invariant to more than 100 initiating cells. CpG error

432  rates varied between CpG sites and were assigned based on the distribution averages of the
433  656 normal individuals from GSE40279[14]. We found that some of the whole blood fCpGs did
434  not appear to have equal methylation and demethylation error rate because their averages
435  tended to always be above or below 50% in multiple individuals. Hence, to better model and
436  match the data, we used a look-up distribution table in the simulations in order to initialize a
437  cell's fCpG parameters, with lower and unequal error rates at CpG sites with average
438  methylation typically found near 0.4 (demethylation > methylation) or near 0.6 (methylation >
439  demethylation) to maintain the variance of the 27,634 fCpG sites around 0.1 during cell
440  divisions. The error rates varied between 0.0001 to 0.001 changes per division, with the
441  highest error rates and more equal methylation and demethylation rates at CpG sites near
442  50% methylation.

443  Cell survival was set at exact replacement (one cell produces one living offspring), and
444  results did not vary much if random replacement was simulated. A proportion of cells
445  underwent replacement at each timestep (Supplemental Table 1). For the neoplastic
446  simulations in Fig 6D in the manuscript, the expansion rate (E) was varied to model either
447  rapid expansion (visible or more than 5% leukemic cells within 1 year or 200 divisions) akin
448  to acute leukemia, modest expansion (visible within 4 years or 12,000 divisions), or very
449  slow expansion (visible within 6 years or 18,000 divisions). The extent of blood involvement
450  was varied between 20% (black lines), 50% (blue lines) and 90% (red lines). These
451  simulations indicated that how clonal expansions change whole blood  fCpG variances
452  depends both on how fast the expansion grows and to what extent it involves the blood.
453  Rapid growth to high levels like acute leukemias results in high fCpG variances and
454  characteristic W-shaped fCpG distributions. Slower growth to lower levels like chronic
455  leukemias results in low fCpG variances and broader distributions that lack the W-shape.
456  Interestingly, very indolent clonal expansions which may occur with CHIP[15] can result in
457  small increases in fCpG variances, which may account for the age-related increase in fCpG
458  variances seen in Fig 6A in the manuscript.

459  More sophisticated modelling with a better selection of whole blood fCpG sites could
460  improve the extraction of ancestral information. For example, a selection of slower fCpG
461  sites may improve the detection and analysis of indolent clonal expansions, where many of
462  the faster fluctuations return to average ~50% methylation by the time the expansion
463  reaches detectable blood levels.

464  The simulation framework can be obtained, along with sample simulation results, on GitHub
465  through https://github.com/MathOnco/flipflopblood.git[16]. A GUI compatible with most
466  operating systems is accompanied to allow for rapid evaluation of different parameters.

467

| Parameter | Description | Values |
|---|---|---|
| N | Number of HSCs, initial population | 100 |
| T | Simulation time | 2,500 |
| S | CpG error rates | 0.0001-0.001 per division |
| λ | Cell survival, exact replacement | 0.6 |
| E | Disease expansion rate | (0.1, 0.005, 0.00225) |
| ω | Final blood frequency | (0.2, 0.5, 0.9) |

**Supplementary Table 1: Parameters of whole blood simulations**

Parameters used in simulations describing how the methylation distribution of well-mixed
hematopoietic stem cells (HSCs) changes in response to the expansion of a single clonal
population.

468

## Supplementary References

469
470

471 1.  Métivier, R. *et al.* Cyclical DNA methylation of a transcriptionally active promoter.
472     *Nature* **452**, 45–50 (2008).

473 2.  Kangaspeska, S. *et al.* Transient cyclical methylation of promoter DNA. *Nature* **452**,
474     112–115 (2008).

475 3.  Rulands, S. *et al.* Genome-Scale Oscillations in DNA Methylation during Exit from
476     Pluripotency. *Cell Syst.* **7**, 63-76.e12 (2018).

477 4.  Parry, A., Rulands, S. & Reik, W. Active turnover of DNA methylation during cell fate
478     decisions. *Nature Reviews Genetics* vol. 22 59–66 (2021).

479 5.  Carpenter, B. *et al.* Stan: A Probabilistic Programming Language. *J. Stat. Software;*
480     *Vol 1, Issue 1* (2017).

481 6.  Gelman, A., Hill, J. & Yajima, M. Why We (Usually) Don't Have to Worry About
482     Multiple Comparisons. *J. Res. Educ. Eff.* **5**, 189–211 (2012).

483 7.  Kozar, S. *et al.* Continuous Clonal Labeling Reveals Small Numbers of Functional
484     Stem Cells in Intestinal Crypts and Adenomas. *Cell Stem Cell* **13**, 626–633 (2013).

485 8.  Nicholson, A. M. *et al.* Fixation and Spread of Somatic Mutations in Adult Human
486     Colonic Epithelium. *Cell Stem Cell* **22**, 909-918.e8 (2018).

487 9.  Ritsma, L. *et al.* Intestinal crypt homeostasis revealed at single stem cell level by in
488     vivo live-imaging. *Nature* **507**, 362–365 (2014).

489 10. Gabbutt, C. *et al.* Cell lineage tracing with molecular clocks based on fluctuating DNA
490     methylation - flipflop. *Zenodo* (2021) doi:10.5281/zenodo.5347259.

491 11. Bravo, R. R. *et al.* Hybrid Automata Library: A flexible platform for hybrid modeling
492     with real-time visualization. *PLoS Comput. Biol.* **16**, e1007635 (2020).

493 12. Lee-Six, H. *et al.* Population dynamics of normal human blood inferred from somatic
494     mutations. *Nature* **561**, 473–478 (2018).

495 13. Watson, C. J. *et al.* The evolutionary dynamics and fitness landscape of clonal
496     hematopoiesis. *Science (80-. ).* **367**, 1449–1454 (2020).

497 14. Hannum, G. *et al.* Genome-wide Methylation Profiles Reveal Quantitative Views of
498     Human Aging Rates. *Mol. Cell* **49**, 359–367 (2013).

499 15. Jaiswal, S. & Ebert, B. L. Clonal hematopoiesis in human aging and disease. *Science*
500     *(80-. ).* **366**, eaan4673 (2019).

501 16. Schenck, R. *et al.* Cell lineage tracing with molecular clocks based on fluctuating DNA
502     methylation – Flip flop blood model. *Zenodo* (2021) doi:10.5281/zenodo.5348301.

503