# Neutral competition within a long-lived population of symmetrically dividing cells shapes the clonal composition of cerebral organoids

Florian G. Pflug[1*], Simon Haendeler[1], Christopher Esk[2], Dominik Lindenhofer[2], Jürgen A. Knoblich[2], Arndt von Haeseler[1,3]

[1] Center for Integrative Bioinformatics Vienna (CIBIV), Max Perutz Laboratories, University of Vienna and Medical University of Vienna, Vienna Bio Center (VBC), Vienna, Austria
[2] Institute of Molecular Biotechnology of the Austrian Academy of Science (IMBA), Vienna Bio Center (VBC), Vienna, Austria
[3] Bioinformatics and Computational Biology, Faculty of Computer Science, University of Vienna, Vienna, Austria
*Correspondence: florian.pflug@univie.ac.at

## Summary

Cerebral organoids model the development of the human brain and have become an indispensable tool for studying neural development and neuro-developmental diseases. Comprehensive whole-organoid lineage tracing has revealed the number of progeny arising from each initial stem cell to be highly diverse, with lineage sizes ranging from one to more than 20,000 cells. This variability exceeds what can be explained by existing stochastic models of corticogenesis, which indicates that an additional source of stochasticity must exist. We propose the quantitative SAN model in which this additional source of stochasticity is neutral competition within a long-lived population of symmetrically dividing cells. In this model, the eventual size of a lineage is determined by its survival time within this population of symmetrically dividing cells, which due to neutral competition varies widely between individual lineages. We demonstrate the SAN model to explain the experimentally observed variability of lineage sizes and use it to derive a formula that captures the quantitative relationship between survival time and lineage size. Finally, we show that our model implies the existence of a mechanism which keeps the size of the population of symmetrically diving cells approximately constants, and that it enables this mechanism to be probed experimentally.
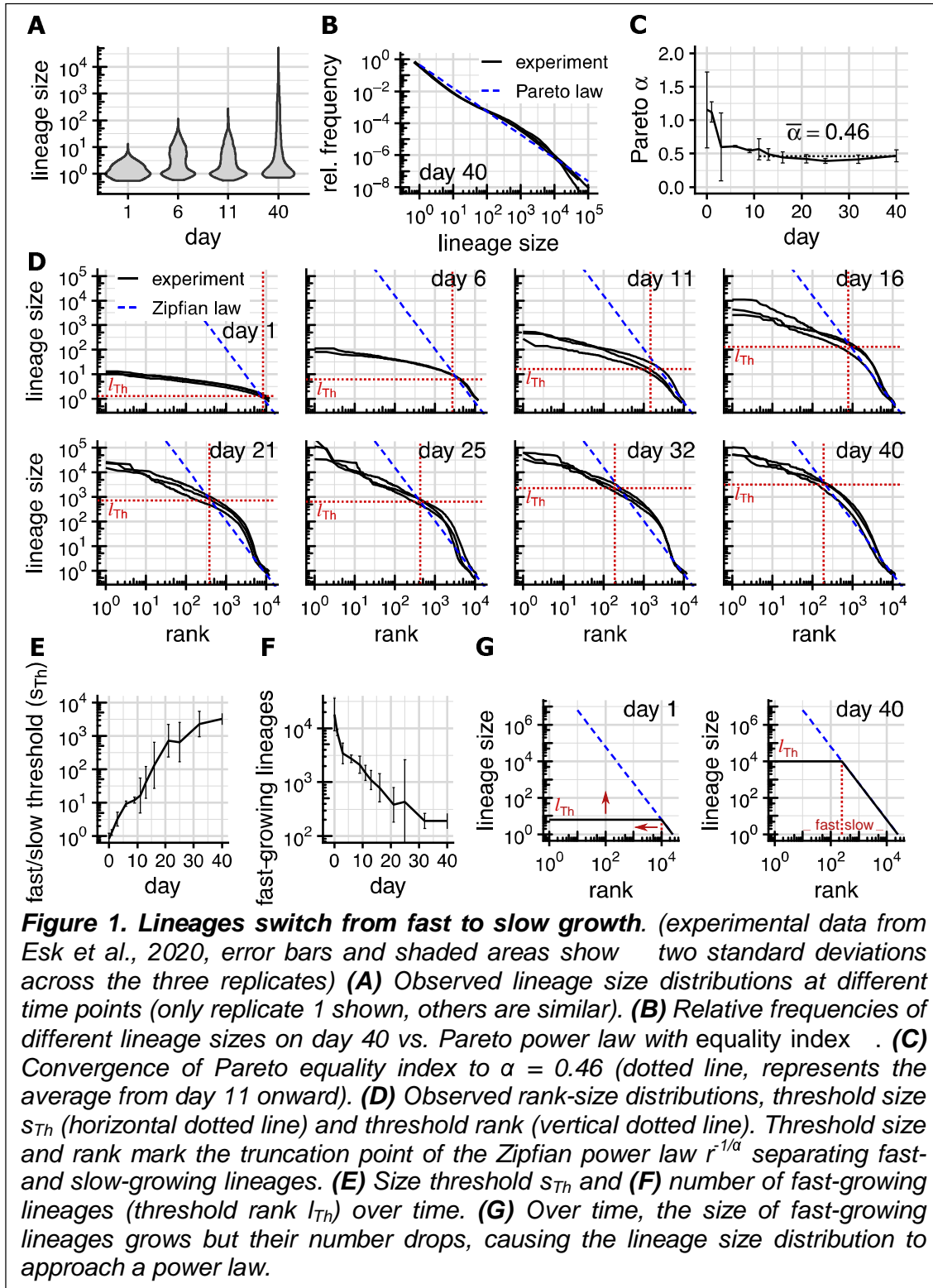
## Introduction

The development and maintenance of the tissues and organs comprising complex organisms rely on sophisticated genetic programs to coordinate the differentiation of cells in both space and time. In many cases, this "program" does not consist of fully deterministic decision chains but instead contains stochastic components; examples include the development of the cortex (Klingler and Jabaudon, 2020; Llorca et al., 2019) and stem cell homeostasis in intestinal crypts (Snippert et al., 2010).

During cortical development, neurons are produced (directly or indirectly) by progenitor cells in the ventricular zone called radial glial cells (RGCs). In mice, the neuronal output of individual RGCs was observed to vary by about one to two orders of magnitude between seemingly identical progenitors, which suggests a stochastic model of cortical neurogenesis (Llorca et al., 2019). Cerebral organoids grown from human stem cells (Lancaster et al., 2017) show even stronger variability of offspring numbers; comprehensive whole-organoid lineage tracing data shows the sizes of individual lineages arising from each ancestral stem cells to vary over up to four to five orders of magnitude (Esk et al., 2020). For RGCs, an alternative explanation of the apparent stochasticity of their number of offspring are hidden variables (Zechner et al., 2020) like transcriptional state within the seemingly homogenous population of progenitors. But in organoids, we expect the pool of ancestral stem cells to be homogenous, and thus conclude that lineage sizes vary predominantly due to stochastic effects. While RGCs output may vary more widely in humans than in mice, varying RGC output alone still cannot account for lineage sizes varying over 4 to 5 orders of magnitude in human cerebral organoids. There is thus likely an additional source of stochasticity in organoid development beyond the stochastic model of neurogenesis proposed by Llorca et al (2019).

In this study, we propose the source of this additional stochasticity to be neutral competition within a long-lived population of roughly 10,000 symmetrically dividing stem cells (S-cells). Neutral competition between stem cells has previously been shown to shape the clonal composition of tissues in homeostasis (Snippert et al., 2010; Corominas-Murtra et al., 2020), and to accurately predict the time until monoclonality (the time until all but a single lineage has died out). We show that in growing tissue like cerebral organoids, neutral competition does not lead to eventual monoclonality. Instead, the tissue records the changing clonal composition of its

65  stem cell population, which causes the sizes of individual lineages to grow
66  increasingly diverse over time. To quantify this effect and its dependence on the size
67  of the S-cell population, we introduce the stochastic SAN model and show that



**Figure 1. Lineages switch from fast to slow growth**. *(experimental data from Esk et al., 2020, error bars and shaded areas show two standard deviations across the three replicates) **(A)** Observed lineage size distributions at different time points (only replicate 1 shown, others are similar). **(B)** Relative frequencies of different lineage sizes on day 40 vs. Pareto power law with equality index . **(C)** Convergence of Pareto equality index to α = 0.46 (dotted line, represents the average from day 11 onward). **(D)** Observed rank-size distributions, threshold size $s_{Th}$ (horizontal dotted line) and threshold rank (vertical dotted line). Threshold size and rank mark the truncation point of the Zipfian power law $r^{-1/\alpha}$ separating fast- and slow-growing lineages. **(E)** Size threshold $s_{Th}$ and **(F)** number of fast-growing lineages (threshold rank $l_{Th}$) over time. **(G)** Over time, the size of fast-growing lineages grows but their number drops, causing the lineage size distribution to approach a power law.*

68 neutral competition within its S-cell population suffices to explain the observed
69 variation of lineage sizes over four to five orders of magnitude.

## Results

**Empirical lineage size distribution**

72 In the experiment conducted by Esk *et al.* (Esk et al., 2020), cerebral organoids were
73 grown from roughly 24,000 stem cells, genetically identical except for a distinct
74 genetic barcode in each cell serving as a *lineage identifier* (LID). To determine the
75 contribution of each initial stem cell to organoids of different ages, organoids were
76 subjected to amplicon high-throughput sequencing. The sequencing reads (after
77 filtering and error-correction) corresponding to each LID were counted, and the per-
78 LID read counts normalized to an approximate number of cells comprising each
79 lineage (see *experimental procedures* for details).

80 The resulting *lineage size distribution* (figure 1A) shows, as expected, small
81 and equally sized lineages for organoids harvested at day 1 (lineages sizes around 1
82 cell). The distribution grows more uneven until day 11 (up to 30 cells/lineage) and
83 extends over 4 to 5 orders of magnitude (up to 100,000 cells/lineage) after 40 days.

84 A common mathematical model for distributions extending over multiple orders
85 of magnitude are so-called (Pareto) power laws where the frequency of objects of
86 size $l$ or larger is proportional to $l^{-\alpha}$. Parameter $\alpha$ is called the (Pareto) equality
87 index because it determines how even (large $\alpha$) or uneven (small $\alpha$) object sizes are
88 distributed. In double-logarithmic frequency vs. size plots, power laws appear as
89 straight lines with slope $\alpha$, which we find matches the lineage size distribution on day
90 40 well for $\alpha \approx 0.46$ (figure 1B). We remark that $\alpha \approx 0.46$ represents a small equality
91 index (i.e. diverse lineage sizes); in applications of Pareto distributions values of $\alpha$
92 often lie between 1 and 2.

93 While the unevenness of the lineage size distribution grows considerably
94 between days 11 and 40 (figure 1A), the equality index stays close to $\alpha \approx 0.46$ from
95 day 11 onwards (figure 1C). The equality index thus fails to capture the large
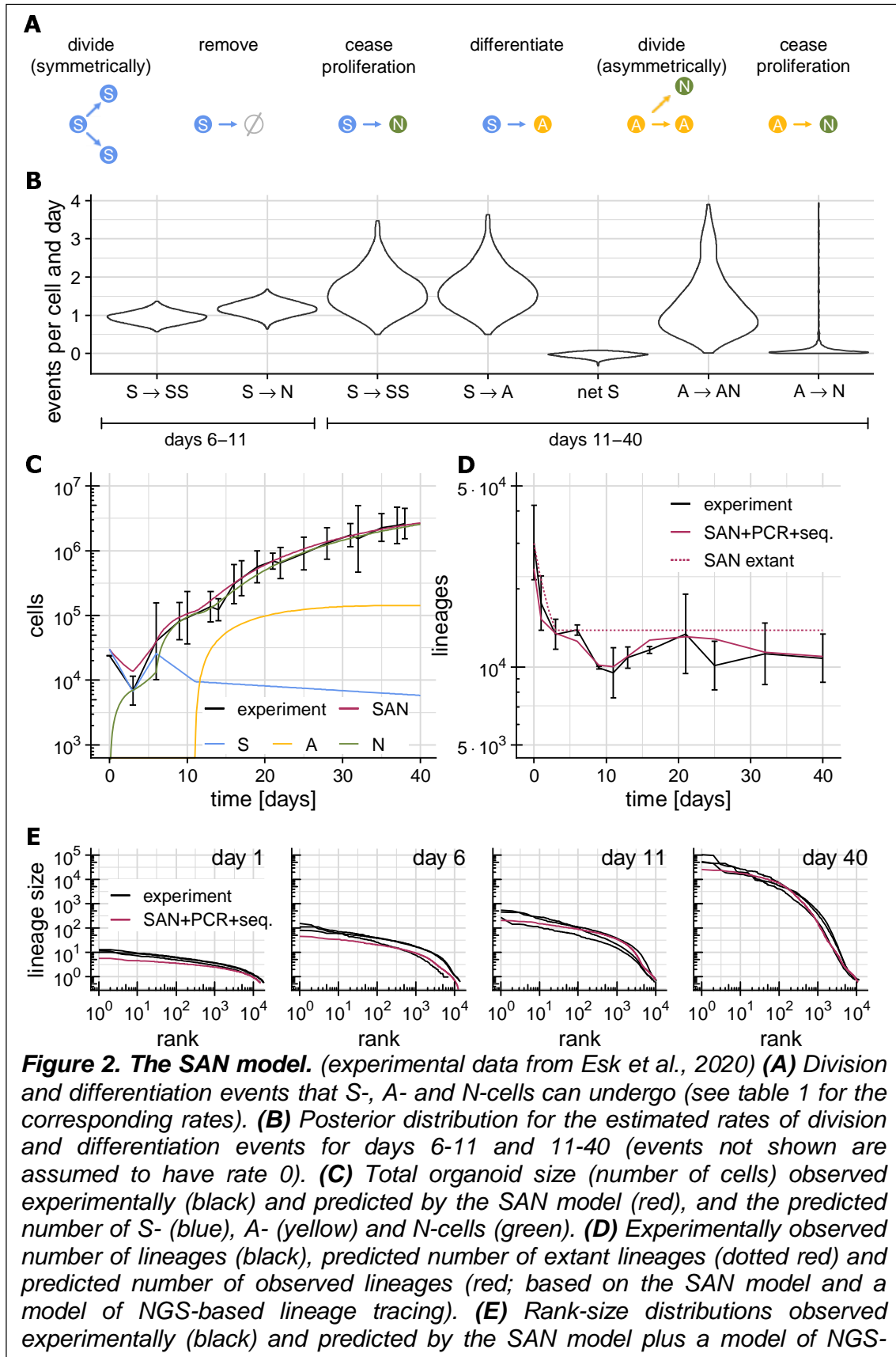96 increase in non-uniformity of lineage sizes between days 11 and 40.

**Truncated Zipfian rank-size distribution**

To describe the evolution of lineage sizes over time, we thus instead rank lineages by size (largest lineage first) and plot the resulting *rank-size distributions* (figure 1D). For lineage sizes governed by a Pareto law with index $\alpha$, the rank-size distribution would be expected to be governed by a Zipfian power law, meaning lineage sizes should decrease proportional to $r^{-1/\alpha}$ with increasing rank $r$ (Adamic and Huberman, 2002).

Instead, we observe *truncated* Zipfian laws in which lineages sizes obey a Zipfian law ($\alpha \approx 0.46$) only up to a certain threshold size $l_{\mathrm{Th}}$ above which lineages are multiple orders of magnitudes smaller and more uniform than the Zipfian law would predict (figure 1D). The threshold size $l_{\mathrm{Th}}$ grows more than 1,000-fold (figure 1E) over 40 days, while the ratio between threshold size and largest lineage size grows only by a factor of 2 (from 8.5 to 21; figure 1D); lineages above the threshold therefore grow roughly uniformly. Lineages below the threshold, in contrast, show no overall shift towards larger lineages sizes over time, indicating that growth has mostly ceased for these lineages.

**Lineages switch from fast to slow growth**

The size threshold $l_{\mathrm{Th}}$ thus partitions lineages according to their growth regime into *fast-growing* and *slow/non-growing*. Of the (on average) 10,851 lineages that contribute to the final organoid 8,389 lineages fall into the fast-growing category on day 1; but on day 11 their number has dropped to 1.496, and on day 40 only 191 (about 2%) fast-growing lineages remain (figure 1F). Lineages thus start out fast-growing, and one by one switch to a regime of slow/no growth as time progresses. The later that switch occurs for a particular lineage, the bigger it has become before its growth ceases, leading to larger and larger lineages in the slow-growing regime and consequently to $l_{\mathrm{Th}}$ increasing as time progresses. In this coarse approximation, lineages are assumed to have the same size as long as they are fast-growing (figure 1G); experimentally we observe a spread of 1.5 orders of magnitude within the sizes of fast-growing lineages versus a spread of 3.5 orders of magnitude within the slow-growing regime (figure 1D).

**Figure 2. The SAN model.** *(experimental data from Esk et al., 2020)* ***(A)*** *Division and differentiation events that S-, A- and N-cells can undergo (see table 1 for the corresponding rates).* ***(B)*** *Posterior distribution for the estimated rates of division and differentiation events for days 6-11 and 11-40 (events not shown are assumed to have rate 0).* ***(C)*** *Total organoid size (number of cells) observed experimentally (black) and predicted by the SAN model (red), and the predicted number of S- (blue), A- (yellow) and N-cells (green).* ***(D)*** *Experimentally observed number of lineages (black), predicted number of extant lineages (dotted red) and predicted number of observed lineages (red; based on the SAN model and a model of NGS-based lineage tracing).* ***(E)*** *Rank-size distributions observed experimentally (black) and predicted by the SAN model plus a model of NGS-*

128      This proposed lineage-specific switching from fast to slow growth can also
129 quantitatively reproduce the observed truncated Zipfian with $\alpha = 0.46$, with one
130 mathematically simple example being fast-growing lineages growing exponentially
131 with rate $\gamma$ and the number of fast-growing lineages declining exponentially with rate
132 $\sigma = \alpha\gamma$. But while this simple example assumes an unspecified biological mechanism
133 behind the lineage-specific growth regime switches, we show in the following that no
134 such mechanism is in fact necessary. Instead, we show that such growth regime
135 switches emerge naturally from a cellular model of organoid growth.

**SAN model**

137 In the SAN model of organoid growth we distinguish between three types of cells
138 based on the proliferation behavior they exhibit (figure 2A). Cells are either
139 *symmetrically* dividing (S-cells), *asymmetrically* dividing (A-cells) or *non-dividing* (N-
140 cells). In this model, S-cells have the ability to self-renew indefinitely through
141 symmetric division and can thus be considered stem cells. They form the initial cell
142 population of an organoid, and apart from dividing symmetrically they differentiate
143 into either A- or N-cells or are removed permanently. A-cells are cells that have
144 committed to a differentiation trajectory and produce N-cells through asymmetric
145 division, while N-cells do not further divide. We emphasize that S, A and N refer
146 solely to a cell's proliferation behavior, not its functional cell type.

147      All these division and differentiation events occur randomly and independently
148 for each cell with specific time-dependent rates (table 1); from a single-lineage
149 perspective, the SAN model is thus stochastic in nature. Any difference between the
150 trajectories of the lineages arising from different ancestral cells is thus assumed to
151 be purely the result of random chance, not of cell fate decisions or spatial
152 configuration. From a whole-organoid perspective, on the other hand, the SAN
153 model is deterministic, because random effects average out over the roughly 10,000
154 lineages comprising an organoid.

**Division and differentiation rates**

156 To find the rates of cell division and differentiation, we split the organoid
157 development into four time intervals (days 0-3, 3-6, 6-11 and 11-40) according to the
158 main phases of the protocol of Lancaster *et al.* (Lancaster et al., 2017). Until day 6,
159 formation of embryoid bodies (EBs) is still ongoing, and organoid development thus

160    does not reflect development *in vivo*. For these time intervals we manually chose

161    rates of S-cell division (S → S S), removal (S → ∅) and death (S → N; dead cells

162    present in the EB are still counted by NGS-based lineage tracing) for which predicted

163    and observed numbers of cells, lineages, and lineage sizes match (table 1). After

164    day 6, EB formation is complete, and no further cells are removed from the organoid.

165    Until day 11 S-cells are then assumed to either divide symmetrically (S → S S) or

166    cease proliferation (S → N), but to not produce A-cells yet. After embedding the

167    organoids into Matrigel droplets on day 11, organoid growth enters the asymmetric

168    division phase where S-cells are assumed to multiply (S → S S) and to differentiate

169    into A-cells (S → A), which then produce N-cells through asymmetric division (A → A

170    N) before they eventually cease to proliferate (A → N).

171    From day 6 onwards organoid development reflects development in vivo and

172    we hence desired to identify the range of likely rates for each event in addition to a

173    single most-likely value. We thus adopted a Bayesian model comprising log-normally

174    distributed measurement inaccuracies on top of the SAN model, and used Markov

175    chain Monte Carlo (MCMC) sampling to find 1,000 likely rate combinations and their

176    (posterior) probabilities (figure 2B). To arrive at a single set of most-likely values for

177    the rates to be estimated, we then computed MAP (maximum a-posteriori) estimates

178    (table 1) from this posterior distribution.

179    Both the posterior distribution (figure 2B) and the MAP estimates (table 1) show

180    the *net* rate of S-cell proliferation (the rate with which the S-cell population grows or

181    shrinks, i.e. the difference between the rates of S → S S and S → A) to lie close to

182    zero. From this, we conclude that the size of the S-cell population changes only

183    slowly from day 11 onwards. The posterior distributions of the individual rates are, on

184    the other hand, much broader. While the MAP estimates are thus arguably the single

185    most likely set of rates, other combinations of rates are possible as well.

186    **Model validation**

187   For the MAP rate estimates (table 1), the predicted organoid sizes between day 0
188   and 40 agree well with the experimentally determined number of cells (figure 2C).
189   Similarly, the lineage size distribution predicted by the SAN model matches the
190   observed lineage size distribution both in the original data of Esk et al. (figure 2E) as
191   well as in independent replicate experiments (figure S1). In particular, the predictions
192   show the same truncated Zipfian distributions as the experimental data recapitulate
193   the spread over 4 – 5 orders of magnitude. The SAN model predicts the number of
194   *extant* lineages (lineages containing at least one S-, A or N-cell) to drop to about
195   ≈13,700 on day 3 where it then remains. This drop in the number of extant lineages
196   is caused by lineages that do not make it into the organoid during EB formation.
197   While the predicted number of remaining lineages slightly exceeds the experimental
198   observation (≈10,900 on day 40 on average), the numbers match closely once we
199   account for non-observed lineages due to the stochastic nature of sequencing (figure
200   2D).

201   **A-cell output**

202   According to the SAN model (table 1) a single A-cell has two options, either to divide
203   asymmetrically (probability $r_{A \to AN}/(r_{A \to AN} + r_{A \to N}) = 91\%$) or to cease proliferation
204   (probability $r_{A \to N}/(r_{A \to AN} + r_{A \to N}) = 9\%$). The likely range of additional N-cells
205   produced over the lifetime of an A-cell is thus 0 to 30 (95% quantile), with an
206   average of $r_{A \to AN}/r_{A \to N} = 10$. This slightly exceeds (about 3x) the observed
207   stochasticity of the neuronal output of RGCs in mice (Llorca et al., 2019) and is
208   consistent with a previously observed increase of intermediate progenitor divisions in
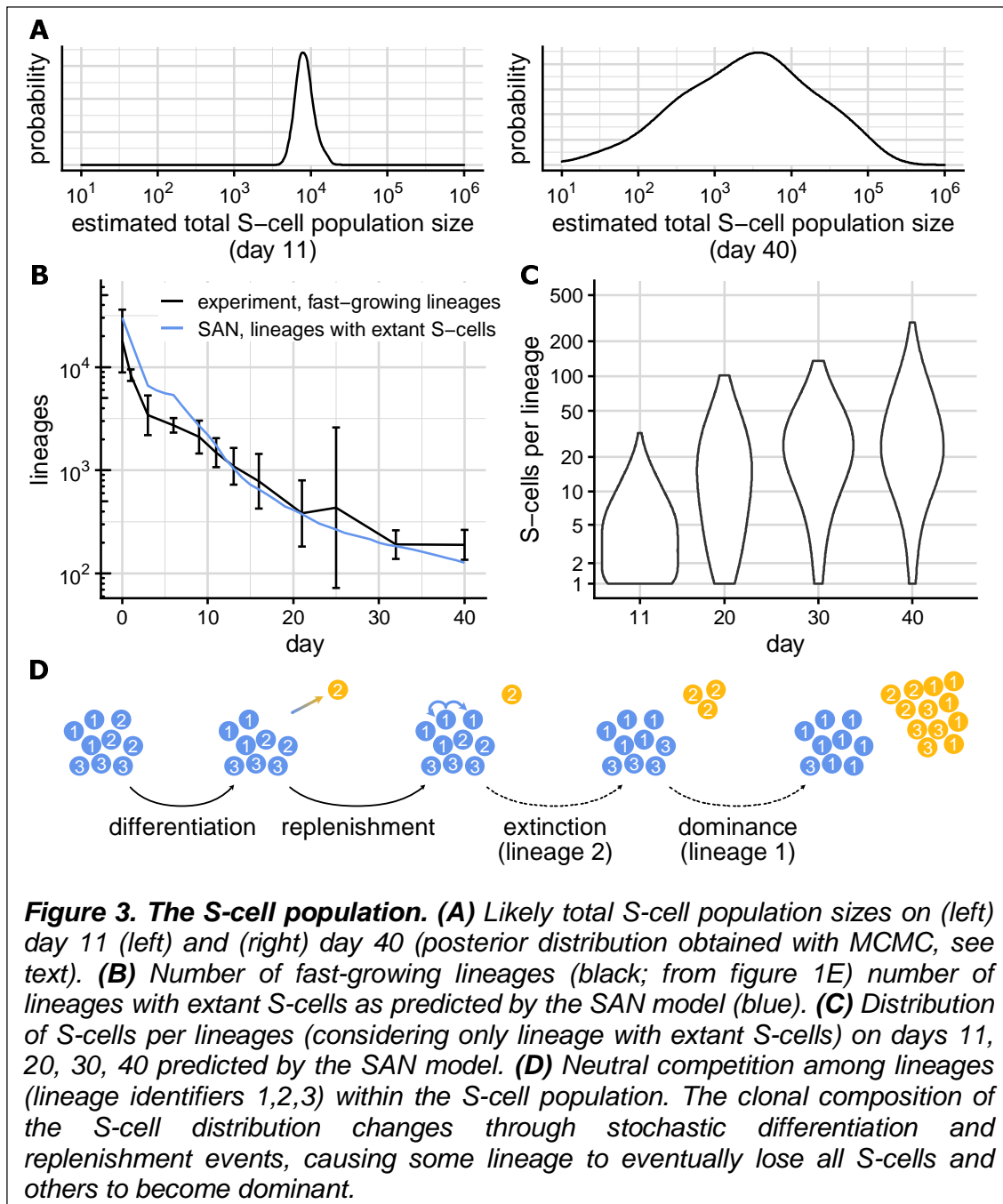209   humans.

210   **Predicted S-cell population size**

211   The MAP rate estimates (table 1) predict that organoids contain ≈9,500 S-cells on
212   day 11 and still ≈5,800 S-cells on day 40. To take the inherent ambiguity of the MAP
213   estimate due to the broadness of the posterior distribution into account, we

| day | S → S S | S → ∅ | S → N | S → A | A → A N | A → N | phase |
|---|---|---|---|---|---|---|---|
| 0-3 | - | 0.35 | 0.15 | - | - | - | EB formation |
| 3-6 | 0.6 | - | 0.15 | - | - | - | EB formation |
| 6-11 | 0.94 | - | 1.14 | - | - | - | neural induction |
| 11-40 | 1.68 | - | - | 1.69 | 0.71 | 0.07 | asymmetric division |

**Table 1. Division and conversion rates in the SAN model.** *Rates specify the number of expected events per cell and day*

214 computed the posterior distribution of these population sizes (figure 3A). We find that

215 the total S-cell population size on day 11 is well-defined up to a factor of at most 2

216 around 10,000 cells. On day 40, the estimates are more dispersed, owning to the

217 large cumulative effect that rates have over 30 days; yet while the exact population

218 size is difficult to estimate, the finding that organoids contain a significant number of

219 S-cells on day 40 is robust.



***Figure 3. The S-cell population.*** *(A) Likely total S-cell population sizes on (left) day 11 (left) and (right) day 40 (posterior distribution obtained with MCMC, see text). (B) Number of fast-growing lineages (black; from figure 1E) number of lineages with extant S-cells as predicted by the SAN model (blue). (C) Distribution of S-cells per lineages (considering only lineage with extant S-cells) on days 11, 20, 30, 40 predicted by the SAN model. (D) Neutral competition among lineages (lineage identifiers 1,2,3) within the S-cell population. The clonal composition of the S-cell distribution changes through stochastic differentiation and replenishment events, causing some lineage to eventually lose all S-cells and others to become dominant.*

**Fast-growing lineages contain S-cells**

While the total size of the S-cell population changes only slowly, its clonal composition changes rapidly. From the ≈13,700 lineages comprising the organoid from day 3 forward, ≈1,700 lineages still contain S-cells on day 11, and until day 40 that number has dropped to ≈100 (figure 3B). This drop in the number of lineages with extant S-cells is offset by an increase in the number of S-cells each of these lineages contains (figure 3C; average grows from ≈5 cells/lineage on day 11 to ≈36 cells/lineage on day 40).

The number of lineages with extant S-cells matches the number of lineages classified as fast-growing by our fast-slow model well (figure 3B). This highlights S-



***Figure 4. Lineage-specific S-cell extinction time determines linage size.*** *(Error bars show two standard deviations across the three replicates of Esk et al, shaded areas show the range between the 2.5% and 97.5% quantile across 5 billion simulations).* **(A)** *Lineage-specific growth trajectories under the SAN model stratified by the lineage's S-cell extinction time $T_S$. Plot show the most likely total lineage size (red), and number of S- (blue), A- (yellow) and N- (green) cells comprising the lineage.* **(B)** *S-cell extinction time $T_S$ vs. final lineage size on day 40. Plot shows simulation results (solid) and the analytical approximation $L(T_S)$ (dotted).* **(C)** *Recovering S-cell extinction times from lineage sizes on day 40. Plot shows the number of lineages reaching S-cell extinction on each day estimated using $L(T_S)$ from experimental data (black) and simulated data (red), and the true number of such lineages according to the SAN model.*

230 cells as being the main driver of lineage growth; as stated above a single
231 differentiating S-cells eventually on average produces 10 additional N-cells, and
232 once a lineage contains no more S-cells its growth will thus slow down and
233 eventually cease.

**Neutral competition shapes S-cell clonal composition**

235 Over time, not only does the average number of S-cells found within lineages with
236 extant S-cells grow, but so does the spread between the lineages' S-cells counts
237 (from about 1 - 30 S-cells per lineage on day 11 to about 1 – 300 S-cells per lineage
238 on day 40). The clonal composition of an organoid's S-cell population thus grows
239 more and more non-uniform over time. Under the SAN model, this change is the
240 result of *neutral competition* (figure 3D) amongst S-cells, a term introduced to
241 describe the population dynamics of stem cells within intestinal crypts (Snippert et
242 al., 2010).

243 Qualitatively, the dynamics of an organoid's S-cell population under neutral
244 competition mimic the population-genetic Moran model (Moran, 1958) in which
245 individuals (cells in our case) carrying different neutral alleles (lineage identifiers in
246 our case) are randomly removed (differentiate) and are replaced (through symmetric
247 division) by offspring of another randomly selected individual. Once the last S-cell of
248 a particular lineage has differentiated, the lineage cannot reappear within the
249 organoid's S-cell population. The observed disappearance (figure 3B) of lineages
250 from the organoids S-cell population is thus a result of more S-cells differentiating
251 than dividing due to random chance. Similarly, the observed growth of the remaining
252 lineages (figure 3C) results from more symmetric divisions than differentiations,
253 again due to random chance.

254 Using the SAN model, we now study the effects of neutral competition between
255 S-cells on the clonal composition quantitatively.

**Lineage-specific S-cell extinction times determine final lineage sizes**

257 Under the population-genetic Moran model, alleles eventually either disappear from
258 a population or become fixed. Tissue homeostasis driven by a stem cell population
259 under neutral competition likewise leads to eventual monoclonality, i.e. to all extant
260 cells being eventually derived from a single ancestral stem cell. In growing neural
261 tissue like cerebral organoids however, the lack of constant cell turn-over restricts

262    eventual monoclonality to S-cells. The clonal composition of the N-cell population

263    instead records the evolution of the S-cell's clonal composition over time; lineages

264    whose last S-cell was lost later and/or which contained more S-cells will contribute

265    more N-cells than lineages which die out quickly from the S-cell population.

266    To study the effects of S-cell extinction on lineage sizes quantitatively, we

267    stratified simulated lineage growth trajectories according to their *S-cell extinction*

268    *time* ($T_S$; the time at which a particular lineage loses the last S-cell). Lineages whose

269    S-cell population goes extinct at day $T_S = 13$ (figure 4A left) respectively day $T_S = 25$

270    (figure 4A middle) show diminished growth and a declining number of A-cells after

271    losing their S-cells at time $T_S$. In contrast, lineages whose S-cell population survives

272    past day 40 (figure 4A right) grow considerably faster and reach a considerably

273    larger size. Comparing the variations in lineage sizes on day 40 between S-cell

274    extinction time strata shows the variation due to $T_S$ to dominate the variations within

275    each stratum (figure 4B). Thus, while other random factors have some influence,

276    their influence on a lineage's sizes on day 40 is negligible compared to the time the

277    lineage loses its last S-cell.

278    Mathematical analysis of the SAN model yields the approximate expression

*(Eq. 1)*    $$L(\Delta T_S) = \left(\frac{1}{3} s_0 \Delta T_S + \frac{r_{S\to SS}}{2\sqrt{6}} \Delta T_S^2\right) r_{S\to A} \left(1 + \frac{r_{A\to AN}}{r_{A\to N}}\right)$$

279    for the final lineage size of a lineage comprising $s_0$ S-cells on day 11 ($s_0 \approx 5$ for

280    rates in table 1; and we assume $r_{S\to SS} \approx r_{S\to A}$) and whose S-cell population goes

281    extinct $\Delta T_S$ days later. We note that *final lineage size* here does not refer to the size

282    on day 40 (or any other particular point in time), but rather to the eventual size a

283    lineage will have reached when its growth ceases. While this does not exactly match

284    our simulation setup (we only simulate up to day 40) the approximate final linage

285    sizes $L(\Delta T_S)$ still matches the simulation results well (figure 4B).

**Recovering S-cell extinction times from final lineage sizes**

287    By solving the equation $L(\Delta T_S) = L_i$, the time at which a lineage lost its last S-cell

288    can be estimated from the final size ($L_i$) of that lineage. To gauge the reliability of

289    this approach, we applied it to a simulate lineage size distribution for day 40, and

290    found that it recovers the number of lineages that reached S-cell extinction on a

291    particular day well (figure 4C). When applied to the experimentally observed lineage

292 sizes on day 40, the estimated number of lineages reaching S-cell extinction lies

293 close to the SAN model prediction, but slightly exceeds it up to about day 30.

294 **Emergence of a Zipfian law**

295 If A- and N-cells are disregarded, the SAN model is equivalent to the well-studied

296 birth-death process, and in particular, the distribution of the S-cell extinction time $\Delta T_S$

297 is known (Feller, 1939). By translating this distribution via $L(\Delta T_S)$ into the

298 corresponding distribution of lineage sizes, the (approximate) distribution of final

299 lineage sizes (i.e. of lineages which have ceased growth) can be found. If we

300 consider only sufficiently large S-cell extinction times $\Delta T_S$, the probability of $\Delta T_S$ is

301 (approximately) proportional to $1/\Delta T_S{}^2$ and translation into lineage sizes via $L(\Delta T_S)$

302 yields a Zipfian law with $\alpha = 0.5$. This theoretical prediction matches the empirical

303 observation that lineage sizes approach a Zipfian law with $\alpha \approx 0.46$.

# Discussion

305 We have empirically observed lineages in cerebral organoids to initially grow fast

306 and roughly uniformly until some lineage-specific stopping time at which growth

307 slows down significantly or ceases altogether; and have found the size of slow or

308 non-growing lineages to follow a Zipfian power law with exponent $-1/\alpha$, $\alpha = 0.46$.

309 While the destructive nature of NGS-based lineage tracing prevents us from directly

310 observing lineages as they switch their growth regime, alternative hypotheses would

311 necessarily involve either very early fate decisions, or lineage-specific proliferation

312 rates to explain the large diversity of observed lineage sizes. Both alternative models

313 seem unlikely given that organoids are grown from a homogenous population of

314 stem cells.

315 To study the cause of the apparently random and lineage-wide switch of growth

316 regime we introduced the cellular *SAN* model. This model accurately recapitulates all

317 experimental data and shows that observed lineage growth dynamics to emerge

318 from neutral competition within a proposal long-lived population of roughly 10,000

319 symmetrically dividing stem cells (S-cells). Under the SAN model the apparently

320 lineage-wide switch of growth regime occurs despite the lack of either direct or

321 indirect (e.g., through spatial colocation) lineage-wide events. Instead, growth of a

322 lineage slows down and eventually ceases as the result of the lineage vanishing

from the S-cell population through neutral competition; lineage survival time within the organoid's S-cell population is thus the major determinant of lineage size.

The relationship between a lineage's survival times within the organoids S-cell population and the size it eventually attains can be expressed by a formula. Inverting this formula allows the history of the organoids S-cell population that was recorded within its clonal composition to be read; doing so we found for days 11-30 a slight excess of lineages reaching S-cell extinction in the experimental data compared to the SAN model. We hypothesize that this might point to gradual reduction of division and differentiation rates in organoids; Since the SAN model assumes constant rates between days 11 and 40, a gradual reduction of rates would cause the model to appear to fall behind at first, and then to catch up once the true rates have fallen below the model's rates.

While we found that we cannot estimate the rates of most division and differentiation events in the SAN model precisely, we could robustly determine the rates of S-cell division and differentiation to be almost identical. This implies that the population of symmetrically dividing cells in cerebral organoids is long-lived, and in particular that organoids still contain a population of symmetrically dividing cells after 40 days.

Furthermore, the similarity of the S-cell division and differentiation rates implies the existence of a mechanism that controls the S-cell population size by linking S-cell differentiation and subsequent replenishment through symmetric divisions. Yet that link must be stochastic in nature; if S-cells simply divided asymmetrically to produce A-cells, or if after symmetric division exactly one offspring always differentiated, no neutral competition between S-cells would occur and the observed large variability of lineage size would remain unexplained. In the terminology of Simons and Clevers (2011), the mechanism must thus be of the *population asymmetric* type.

Given the similarity between the population-level link of S-cell division and differentiation and the dynamics within stem cell niches in intestinal crypts (Snippert et al., 2010), we conjecture that similar structures located within proliferation centers called *neural rosettes* (Esk et al., 2020). might be responsible for balancing division and differentiation of S-cells.

To study the mechanism controlling S-cell population size in more detail, it needs to be probed experimentally by perturbing organoids at specific points in time and observing their response. If a fraction of cells is killed, different mechanism

would respond differently: A mechanism that relies on spatial constraints (i.e. stem cells being pushed out of a niche) would be expected to show a reduced rate of differentiations until the population has recovered. A regulatory mechanism which more directly links S-cell division to differentiation would respond differently; there we might expect the S-cell population to never reach its original size, but to instead increase its overall cell turn-over to make up for lost S-cells.

Since the SAN model accurately predicts the lineage sizes observed for cerebral organoids grown from wildtype cells, it is also useful both when planning organoid-based perturbation screens, and when analyzing the resulting data. During the planning phase the model makes it possible to judge the effect of proliferation phenotypes on final lineage size, and thus to estimate the statistical power of different screen designs. During statistical analysis of screening data, the model provides a baseline (null model), against which the sizes of (genetically) perturbed lineages can be compared.

To facilitate the adoption of the SAN model, we offer an implementation of the model both as an interactive online service (URL to be determined) as well as a R package (http://github.com/Cibiv/SANjar).

**Funding**

# Methods

## Total organoid sizes

For days 0 through 21, organoid sizes were measured using fluorescence-activated cell sorting (FACS). For days 11 through 40, organoid volumes were estimated from microscopy images, and translated into cell counts using the average number of cells per volume for days 11 through 21 where both FACS and volume measurements were available.

## NGS data processing

The lineage tracing data of Esk et al. (2020) was obtained from GEO (accession GSE151383, supplementary file GSE151383_LT47.tsv.gz), and organoids "H9-day06-03" and "H9-day09-01" removed as outliers. Based on the assumption that in all samples the most common lineage size is 1 cell, we located the mode of the log-transformed read count distribution for every sample and used it to normalize relative lineage sizes (reads) to absolute cell counts. The validity of the underlaying assumption is confirmed by the good agreement the sum of absolute lineage sizes and the FACS and area-derived estimates of total organoid size.

## Pareto index and fast-slow threshold estimation

For each organoid, we used the observed lineage sizes $l_1, \dots, l_n$ to estimate the Pareto equality index $\alpha$ and minimal lineage size $m$ with the maximum-likelihood estimator

$$\hat{m} = \min_i l_i, \qquad \hat{\alpha} = n \left( \sum_i \log \frac{l_i}{\hat{m}} \right)^{-1},$$

and computed the steady-state average $\bar{\alpha}$ from the alpha estimates of all organoids sequenced on day 11 or later. To find the fast-slow threshold $l_{\text{Th}}$ for a particular organoid, we first found intersect $d^{\text{Pareto}}$ such that the Pareto-induced rank-size powerlaw $\log_{10} L^{\text{Pareto}}(r) = -\bar{\alpha}^{-1} \log_{10} r + d^{\text{Pareto}}$ fits the size of the smallest observed lineage, and determined the smallest rank $R$ for which the actual lineage size $l_{(r)}$ matches or exceeds the power law $L^{\text{Pareto}}(r)$. We then fit a separate log-log-linear model $\log_{10} L^{\text{Large}}(r) = k \log_{10} r + d^{\text{Large}}$ to lineages with ranks $1, \dots, \sqrt{R}$ (which we

403  assume are surely not governed by the Pareto law), and set $l_{\text{Th}}$ to the size at which

404  the two laws intersect (meaning $l_{\text{Th}} = L^{\text{Pareto}}(r) = L^{\text{Large}}(r)$.

### SAN model simulation

406  The total number of S-, A- and N-cells that an organoid is predicted to comprise at

407  time $t$ is computed based on the deterministic SAN model (with rates $r_{S \to SS}$, $r_{S \to \emptyset}$,

408  $r_{S \to A}$, $r_{S \to N}$, $r_{A \to AN}$, $r_{A \to N}$ of these events occurring per cell and per day). The

409  deterministic SAN model is described by the ordinary differential equations (ODE),

$$s(0) = 30{,}000, \qquad a(0) = 0, \qquad n(0) = 0,$$

$$\dot{s} = (r_{S \to SS} - r_{S \to \emptyset} - r_{S \to A} - r_{S \to N})s,$$

$$\dot{a} = r_{S \to A}s - r_{A \to N}a,$$

$$\dot{n} = r_{S \to N}s + (r_{A \to AN} + r_{A \to N})a,$$

410  which can be solved analytically (for the time-homogenous case) and is then

411  evaluated separately for each time interval within which rates are constant  (The

412  initial number $s(0)$ of S-cells is set to 30,000 instead of 24,000 to account for a slight

413  excess in the number of observed lineages on day 0, likely due to a combination of

414  multiple labelling and sequencing artefacts).

415  To find the predicted lineage size distribution at time $t$, the stochastic SAN

416  model is simulated independently for each of the 30,000 lineages in an organoid.

417  The simulation proceeds in discrete time steps $\Delta t$, which are chosen small enough to

418  make the probability of a single cell undergoing two events negligible ($< 10^{-3}$).

419  Given the numbers $S_i(t), A_i(t), N_i(t)$ of S-, A-, N-cells comprising lineage $i$ at time $t$,

420  the number of cells $\Delta_e$ undergoing event $e$ is chosen from a Poisson distribution.

421  Specifically,

$$\Delta_{S \to SS} \sim \text{Poisson}(r_{S \to SS}S_i(t)\Delta t), \qquad \Delta_{S \to \emptyset} \sim \text{Poisson}(r_{S \to \emptyset}S_i(t)\Delta t),$$

$$\Delta_{S \to A} \sim \text{Poisson}(r_{S \to N}S_i(t)\Delta t), \qquad \Delta_{S \to N} \sim \text{Poisson}(r_{S \to N}S_i(t)\Delta t),$$

$$\Delta_{A \to AN} \sim \text{Poisson}(r_{A \to AN}A_i(t)\Delta t), \quad \Delta_{A \to N} \sim \text{Poisson}(r_{A \to N}A_i(t)\Delta t),$$

422  and the number of S-, A-, N-cells at time $t + \Delta t$ is then set to be

$$s_i(t + \Delta t) = s_i(t) + \Delta_{S \to SS} - \Delta_{S \to \emptyset} - \Delta_{S \to A} - \Delta_{S \to N},$$

$$a_i(t + \Delta t) = a_i(t) + \Delta_{S \to A} - \Delta_{A \to N},$$

$$n_i(t + \Delta t) = n_i(t) + \Delta_{S \to N} + \Delta_{A \to AN} + \Delta_{A \to N}.$$

423  Finally, the lineage size distribution $l_1, \ldots, l_{30{,}000}$ at time $t$ is found by summing up the

424  number of S-, A- and N-cells, $l_i(t) = s_i(t) + a_i(t) + n_i(t)$.

**Technical noise simulation**

The effect of PCR amplification and sequencing on the observed lineage sizes was simulated using a stochastic model of PCR amplification and sequencing (Pflug and von Haeseler, 2018) with parameters *PCR efficiency* and average *reads per molecule* (in our case per *lineage*). For every sampling time $t$, we simulated one read count (normalized to one read per cell on average) per lineage; parameters were *PCR efficiency* 35% (estimated from the day 0 data) and average *reads per lineage* $W l_i / \sum_i l_i$ for a linage comprising $l_i$ cells ($W$ is the median experimental library size for time $t$). The simulated read counts where then normalized to cells by division by the average number of reads per cell ($W / \sum_i l_i$).

**SAN rate estimation**

For days 0-3 and 3-6, rates which replicate the experimental data well were found by trial and error. For the remaining 6 biologically relevant rates (of S → S S and S → N between 6 and 11, and S → S S, S → A, A → A N and A → N between days 11 and 40) we computed the posterior distribution given experimentally observed total organoid sizes $\hat{c}^{(dayt)}$ (on days $t \in \mathcal{D} = \{0, 3, 6, 9, 10, 13, 14, 16, 17, 19, 21, 22, 25, 28, 31, 32, 35, 37, 38\}$) and ranked lineage sizes $\hat{l}_{(r)}^{(\text{day }11)}$, $\hat{l}_{(r)}^{(\text{day }40)}$ (on days 11 and 40, for ranks $r \in \mathcal{R} = \{1, 2, 5, 10, 15, 25, 40, 60, 100, 150, 250, 400, 600, 1000, 1500, 2500, 4000, 6000, 10000, 15000, 25000\}$). To account for biological differences between replicates we assumed that experimental observations are log-normally distributed around the SAN model predictions $c^{(dayt)}$ and $l_{(r)}^{(\text{day }11)}$, $l_{(r)}^{(\text{day }40)}$; the likelihood of the rate vector $\theta$ (comprising the 6 rates mentioned above) given the experimental data is thus

$$l(\theta) = -\frac{1}{2} \sum_{t \in \mathcal{D}} \left( \frac{\mu\left[\hat{c}^{(\text{day } t)}\right] - c^{(\text{day } t)}}{\sigma\left[\hat{c}^{(\text{day } t)}\right]} \right)^2 - \frac{1}{2} \sum_{t \in \{11,40\}} \sum_{r \in \mathcal{R}} \left( \frac{\mu\left[\hat{l}_{(r)}^{(\text{day } t)}\right] - l_{(r)}^{(\text{day } t)}}{\sigma\left[\hat{l}_{(r)}^{(\text{day } t)}\right]} \right)^2$$

where $\mu[\ldots]$ and $\sigma[\ldots]$ denote the mean respectively standard deviation across biological replicates. Rates were restricted to lie between 0 and 4 and *a priori* assumed to be equally probable; the posterior probability of $\theta$ is thus proportional to $l(\theta)$. To find this posterior distribution, we sampled 1,000 random rate vectors according to their likelihoods by simulating 1,000 Markov chains using pseudo-marginal Metropolis-Hastings Markov chain Monte Carlo sampling (Beaumont 2003;

454  Andrieu & Roberts, 2009; Warne et al., 2020). We then computed the maximal mode

455  of the (joint) posterior distribution with the mean-shift algorithm to obtain the MAP

456  estimates (table 1).

**Mathematical Analysis**

458  If we consider only S-cells, the SAN model corresponds to the well-known birth-

459  death process (Feller, 1939). We consider the diffusion approximation of this process

460  and restrict our mathematical treatment to day 11 and later where the rates of

461  symmetric division ($r_{S \to SS}$; birth) and of differentiation ($r_{S \to A}$; amounts to death since

462  we consider only S-cells) are similar enough to be considered identical ($r_{S \to SS} =$

463  $r_{S \to A} = \lambda/2$). The number of S-cells within a lineage at time $t$ (where $t = 0$ represents

464  day 11) is then governed by the stochastic differential equation (SDE)

$$ds(t) = \sqrt{\lambda s(t)}\, dW(t).$$

465  Using Onsager-Machlup theory (Onsager & Machlup, 1953; Dürr & Bach 1978) we

466  find the most probably trajectory of a linage that contains $s_0$ cells at $t = 0$ and loses

467  its last S-cell $\Delta T_S$ days later,

$$s_{\text{ext}}(t \mid s_0, \Delta T_S) = s_0 \left(1 - \frac{t}{\Delta T_S}\right)\left(1 + \rho \frac{t}{\Delta T_S}\right) \quad \text{where} \quad \rho = \frac{\Delta T_S}{s_0}\lambda\sqrt{\frac{3}{8}} - 1.$$

468  On average, a lineage grows by $\lambda/2$ A-cells per S-cell and per day, and over its

469  lifetime every A-cell will eventually produce $r_{A \to AN}/r_{A \to N}$ additional N-cells through

470  asymmetric division. Eventually, a lineage that starts out with $s_0$ S-cells and loses its

471  last S-cell $\Delta T_S$ days later will thus approximately grow to size

$$L(\Delta T_S) = \frac{\lambda}{2}\left(1 + \frac{r_{A \to AN}}{r_{A \to N}}\right)\int_0^{\Delta T_S} s_{\text{ext}}(t \mid s_0, \Delta T_S)\, dt.$$

472  Integration of this expression yields Eq. (1).

# References

474  Adamic, L. A., and Huberman, B. (2002). Zipf's law and the Internet. Glottometrics *3*, 143–150.

475  Andrieu, C., and Roberts, G. (2009). The Pseudo-Marginal Approach for Efficient Monte Carlo Computations.
476  The Annals of Statistics 37 (2).

477  Corominas-Murtra, B., Scheele, C.L.G.J., Kishi, K., Ellenbroek, S.I.J., Simons, B.D., van Rheenen, J., and
478  Hannezo, E. (2020). Stem cell lineage survival as a noisy competition for niche access. Proceedings of the
479  National Academy of Sciences *117*, 16969–16975.

480     Beaumont, M. (2003). Estimation of Population Growth or Decline in Genetically Monitored Populations. Genetics
481     164 (3), 1139-1160.

482     Dürr, D., and Bach, A. (1978). The Onsager-Machlup function as Lagrangian for the most probable path of a
483     diffusion process. Communications in Mathematical Physics *60*, 153–170.

484     Esk, C., Lindenhofer, D., Haendeler, S., Wester, R.A., Pflug, F., Schroeder, B., Bagley, J.A., Elling, U., Zuber, J.,
485     von Haeseler, A., et al. (2020). A human tissue screen identifies a regulator of ER secretion as a brain-size
486     determinant. Science *370*, 935–941.

487     Feller, W. (1939). Die Grundlagen der Volterraschen Theorie des Kampfes ums Dasein in
488     wahrscheinlichkeitstheoretischer Behandlung. Acta Biotheoretica *5*, 11–40.

489     Klingler, E., and Jabaudon, D. (2020). Do progenitors play dice? ELife *9*, e54042.

490     Lancaster, M.A., Corsini, N.S., Wolfinger, S., Gustafson, E.H., Phillips, A.W., Burkard, T.R., Otani, T., Livesey,
491     F.J., and Knoblich, J.A. (2017). Guided self-organization and cortical plate formation in human brain organoids.
492     Nature Biotechnology *35*, 659–666.

493     Llorca, A., Ciceri, G., Beattie, R., Wong, F.K., Diana, G., Serafeimidou-Pouliou, E., Fernández-Otero, M.,
494     Streicher, C., Arnold, S.J., Meyer, M., et al. (2019). A stochastic framework of neurogenesis underlies the
495     assembly of neocortical cytoarchitecture. ELife *8*, e51381.

496     Moran, P.A.P. (1958). Random processes in genetics. Mathematical Proceedings of the Cambridge Philosophical
497     Society *54*, 60–71.

498     Onsager, L., and Machlup, S. (1953). Fluctuations and irreversible processes. Physical Review *91*, 1505–1512.

499     Pflug, F.G. and von Haeseler, A. (2018). TRUmiCount: Correctly counting molecules using unique molecular
500     identifiers. *Bioinformatics* 34-18, 3137–3144.

501     Simons, B. D. and Clevers, H. (2011). Strategies for homeostatic stem cell self-renewal in adult tissues. Cell *145*,
502     851–862.

503     Snippert, H.J., van der Flier, L.G., Sato, T., van Es, J.H., van den Born, M., Kroon-Veenboer, C., Barker, N.,
504     Klein, A.M., van Rheenen, J., Simons, B.D., et al. (2010). Intestinal Crypt Homeostasis Results from Neutral
505     Competition between Symmetrically Dividing Lgr5 Stem Cells. Cell *143*, 134–144.

506     Warne, D. J., Baker, R. E., and Simpson M. J. (2020). A Practical Guide to Pseudo-Marginal Methods for
507     Computational Inference in Systems Biology. Journal of Theoretical Biology 496.

508     Zechner, C., Nerli, E., and Norden, C. (2020). Stochasticity and determinism in cell fate decisions. Development
509     *147*, dev181495.

510

511 ## Supplemental Information

512



**Figure S1. Replicate experiments.** *Replicate experiments based on the same organoid protocol show similar lineage size distributions as the data from Esk et al. (2020). Ranks of the Esk et al. data and SAN model predictions were scaled to account for an 1.7-fold increase in the number of detected lineages in the replicate experiments.*