# A meta-modelling assessment of "stochastic process mitotic mode" explanations for zebrafish retinal progenitor lineage outcomes

Michael Mattocks[1], Vince Tropepe[1*]

**1** Department of Cell and Systems Biology, University of Toronto, Toronto, ON, Canada

* vince.tropepe@utoronto.ca

## Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Curabitur eget porta erat. Morbi consectetur est vel gravida pretium. Suspendisse ut dui eu ante cursus gravida non sed sem. Nullam sapien tellus, commodo id velit id, eleifend volutpat quam. Phasellus mauris velit, dapibus finibus elementum vel, pulvinar non tellus. Nunc pellentesque pretium diam, quis maximus dolor faucibus id. Nunc convallis sodales ante, ut ullamcorper est egestas vitae. Nam sit amet enim ultrices, ultrices elit pulvinar, volutpat risus.

## Introduction

Mechanistic explanations (MEx) derive their utility from the resemblance of the conceptual mechanism's output to empirically observed outcomes. As maps to biological territories, biological MEx are naturally identified with the living systems they represent. Most well-developed MEx take the form of a model, whose internal structure is taken to reflect the underlying causal structure of a biological phenomenon. The nature of the causal relationship between a mechanistic model and the phenomenon it purports to explain remains a topic of active dispute in the philosophy of biology [1]. Biologists, nonetheless, usually accept that a model which explains empirical observations well (usually measured by statistical or information theoretic methods), and reliably predicts the results of interventions, bears a meaningful structural resemblance to the actual causal process giving rise to the modelled phenomenon.

Biological phenomena are notable for exhibiting both complex order and unpredictable variability. A significant challenge for MEx in multicellular systems is to explain how complex, highly ordered tissues, like those produced by neural progenitors, can arise from cellular behaviours with unpredictably variable outcomes. Stem cell biologists have traditionally resorted to Simple Stochastic Models (SSMs) in order to explain the observed unpredictable variability in clonal outcomes of putative stem cells [2,3]. Because SSMs are susceptible to Monte Carlo numerical analysis as Galton-Watson branching processes, they have been convenient explanatory devices, appearing in the literature for more than half a century. By specifying the probability distributions of symmetric proliferative (PP), symmetric postmitotic (DD), and asymmetric proliferative/postmitotic (PD) mitotic modes, SSMs allow cell lineage outcomes to be simulated.

SSMs are "stochastic" insofar as they incorporate random variables with defined probability distributions. As Jaynes has noted, "[b]elief ... that the property of being

'stochastic' rather than 'deterministic' is a real physical property of a process, that exists independently of human information, is [an] example of the mind projection fallacy: attributing one's own ignorance to Nature instead." [4] Despite this, macromolecular processes are often described as "stochastic" in the stem cell literature. Generally speaking, the behaviour of an SSM's random variable is taken to represent sequences of outcomes that are produced by multiple, causally independent events. Recently, the influence of "transcriptional noise" on progenitor specification has been identified as a candidate macromolecular process that may produce unpredictable variability in cellular fate specification. Therefore, one explanatory strategy for stem and progenitor cell function compares SSM model output to observed lineage outcomes, in order to argue that "noisy", causally independent events give rise to the proliferative and specificative outcomes of progenitor lineages.

In this report, we evaluate one of the best-developed of these explanations, proffered by William Harris' retinal biology group. We have dubbed this the "Stochastic Process Mitotic Mode Explanation" (SPMME) for zebrafish retinal progenitor cell (RPC) function. The SPMME is noteworthy because it claims: (1) to "provid[e] a complete quantitative description of the generation of a CNS structure in a vertebrate in vivo" [5]; (2) to have established the functional equivalency of embryonic RPCs and their descendants in the postembryonic circumferential marginal zone (CMZ) [6]; and (3) to have established the involvement of causally independent transcription factor signals in the production of unpredictably variable RPC lineage outcomes [7]. Moreover, the SPMME is taken to explain histogenetic ordering (the specification of some retinal neural types before others, notably the early appearance of retinal ganglion cells (RGCs)) in vertebrates, without reference to the classical explanation of a temporal succession of competency states [8]. These would be significant achievements with important consequences for both our fundamental understanding of CNS tissue morphogenesis and for retinal regenerative medicine. However, these models were not subjected to typical model selection procedures, nor has their output been examined using modern information theoretic methods, as advocated by mainstream model selection theorists [9].

In order to facilitate the critical evaluation of the two SSMs which form the SPMME's MEx for RPC function, we have re-expressed the models as cellular agent simulations conducted using the modular, open source CHASTE cell-based simulation framework [10]. We explored the structure of the SSMs, dubbed the He and Boije SSM respectively, compared to their explanatory forebear, dubbed the Gomes SSM [11]. This analysis strongly suggested that the introduction of an unexplained progression of temporal "phases" was largely responsible for the SPMME model fits to data. In order to investigate this possibility, we built an alternative model with a deterministic mitotic mode and compared its output to the He SSM. We found that the deterministic alternative model was a better explanation for the observations, demonstrating that SPMME fails when compared to alternatives. We therefore suggest that an explanatory approach based on the use of SSMs is incapable of distinguishing between theoretical alternatives for the causal structure of RPC lineage behaviours. Furthermore, by comparing the output of the models with novel postembryonic measurements of proliferative actvity, we find that the SPMME explanation cannot account for quantitative majority of retinal growth in the zebrafish, driven by the CMZ. Finally, we discuss the place of SSMs, the concept of "mitotic mode", and the role of "noise" in explaining RPC behaviour, and suggest ways to avoid the modelling pitfalls exemplified by the SPMME.

# Results and Discussion

The SPMME for zebrafish RPC function has been advanced using two SSMs. One first appears in He et al., and again, unmodified, in Wan et al. [5,6]; it is concerned primarily to explain lineage population statistics and time-dependent rates of the three generically construed mitotic modes (that is, PP, PD, DD). We have called this the He SSM. The second appears in Boije et al. [7]; its intent is both to introduce the role of specified macromolecules into the mitotic mode process, and to explain neuronal fate specification in terms of the process. We have called this the Boije SSM. The He model is directly descended an SSM advanced to investigate causally independent fate specification in late embryonic rat RPCs, formulated in Gomes et al. [11]. The Boije SSM differs substantially from the He and Gomes models, but inherits its general structure from the He SSM.

The metascientific analysis of biological explanations developing in time remains understudied. Perhaps most refined tool for global evaluation of biological theories (Schaffner's "Extended Theories"), treats biological explanations as hierarchically organised logical structures, after the fashion of Imre Lakatos, and proposes the use of Bayesian logic to distinguish between them [12]. However, biologists rarely offer explanations in this form; we rather prefer MEx, expressed in diagrammatic form or as mathematical model-objects.

In this report, we accept Fagan's view that MEx for the behaviour of stem and progenitor cells consist of assemblages ("mechanisms") of explanatory components which are understood to be causally organised by virtue of their intermeshing properties [1]. While Fagan treats SSMs seperately from macromolecular MEx, and we find that the Gomes SSM was not deployed in this role, the He and Boije SSMs were used as explanations for the behaviour of RPCs. As Feyerabend famously observed, scientists often operate as epistemological anarchists; the development of our explanatory logic is not bound by an identifiable set of rules, but rather arises organically from our scientific objectives, extrascientific context, etc [13].

We have therefore chosen to examine the structure of the Gomes, He, and Boije SSMs arranged in chronological order, to highlight how the explanatory logic of the zebrafish SPMME SSMs differs from their immediate ancestor, and from the traditional uses of SSMs in stem cell biology. We have diagrammatically presented these SSMs as MEx consisting of components describing the proliferative and fate specification behaviour of cellular agents. The proliferative and specificative components of the MEx are causally organised by their Faganian intermeshing property, mitotic mode. We have used abbreviations to denote important classes of model components and inputs, informed by the emphasis of Feyerabend on the persuasive role of metaphysical ingredients and auxiliary scientific material; these are as follows:

- MI - Model ingredient, making reference to some conceptual or metaphysical construct

- AS - Auxiliary scientific content

- RV - Random variable

- PM - Parameter measurement, model parameter set by measurements, independently of model output considerations

- PF - Parameter fit, model parameter set without reference to measurement, in order to produce model output agreement with observations

Numerous theoretical options to explain observed variability in RPC lineage outcomes have historically been considered. Harris has previously argued that these

belong to three cell-autonomous categories of process, in addition to extracellular ₁₂₄
influences: (1) a linear temporal progression of competency states, (2) asymmetric ₁₂₅
segregation of determinants during mitoses, and (3) intracellular "stochastic ₁₂₆
events" [14, 15]. All of the SSMs are used to argue for the predominant influence of the ₁₂₇
third type of process in RPC lineage outcomes, essentially by demonstrating that the ₁₂₈
output of the SSM resembles observations. ₁₂₉

## Gomes SSM: Ancestral Model of the SPMME SSMs ₁₃₀

The Gomes SSM is presented in Fig 1. The model's structure is straightforward; there ₁₃₁
are three independent random variables, drawing from empirically-derived probability ₁₃₂
distributions for the time each cell takes to divide, the mitotic mode of the division, and ₁₃₃
the specified neural fate of any postmitotic progeny. The random mitotic mode variable ₁₃₄
functions as the "intermeshing property" linking cycle behaviour to fate specification. ₁₃₅
The model thus represents a scenario where the processes governing proliferation, ₁₃₆
leading to cell cycle exit, and governing cell fate commitment are, at every stage, ₁₃₇
causally independent of each other and of their foregoing history of outcomes. The ₁₃₈
objective of the Gomes et al. study is to compare the lineage outcomes of dozens of ₁₃₉
individual E20 rat RPCs in clonal-density dissociated culture with the model output, ₁₄₀
developing earlier work in this system [16]. Although the model incorporates ₁₄₁
conventional proper time (clock-time), none of the RVs reference it to determine their ₁₄₂
values. The abstract "cells" represented by this model do not have any timer or any ₁₄₃
source of information about their relative lineage position. Fate specification is ₁₄₄
construed in terms of conventional histochemical markers of stable cell fates; only the ₁₄₅
neural types generated late in the retinal histogenetic order are represented, as these are ₁₄₆
the only neurons specified by E20 rat RPCs. ₁₄₇

### Fig 1. Structure of Gomes SSM
Structure of the Gomes SSM. pRPh, probability of rod photoreceptor specification. pBi,
probability of bipolar cell specification. pAm, probability of amacrine cell specification.
pMü, probability of Müller glia specification.

This SSM is explicitly built as a null hypothesis. It represents an extreme case in ₁₄₈
which all of the specificative and proliferative behaviours of an RPC are totally ₁₄₉
independent of all other RPCs and events in its clonal lineage. The stated purpose of ₁₅₀
this model is "to calibrate the data", serving "as a benchmark" [11], a purely ₁₅₁
hypothetical apparatus to produce sequences of causally unrelated lineage outcomes. ₁₅₂
Substantial deviations of the observed data from the fully independent events of the ₁₅₃
SSM may then be interpreted as causal dependencies between RPC outcome and their ₁₅₄
relative lineage relationships, as might be observed in a developmental "program" or ₁₅₅
algorithmic process. Finding that, with some exceptions, observed proliferative and ₁₅₆
specificative outcomes generally fall within the plausible range of Gomes SSM output, ₁₅₇
Gomes et al. conclude that RPC lineage outcomes in late embryonic rat RPCs seem to ₁₅₈
be dominated by causally independent events, suggesting that causally isolated "noisy" ₁₅₉
macromolecular processes may give rise to the unpredictable variability in the behaviour ₁₆₀
of these cells. ₁₆₁

## He SSM: Explaining variability in zebrafish neural retina ₁₆₂
## lineage size ₁₆₃

The He SSM, shown in Fig 2, is deployed in a explanatory role rhetorically identical to ₁₆₄
the Gomes SSM. The extent to which model output "captures.. aspects of the data" is ₁₆₅
taken to obviate the need for explicit "causative hypothes[e]s". On this account, only ₁₆₆

residual error between model output and observations may be ascribed to                              167
non-stochastic processes like "histogene[tic ordering] of cell types or a signature of early          168
fate specification". That is, if the model can be fit to observations, this is taken to                169
exclude the presence of any cell-autonomous temporal program, asymmetric segregation                  170
of fate determinants, extracellular influences, and the like.                                          171

**Fig 2. Structure of He SSM**
Structure of the He SSM. TiL, Time in Lineage.

    There are fewer direct empirical inputs to the He model's parameters than the         172
Gomes SSM, limited to the duration of the cell cycle. The close synchrony of zebrafish                173
RPC divisions is modelled by assigning sister RPCs the same cycle length, shifted by a                174
normal distribution of one hour width, in contrast to Gomes RPCs, which are treated as                175
fully independent. The lineage outcomes the He SSM is called on to explain are also                    176
notably different from those referred to by the Gomes SSM. Most notably, the Gomes                     177
SSM does not account for the early appearance of RGCs, which are not produced by the                   178
late E20 progenitors examined in that study. The zebrafish RPC lineages studied by He                 179
et al. produce all of the retinal neural types, including RGCs, which are typically                    180
produced by PD-type divisions. The He SSM does not model particular cell fates,                        181
supposing that mitotic mode is "decoupled" from fate specification. This decoupling                    182
cannot be anything more than a model construct, since, as noted (and as explicitly                     183
modelled by the Boije SSM), particular mitotic modes are definitely associated with                    184
particular neural fate outcomes. The mitotic mode RV linking cell cycle to fate                        185
specification in the Gomes SSM has thus subsumed the specification outcomes of RPC                     186
lineages entirely, and the model is concerned only to explain the sizes of lineages                    187
(marked by an inducible genetic marker at various times), and the observed progression                188
of mitotic modes in these early RPCs.                                                                  189
    There is, therefore, a temporal structure to the proliferative and specificative            190
behaviour of early zebrafish RPCs that dissociated late rat RPCs do not exhibit. A                     191
Gomes-type SSM, in which the RVs determining RPC behaviour are independent of any                      192
measure of time, cannot account for this temporal structure. In order to address this,                193
He et al., assume a linear progression of three phases which cells in each lineage pass                194
through, the timing of these phases being determined relative to the first division of the             195
RPC lineage, called here "Time in Lineage" or TiL. The values the mitotic mode RV                      196
may take on is determined by these TiL phases. The temporal structure of the phases                    197
and their effect on the mitotic mode RV are selected to produce a model fit. While He                  198
et al. acknowledge that the model therefore represents a "combination of stochastic and               199
programmatic decisions taken by a population of equipotent RPCs," no effort is made                    200
to determine the relative contribution of the model's stochastic vs. linear programmatic              201
elements. Instead, the purportedly stochastic nature of mitotic mode determination is                 202
emphasized throughout the report.                                                                      203

## Boije SSM: Explaining variability in zebrafish RPC fate outcomes                                    204 205

The second SPMME model advanced to explain zebrafish RPC behaviour, the Boije                          206
SSM, is diplayed in Fig 3. This model is primarily concerned with the lineage fate                     207
outcomes that the He SSM does not treat, while abandoning the explicit proper time of                  208
the He and Gomes SSMs in favour of abstract generation-counting. The mitotic mode                     209
model-ingredient now subsumes all RPC behaviours. Boije et al. make a laudable effort                  210
to specify the particular macromolecules ostensibly involved in determining mitotic                    211
mode, nominating the transcription factors (TFs) Atoh7, known to be involved in RGC                    212
specification, and Ptf1a, known to be involved in the specification of amacrine and                    213

horizontal cells, as primary candidates. The contribution of vsx2 is taken to determine the balance between PP and DD divisions late in the lineage, with the latter resulting in the specification of bipolar or photoreceptor cells. The binary presence or absence of these signals is determined by independent RVs structured by the phase structure present in the He SSM, translated into generational time from proper time.

**Fig 3. Structure of Boije SSM**
Structure of the Boije SSM. RPC, retinal progenitor cell. RGC, retinal ganglion cell. AC, amacrine cell. HC, horizontal cell. BC, bipolar cell. PR, photoreceptor. pNG-parameter representing contributions of Vsx1 and Vsx2 to specification of BC & PR fates in absence of Atoh7 or Ptf1a signals.

By this point, the SPMME has become a very different type of explanation from its Gomes SSM forebear. We are no longer dealing with RVs that model causally and temporally independent processes for different aspects of RPC behaviour. There is, rather, one temporally dependent process, the determination of mitotic mode, which is explained by random subsets of each lineage generation expressing particular TF signals. Where the Gomes SSM takes its parameters directly from empirical measurements of lineage outcomes, asking whether it is sufficient to assume that these are independently determined, the Boije SSM's parameters are derived solely from model fit considerations. In spite of these considerable differences, the explanatory role of the SSM is effectively the same: the model's fit to observations is taken as evidence of the predominant influence of "stochastic processes" determining mitotic mode on RPC behaviour. While Boije et al. acknowledge that whether some process is called "deterministic" or "stochastic" is "a matter of the level of description" [7] (i.e. is a property of the model-description and not of the physical process), the explanatory role of stochasticity for RPC lineage outcomes is, once again, emphasized throughout.

## Model selection demonstrates the SPMME is not the best available explanation for RPC lineage outcomes

From a modeller's perspective, it is notable that neither of the reports which use the He SSM [5,6], nor that using the Boije SSM [7] report in any detail their fitting procedures, nor do any of the above report any statistical measures of goodness-of-fit. Additionally, no models representing the alternative "theoretical options" available to explain variability in RPC lineage outcomes are compared to those advanced as evidence for "stochastic processes". Given that the He SSM and Boije SSM depart from the Gomes SSM by the addition of an unexplained temporal structure to the mitotic mode model ingredient, it is striking that the overall argument remains similar to Gomes et al.'s, despite the force of the latter deriving from the lack of such structures. While He et al. and Boije et al. acknowledge that their models involve both stochastic and linear programmatic elements, no effort is made to quantify their relative influence on model output, and macromolecular explanation is only applied to the stochastic elements. Moreover, the emphasis on this mitotic mode model construct increases with each successive model, to the extent that in the Boije model there are no other elements that are used to explain RPC behaviours. All cellular behaviours are, in effect, progressively collapsed into this stochastic mitotic mode concept. Finally, the He and Boije SSMs contain more parameters, which are determined by fewer empirical measurements than the Gomes SSM. Clearly, then, the SPMME SSMs are at significant risk of representing trivial overfits to their data, the modeller's equivalent of the "just-so" story, not representing any actually-existing macromolecular or cellular process.

Fortunately, we may employ standard statistical model optimisation and selection techniques to adjudicate this possibility. The most straightforward way to do so is to

reexamine the He SSM alongside a competing alternative model. The data to which the He SSM has been applied is particularly amenable to a scheme in which the models are fit to a "training" dataset, followed by a "test" dataset. In this case, the training data are the induced lineage size and mitotic mode rate data from He et al., and the test dataset are the Atoh7 morpholino observations from He et al., and the CMZ lineage size data from Wan et al. This allows us to test how well the He SSM, and an alternative, hold up under novel experimental conditions, without relying on a trivial ordering of goodness-of-fit to training data. Since the Boije SSM draws straightforwardly on the He SSM for its temporal phase-parameterisation and stochastic mitotic mode model ingredient, this analysis of the He SSM also bears directly on the validity of the Boije SSM.

As an alternative to the SPMME He SSM, we constructed a model that differs only in its construal of the process governing the RPC mitotic mode. Rather than variability arising from a stochastic-process mitotic mode changing across phases of fixed length, we simply supposed that mitotic mode is deterministic in each phase (guaranteed PP mitoses in the first phase, PD in the second, and DD in the third), but the phase lengths are variable between lineages and shift slightly between sister cells. That is, we represented this linear progression of deterministic mitotic mode phases using the same type of statistical construct the He SSM applies to model cell cycle length, to avoid introducing any novel or contentious elements into the model comparison. More specifically, each lineage has a first PP phase length drawn from a shifted gamma distribution, followed by a second phase length drawn from a standard gamma distribution. Upon mitosis, these phase lengths are shifted in sister cells by a normally distributed time period, exactly like cell cycle lengths in the He SSM.

We take this to be a reasonable representation of the classic suggestion that RPCs step through linear succession of competency phases, given the conceptual tools that the He SSM uses to model mitotic phenomena. If we suppose that this temporal program is governed by RPC lineages passing through a stereotypical series of chromatin configurations which allow for PP, then PD, then DD mitoses in turn, along with the associated competence to produce the particular cell fates associated with PD and DD mitoses, it seems entirely plausible to suggest that lineages differ in the lengths of time they occupy each state. This is particularly true if we concede that an SSM, foregoing any spatial modelling whatsoever, must necessarily abstract extracellular and spatial influences on these processes. Moreover, since these chromatin configurations must be broken down and rebuilt with each mitosis, the re-use of the "sister shift" model ingredient from the He SSM's cell cycle RV is congenial, representing the same sort of cell-to-cell variability that results in the small differences between sister cells in cycle timing.

While the code used to implement the SSMs mentioned above has not been published, the relevant reports provide enough detail to reconstruct these models in full, which we did using the CHASTE cell-based simulation framework, in order to provide transparent and reproducible implementations of the SSMs. Because the values selected for the He SSMs' parameters seem to have been selected as a series of rough estimates, with only the value for the probability of PD-type mitoses in the second model phase being varied to produce the fit, we suspected that the fit would not be at or near the local minimum for any reasonable loss function. That is, a model fit produced in this manner is likely to be located in a region of the parameter space that is highly sensitive to small perturbations, and therefore may depend strongly on implementation-specific idiosyncracies. He et al. report that they experienced difficulty in obtaining a good fit to their 32 hour induction data, with changes in the phase two PD probability producing large differences in fit quality. Unsurprisingly, when we rebuilt the He SSM with the original fit parameterisation, we substantially reproduced the original fit,

except for the 32 hour data, where the model output diverges substantially from that <sup>310</sup> reported in He et al. (see S1 Fig). Since this is clearly not the best fit available for the <sup>311</sup> He SSM, and we wish to directly compare the best fits (i.e. the parameterisations at <sup>312</sup> minima of some loss function) for the He SSM and our putative alternative model, we <sup>313</sup> used the simultaneous perturbation stochastic approximation (SPSA) algorithm [17] to <sup>314</sup> optimise both model fits. SPSA is particularly convenient for complex, multi-phase <sup>315</sup> models like the He SSM, because no knowledge of the relationship between the model's <sup>316</sup> parameters and the loss function is required. We used Akiake's information criterion <sup>317</sup> (AIC) as the loss function to be minimised, in order to provide a rigorous comparison <sup>318</sup> between the two differently-parameterised models (the deterministic alternative has two <sup>319</sup> fewer parameters than the He SSM). <sup>320</sup>

After (re)-fitting to the training dataset, we calculated AIC for the He SSM <sup>321</sup> (hereafter SM, for stochastic model) and our deterministic mitotic mode alternative <sup>322</sup> (hereafter DM), for both training and test datasets. The combined results are displayed <sup>323</sup> in Fig ??, with the output of the two models being presented separately in S2 Fig and <sup>324</sup> S3 Fig. Remarkably, the DM closely recapitulates the output of the SM for both <sup>325</sup> datasets. Moreover, while the more highly-parameterised SM permits a better fit to the <sup>326</sup> training data, the DM proves to be a better fit to the test dataset, as summarised in <sup>327</sup> Table 1. This is, therefore, a classic case of model overfitting: a higher-parameter model <sup>328</sup> fits some training dataset better than a simpler model, but fails upon challenge with a <sup>329</sup> new dataset. Given this model selection scheme, it is plain that we should choose the <sup>330</sup> DM over the SM as a superior explanatory model, both on the basis of explanatory <sup>331</sup> power and of Occam's Razor. <sup>332</sup>

### Fig 4. Model comparison: the SPSA-optimised He SSM and a deterministic alternative

Empirical observations (black crosses and bars) and SPSA-optimised model output (magenta, stochastic mitotic mode model; green, deterministic mitotic mode model). Model output in panels A-G is displayed as mean $\pm$ 95% CI. Panels A-F: Training dataset, to which models were fit. Panels A-C: Probability of observing lineages of a particular size ("count") after inducing single RPC lineage founders at (A) 24, (B) 32, and (C) 48 hours with an indelible genetic marker, tallying their size at 72hpf. Panels D-F: Probability density of mitotic events of modes (D) PP, (E) PD, and (F) DD over the period of retinal development observed by He et al. Panel G: Probability of observing lineages of a particular size ("count") originating from the CMZ; RPCs are taken to be resident in the CMZ for 17 hours before being forced to differentiate, lineage founders are assumed to be evenly distributed in age across this 17 hour time period. Panel H: Average clonal lineage size of wild type and Ath5 morpholino-treated RPCs. Ath5mo treatment is taken to convert 80% of PD divisions, which occur in the second phase of the He model, to PP divisions, resulting in larger lineages. Panel I: Ratio of even to odd sized clones in wild type and Ath5 morpholino-treated RPCs.

**Table 1. AIC values for models assessed against training and test datasets**

| Model | Training AIC | Test AIC |
|---|---|---|
| Stochastic (He fit) | -93.26 | 255.64 |
| Stochastic (SPSA fit) | **-93.80** | 92.24 |
| Deterministic (SPSA fit) | -69.33 | **86.60** |

AIC: Akiake's information criterion. Lower values reflect better model explanations of the noted dataset. Lowest values are noted in bold.

We therefore conclude that, had the Harris group employed standard model selection <sup>333</sup> procedures, comparing plausible alternatives to their favoured explanation, they would <sup>334</sup>

have been forced to conclude that the SPMME is not the best available explanation for variability in zebrafish RPC lineage outcomes. It is clear from this analysis that the assumed, but unexplained, linear succession of mitotic mode phases is what provides the overall structure of the model output. Variability in lineage outcomes may be supplied by entirely different model ingredients without any loss of explanatory power (indeed, with some improvement, in our case). The rhetorical emphasis on a stochastic process governing mitotic mode, ostensibly ruling out other types of explanation, is unjustified. It is likely that any number of different types of SSMs, representing other sorts of processes (such as asymmetric segregation of fate determinants, or differential spatial exposure to extracellular signals), can produce identical model output. Given this, we conclude that the SSM model type is simply inadequate for the task of locating the source of variability in RPC lineage outcomes.

## SPMME SSMs cannot explain the post-embryonic phase of CMZ-driven zebrafish retinal formation

The zebrafish retina, like other fish retinas, and unlike the mammalian retina, continues to grow long after the early developmental period, indeed, well past the organism's sexual maturity. This may be unsurprising, given that zebrafish increase in length almost ten-fold over the first year of life [18], necessitating a continuously growing retina during this period. In fact, the quantitative majority of zebrafish retinal growth occurs post-metamorphosis, outside of the early developmental period. This growth occurs due to the persistence of a population of proliferative RPCs present in an annulus at the periphery of the retina, called the ciliary or circumferential marginal zone (CMZ), which plate out the retina in annular cohorts. A typical "tree-ring" analysis from our studies, marking the DNA of cohorts of cells contributed to the retina at particular times with indelible thymidine analogues, is shown in Fig 5. It is immediately obvious from such experiments that the structure of the zebrafish retina is quantitatively dominated by contributions from the CMZ in the period between one and three months of age. Why peripheral RPCs in zebrafish remain proliferative, while those in mammals are quiescent [19], and whether and how their behaviour might differ from embryonic RPCs, remain unresolved. Answers to these questions may have significant fundamental and therapeutic implications, especially given e.g. the possibility of harnessing endogenous, quiescent, peripheral RPCs in humans for regenerative retinal medicine.

**Fig 5. Most zebrafish retinal neurons are contributed by the CMZ between one and three months of age**
Maximum intensity projection derived from confocal micrographs of a zebrafish whole retina dissected at 5 months post fertilisation (mpf). The animal was treated with BrdU at 1, 2, 3, 4, and 5 mpf for 24hr. Anti-BrdU staining of the whole retina reveals the extent to which the CMZ (which is responsible for retinal growth after approximately 72hpf) contributes new neurons in tree-ring fashion, extending out from the center of the retina (CR), which is formed before 72hpf. Monthly cohorts are labelled appropriately. Scale bar, 100 $\mu$m.

Indeed, the SPMME SSMs were originally brought to our attention when the He SSM was deployed in Wan et al. [6], with the claim that the He SSM explains the behaviour of CMZ RPCs. Wan et al. argue that a slowly mitosing population of bona fide stem cells, at the utmost retinal periphery, divides asymmetrically to populate the CMZ with He-SSM-governed RPCs. In other words, the usual suggestion that CMZ RPCs undergo a somewhat different process than embryonic RPCs, perhaps recapitulating across the peripheral-central axis some progression of states or lineage phases that embryonic RPCs pass through in time [20], is repudiated in favour of one

model which describes the behaviour of all RPCs throughout the life of the organism, with the addition of a small population of stem cells to keep the CMZ stocked with RPCs. If so, this might suggest that the problem of activating quiescent stem cells in the retina is simply that- one need only sort out how to throw the proliferative switch in these cells, since the proliferative and fate specification behaviours of the resultant RPCs will reliably be the same as those observed in development.

While we determined that the SPMME is not a particularly good explanation for RPC lineage outcomes, we still felt that the SSMs associated with this explanation might be used to elucidate this point. In particular, if it is the case that the He SSM provides good estimates of lineage size and proliferative dynamics in the early zebrafish retina, it should be possible, using this model, and the estimates of putative stem cell proliferative behaviour provided by Wan et al., to simulate the population dynamics of the CMZ, at least through the first few weeks of the organism's life.

In pursuing this point, we noted a peculiar feature of the He SSM not documented by any of the SPMME reports: the cell cycle model overstates the *per-lineage* rate of mitoses by as much as a factor of 3. That is, the mitotic mode rate data presented in He et al., recapitulated here in Fig 4, panels D, E, and F, and used to optimise both the SM and DM, are probability density functions that are not standardised on a per-lineage basis. These data simply indicate the distribution of mitotic events of a particular type. When we take all of the mitotic events documented by He et al. and calculate the probability of any such event occuring per lineage, per hour, we obtain the values presented in Fig 6.

### Fig 6. Per-lineage probabilities of mitoses, He et al. observations compared to model output
Probability of observing a mitotic event over the period studied by He et al., per hour, per lineage. Model output is presented as mean ± 95% CI.

Similarly, when we performed cumulative thymidine analogue labelling of the 3dpf CMZ, (using the assumptions of Nowakowski et al. [21], which treat the proliferating population as a homogenous, and dividing asymmetrically; these assumptions are flawed but adequate for a rough estimate) we obtain an average cell cycle length of approximately 15 hours, more than twice as long as the He SSM's mean cycle length. These data are displayed in . Therefore, the He SSM (in both its original and refit parameterisation, as well as the deterministic alternative, since they all rely on the same proliferative model elements) substantially overstates the proliferative potential of both embryonic and early CMZ RPCs. Still, because we observed a massive build-up of proliferating RPCs between two the first few weeks of post-embryonic CMZ activity involve a massive build-up of proliferating RPCs between two and four weeks post-fertilisation, we thought this might actually suit this later context better.

In order to produce estimates of total annular CMZ population in zebrafish retinas over time, we counted proliferating RPCs present in central coronal sections of zebrafish retinas throughout the first year of life, treating these as samples of the annulus, and calculated the total number of cells that would be present given the diameter of the spherical lens measured at these times. Our simulated CMZ populations were constructed at 3dpf by drawing an initial population of RPCs governed by the He SSM (using the original fit parameters, which further exaggerate the proliferative potential of these lineages) from the observed distribution. An additional number of immortal stem cells amounting to one tenth of this total was added. This is likely an overestimate, given that these putative stem cells are thought to be those in the very peripheral ring of cells around the lens, of which typically two to four may be observed in our central sections with an average of over one hundred proliferating RPCs. Moreover, these simulated stem cells were given a mean cycle time of 30 hours, proliferating about twice

as quickly as Wan et al. suggest. Finally, to reflect the fact that these stem cells 422
contribute to the retina in linear cohorts [22], and more of them are therefore required 423
as the retina grows, the stem cells were permitted to divide symmetrically when 424
necessary to maintain the same density of stem cells around the annulus of the lens. 425
Two hundred such CMZ populations were simulated across one year of retinal growth. 426

The results of these simulations are displayed in Fig 7, overlaid over CMZ 427
population estimates derived from observations. Given the generous parameters of the 428
population model, consistently overestimating the proliferative potential of embryonic 429
and early CMZ RPCs, it is surprising that the He SSM proves completely unable to 430
keep up with the growth of the CMZ in the first two months of life. Indeed, the 431
unrealistically active stem cell population the simulated CMZs are provided with is 432
unable to prevent a near-term collapse in RPC numbers, only catching up after the 433
months later as the number of stem cells increases with the growth of the lens. Thus, 434
even given permissive model parameters, the He SSM is not able to recapitulate the 435
quantitatively most important period of retinal growth in the zebrafish. 436

**Fig 7. CMZ population of proliferating RPCs: estimates from observations
and simulated Wan-type CMZs**
Toroidal CMZ population was estimated from empirical observations of central coronal
cryosections stained with anti-PCNA, a marker of proliferating cells, standardising
by the diameter of the spherical lens observed in these animals. "Wan-type" CMZs
were initialised with a number of RPCs drawn from from the observed 3dpf population
distribution, and given an additional $1/10^{\text{th}}$ of this number of immortal, asymmetrically
dividing stem cells, as described in the text, and simulated for the first year of life. All
data is presented as mean $\pm$ 95% CI.

This analysis strongly suggests that observations of RPCs in embryogenesis and 437
early larval development are unlikely to provide a good quantitative model of the 438
development of the zebrafish retina, even abstracting away spatial and extracellular 439
factors as SSMs necessarily do. Our data point to a second, quantitatively more 440
important phase of retinal development, between approximately one and four months of 441
age, in which RPCs are far more proliferatively active than in early development. It is 442
likely that models that closely associate proliferative behaviour with fate specification, 443
like those of the SPMME, will necessarily be unable to explain this period. Recent 444
evidence suggests that mitotic and fate specification behaviours in RPCs may be 445
substantially uncoupled [23]. This would permit the CMZ population to scale 446
appropriately with the growing retina. Alternatively, it is also possible that RPCs are 447
substantially heterogenous with respect to their proliferative behaviour, and that this 448
heterogeneity is mainly apparent later in development. 449

# Conclusion 450

Simple stochastic models are familiar tools for stem cell biologists. Introduced by Till, 451
McCulloch, and Siminovitch in 1964, they proved their utility in describing variability 452
in clonal lineage outcomes of putative stem cells, originating from macromolecular 453
processes beyond the scope of cellular models (and beyond the reach of the molecular 454
techniques of the time). More prosaically, their simplicity afforded computational 455
tractability in an era when processing time was relatively scarce, allowing early access 456
to Monte Carlo simulation techniques. That said, their abstract nature necessarily 457
emphasizes an aspatial, lineage-centric view of tissue development, which cannot 458
account for the generation of structurally complex tissues like retinas beyond cell 459
numbers, and, perhaps, fate composition. As a result of this relatively loose relationship 460

to the complex morphogenetic environment, there are few constraints on the model ₄₆₁
configurations that may produce similar outputs. As we hope has been made plain by ₄₆₂
the analyses herein, this can result in modellers being led astray by their apparently ₄₆₃
good fits to observations. This is one reason why it has long been unacceptable in other ₄₆₄
biological modelling communities (particularly, among ecologists) to present the fit of ₄₆₅
one highly parameterised model as evidence for some theory; it is very rare that there is ₄₆₆
not a model representing an alternative theory that cannot be made to produce a ₄₆₇
reasonable model fit. Indeed, as we demonstrate here, the use of standard model ₄₆₈
selection techniques may make plain that some model form is unable to distinguish ₄₆₉
between competing theories. ₄₇₀

To some extent, this can be ameliorated by careful attention to the particular ₄₇₁
macromolecular or cellular referents that particular model constructs represent. The ₄₇₂
"mitotic mode" construct present in all SSMs is a particularly ambiguous and ₄₇₃
problematic one in this regard. "Mitotic mode" is not, itself, a property of a mitotic ₄₇₄
event, but is rather a retrospective classification of the event after an experimenter ₄₇₅
observes whether progeny resulting from the event continue to proliferate. Its original ₄₇₆
appearance in SSMs was simply to allow calculation of clonal population sizes; it was ₄₇₇
never intended to represent a particular type of process or "decision" made by cells at ₄₇₈
the time of mitosis to continue proliferating or not. Any number of pre- or post-mitotic ₄₇₉
signals and processes may result in a particular mitosis being classified as PP, PD, or ₄₈₀
DD, without anything about the event itself determining this. The use of such a ₄₈₁
retrospective classification, rather than the identification of some physical property of ₄₈₂
the mitotic event (such as the asymmetric inheritance of fate determinants), is ₄₈₃
straightforwardly a concession that the actual macromolecular determinants of cellular ₄₈₄
fate are outside the scope of the model. ₄₈₅

The SPMME represents an attempt to connect this abstract, retrospective model ₄₈₆
construct to observations of noisy gene transcription [24]. Motivated by the observation ₄₈₇
that momentary mRNA transcript expression in RPCs is highly variable from ₄₈₈
cell-to-cell [25], the suggestion is that this transcriptional variability may be responsible ₄₈₉
for the observed variability in RPC lineage outcomes. Since the extent of transcription ₄₉₀
noise may be "tuned" to a degree by e.g. promoter sequences [26], this is a legitimate ₄₉₁
target for Darwinian evolution. There are two significant problems with identifying the ₄₉₂
SPMME's stochastic mitotic mode with this type of process. The first we have ₄₉₃
demonstrated by constructing an alternative model with deterministic mitotic mode but ₄₉₄
variable phase lengths: the source of variability may be located elsewhere without ₄₉₅
compromising the explanatory power of the model. It is therefore impossible to ₄₉₆
determine what sort of process might give rise to variability in RPC lineage outcomes ₄₉₇
by using SSMs. The second, more fundamental problem is that a causal explanation of ₄₉₈
the presence of the signal, for which noise is a property, is elided entirely in favour of ₄₉₉
emphasis on "stochasticity". This is most apparent in the Boije SSM. In this model, ₅₀₀
Atoh7 and Ptf1a TFs are available to provide their noisy signal in the 4th and 5th ₅₀₁
lineage generations, and at no other time. We do not dispute that a noisy signal may ₅₀₂
contribute to variability in that signals' effects; rather, we suggest that it is the ₅₀₃
temporal structure of such a signal (assuming this structure can be empirically ₅₀₄
demonstrated, rather than assumed) that calls for causal explanation. The SPMME ₅₀₅
reports explicitly disclaim the necessity for "causative hypotheses" in the case that a ₅₀₆
stochastic model provides a good fit to observations. As Jaynes remarked in his classic ₅₀₇
text on probability theory, "[stochasticity] is always presented in verbiage that implies ₅₀₈
one is describing an objectively true property of a real physical process. To one who ₅₀₉
believes such a thing literally, there could be no motivation to investigate the causes ₅₁₀
more deeply ... and so the real processes at work might never be discovered." [4] ₅₁₁

In identifying the model construct with the physical processes determining RPC ₅₁₂

outcomes, the SPMME obscures what seems to us to be the primary lesson to be drawn ₅₁₃
from the Harris groups' beautiful in vivo studies of zebrafish RPCs. That is, compared ₅₁₄
to late rat RPCs in dissociated clonal culture, RPCs in intact zebrafish retinas produce ₅₁₅
far more orderly outcomes. To the extent that these outcomes are variable, the source ₅₁₆
of variability remains unidentified. We suggest that, in order to produce good ₅₁₇
explanations for both the order and unpredictable variability exhibited in RPC ₅₁₈
behaviours, more sophisticated models and more rigorous modelling practices are ₅₁₉
required. Resort to "stochasticity" as an explanatory element should not be made in the ₅₂₀
absence of model comparisons that rule out alternatives with well-defined causal ₅₂₁
structures, lest we fall into the trap Jaynes warned us about. ₅₂₂

# Materials and methods ₅₂₃

## Zebrafish husbandry ₅₂₄

Zebrafish used in this study were of a wild type AB genetic background. Embryos were ₅₂₅
derived from pairwise crossings. Larvae were collected upon crossing and held in a dark ₅₂₆
incubator at 28°C. At 1dpf, embryos were treated with $100\mu$l of bleach diluted in 170ml ₅₂₇
of embryo medium for 3 minutes, rinsed and dechorionated. Embryo medium was ₅₂₈
changed at 3dpf. After this, all animals were maintained at 28°C on a 14-hour ₅₂₉
light/10-hour dark cycle (light intesity of 300 lux) in an automated recirculating ₅₃₀
aquaculture system (Aquaneering). Animals were reared using standard protocols [27]. ₅₃₁
Fish were sacrificed by tricaine overdose at the appropriate timepoints indicated in the ₅₃₂
text. All animal experiments were performed with the approval of the University of ₅₃₃
Toronto Animal Care Committee in accordance with guidelines from the Canadian ₅₃₄
Council for Animal Care (CCAC). ₅₃₅

## Proliferative RPC Histochemistry ₅₃₆

### Anti-PCNA histochemistry ₅₃₇

In order to assay the number of proliferating cells in zebrafish CMZs, we used ₅₃₈
anti-PCNA histochemical labelling, as PCNA is, in zebrafish, detectable throughout the ₅₃₉
cell cycle in proliferating cells. After sacrifice as described above, fish (or their ₅₄₀
razor-decapitated heads, in the case of fish older than 30dpf) were fixed in 1:9 37% ₅₄₁
formaldehyde:95% ethanol at room temperature (RT) for 30 minutes, followed by ₅₄₂
overnight incubation in a fridge at 4°C. Subsequently, the samples were removed from ₅₄₃
fix and washed 3 times in PBS. They were then cryoprotected by successive 30 minute ₅₄₄
rinses at RT on a rocker in 5%, 13%, 17.5%, 22%, and 30% sucrose in PBS. Samples ₅₄₅
were allowed to incubate overnight at 4°C in 30% sucrose. The next day, samples were ₅₄₆
removed from the sucrose rinse and infiltrated with 2:1 30% sucrose:OCT compound ₅₄₇
(TissueTek) for 30 minutes at RT. Samples were then embedded for cryosectioning and ₅₄₈
frozen. $14\mu$m coronal cryosections were cut through the fish's heads and collected on ₅₄₉
Superfrost slides (Fisherbrand). These were stored in a freezer at -20°C until staining. ₅₅₀
  Staining was begun by allowing the slides to dry briefly at RT. Sections were ₅₅₁
outlined with a PAP pen, then rehydrated in PBS for 30 minutes at RT. Subsequently, ₅₅₂
sections were blocked in 0.2% Triton X-100 + 2% goat serum in PBS for 30 minutes at ₅₅₃
RT. This was followed by incubation in mouse monoclonal (PC10) anti-PCNA primary ₅₅₄
antibody (Sigma), diluted 1:1000 in blocking solution, at 4°C overnight. The next day, ₅₅₅
primary antibody was removed, and slides were rinsed five times in PDT (PBS + 1% ₅₅₆
DMSO + 0.1% Tween-20). Cy5-conjugated goat-anti-mouse secondary antibody ₅₅₇
(Jackson Laboratories), diluted 1:100 in blocking solution, was applied to the sections, ₅₅₈
which were incubated at 37°C for 2 hours. Secondary antibody was removed, followed by ₅₅₉

five rinses in PDT as above. Sections were counterstained with Hoechst 33258 diluted to 100 $\mu$g/mL in PBS for 15 minutes at RT. Counterstain was then removed, five rinses in PDT were performed as above, followed by a final rinse in PBS. Slides were then mounted in ProLong Gold antifade mounting medium (ThermoFisher Scientific), coverslipped, and sealed with clear nail polish. Slides were kept at 4°C until imaging.

### Cumulative thymidine analogue labelling for estimation of CMZ RPC cell cycle length

In order to obtain an estimate of cell cycle length in CMZ RPCs at 3dpf, we used cumulative thymidine analogue labelling following the method of Nowakowski et al. [21]. 3dpf zebrafish were divided into ten groups of n5 animals and held in 10 mM EdU for 0.5, 1.5, 2.5, 3.5, 4.5, 5.5, 7.5, 8.5, 9.5, and 10.5 hours, after which they were sacrificed as described above. Samples were processed as described under "Anti-PCNA histochemistry", except that goat-anti-mouse Cy2-conjugated secondary antibody (Jackson Laboratories) was used to label anti-PCNA, after which the Alexa Fluor 647-conjugated azide from a Click-iT kit (ThermoFisher Scientific) was added to the sections for 30 minutes at room temperature to label EdU-bearing nuclei. This was followed by five PDT rinses as described above, and resumption of the protocol with counterstaining, mounting, etc. The raw data is available in \empirical_data\cumulative_edu.xlsx

### Whole retina thymidine analogue labelling of post-embryonic CMZ contributions

We used thymidine analogue labelling as an indelible marker of CMZ contributions to the retina in Fig 5. Because thymidine analogues are incorporated into DNA during S-phase, postmitotic progeny of proliferating RPCs which take up the label may be detected at any point after this treatment. Dosing zebrafish with a thymidine analogue for a limited period of time ensures that the label is diluted to undetectable levels in cells which continue to proliferate. Repetitive dosing thus gives rise to labelled cohorts which represent a limited set of generations of cells which become postmitotic after the dose is provided. n=3 fish were held in 10mM BrdU (Sigma) diluted in facility water for 8 hours at 30, 60, 90, 120, and 150 days post fertilisation. 10 days after the last BrdU incubation, the fish were sacrificed as described above. One representative whole retina dissection is presented.

## Confocal microscopy and image analysis

All histochemically processed samples (coronal sections of retinas and whole retinas) were imaged on a Leica TCS SP5 II laser confocal imaging microscope, using the Leica LAS AF software for acquisition. For coronal sections, central sections were identified by their position in the ribbon of cryosections. For Fig 7, the central section was imaged together with the flanking section on either side of it. Confocal data was analysed using Imaris Bitplane software (v. 7.1.0). Cell counting analyses were performed by segmenting the relevant channels using the Surfaces tool to produce counts of PCNA positive (Figs 7, S4 Fig) or PCNA/EdU double-positive (S4 Fig) nuclei. Lens diameters were measured in the DIC channel using the linear measurement tools, across the widest point of the section of spherical lens.

## Estimation of CMZ toroidal population size

In order to estimate the total population of the CMZ in zebrafish eyes sampled throughout the first year of life, we counted PCNA-positive cells in the dorsal and

ventral CMZ of central and flanking sections of these eyes as described above. We added the dorsal and ventral populations and took the average of the central and flanking section per-section populations to reduce the possibility of sampling error. We treated this as a sample of a torus mapped to the $14\mu$m sphere segment of lens intersected by the average section. As the surface area of this zone is given by $S = \pi \times d_{lens} \times h_{section}$, and the total surface area of the lens by $A = \pi \times d_{lens}^2$, where $d$ is the diameter of the lens and $h$ the section thickness, we may calculate an estimate of the total population by $pop_{total} = \frac{pop_{section}}{S} * A$. We therefore obtained our toroidal CMZ population estimates using our measurements with the simplified formula $pop_{total} = \frac{pop_{section} \times d_{lens}}{h_{section}}$, calculating $pop_{total}$ on a per-eye basis to account for covariance of the measurements, and subsequently averaging across n=6 eyes per sampled timepoint (3, 5, 8, 12, 17, 23, 30, 60, 90, 120, 180, and 360 dpf). These data are presented in Fig 7.

## Modelling lens growth

To simulate the toroidal CMZ population, we wanted to be able to simulate the stem cells at the periphery of the retina occasionally undergoing symmetric divisions, so as to keep up the same density of stem cells in the first ring around the lens. We did this by normalising our lens diameter measurements to the 72hpf value, to produce a measure of relative increase in lens size over the first year of CMZ activity. We performed a linear regression of the logarithm of this relative increase vs the logarithm of time using the Python statsmodels library, using the constants derived from this regression for a power-law model of lens growth. This model was used to supply the `WanStemCellCycleModel` with a target population value as described below.

## Estimation of 3dpf CMZ cell cycle length

To produce an estimate of the length of the cell cycle in 3dpf RPCs, we first took the total count of EdU-positive cells in dorsal and ventral CMZ from our cumulative labelling experiment described above, and divided by the total dorsal and ventral CMZ RPC population, as measured by the number of PCNA positive cells. This gave the labelled fraction measurements displayed in S4 Fig. Following Nowakowski et al. [21], we performed an ordinary least squares linear regression on these data using the Python statsmodels library. We then calculated the total cell cycle time $T_c = \frac{1}{slope}$ and $T_s = T_c \times intercept_y$. While this method is not suited for heterogenous populations of proliferating cells, its use is justified here, as the He SSM to which the estimate is being applied makes similar assumptions as Nowakowski et al., in particular the homogeneity of the population. It should not be considered more than a rough estimate.

## CHASTE Simulations

### Project code

All simulations were performed in an Ubuntu 16.04 environment using the CHASTE C++ simulation package version 2017.1 (git repository at `https://chaste.cs.ox.ac.uk/git/chaste.git`). The CHASTE package is a modular simulation suite for computational biology and has been described previously [10]. All of the code, simulator output, and empirical data used in this paper is available in the associated CHASTE project git repository at `http://github.com/mmattocks/ISP`, as well as in the supplementary archive S1 File. In brief, the approach taken was to encapsulate the SSMs examined here as separate concrete child classes inheriting from the CHASTE `AbstractSimpleCellCycleModel` class. An additional model, representing the simple immortal stem cell proposed in Wan

et al. [6], was similarly produced. Each model is generic and provided with public 652
functions to set their parameters and output relevant simulation outcomes (and 653
per-lineage debug data, if necessary). While these may be used in the ordinary 654
CHASTE unit test framework, permitting their use in any kind of cell-based simulation, 655
we wrote standalone simulator executables (apps, in CHASTE parlance) to permit 656
Python scripting of the simulator scenarios used in this paper. This allows for the 657
multi-threaded Monte Carlo simulations performed herein. Therefore, the Python 658
fixtures available in the project's `\python_fixtures\` directory were used to operate 659
the single-lineage simulators (`GomesSimulator`, `HeSimulator`, and `BoijeSimulator`) 660
and single-toroidal-CMZ simulator (`WanSimulator`) in `\apps\`, including the 661
CellCycleModels and related classes in `\src\`; the simulators output their data into 662
`\python_fixtures\testoutput\`. These data, along with the empirical observations 663
(both from the Harris groups' papers and our own described herein) available in 664
`\empirical_data\`, were processed by the analytical Python scripts in 665
`\python_fixtures\figure_plots\` to produce all of the figures used in this paper. A 666
guide to obtaining and using the simulators is available in S1 Appendix. 667

We have attempted to make the project fully reproducible and transparent. 668
CHASTE uses the C++ boost implementation of the Mersenne Twister random number 669
generator (RNG) to provide cross-platform reproducibility of simulations. All of our 670
simulators make use of this feature, permitting user-specified ranges of seeds for the 671
RNG. We have also used the numpy libary's RandomState Mersenne Twister container 672
to provide reproducible results for the Python fixtures, notably the SPSA optimisation 673
fixture. The output of the project code will, therefore, be identical every time it is run, 674
on any platform. 675

It should be noted that none of the code Harris' group used in their reports has been 676
published. We have made every effort to replicate the functional logic of the models as 677
closely as possible. Nonetheless, it is not possible to determine precisely why, for 678
instance, the original He SSM parameterisation failed so notably in our implementation. 679

## SPSA optimisation of models 680

In order to make a fair comparison between the He SSM and our deterministic mitotic 681
mode alternative model, we coded an implementation of the SPSA algorithm described 682
by Spall [17]. This is available in `\python_fixtures\SPSA_fixture.py`. This 683
algorithm, properly tuned, will converge almost-surely on a Karush-Kuhn-Tucker 684
optimum in the parameter space, given sufficient iterates, even with noisy output (eg. 685
due to RNG noise from low numbers of Monte Carlo samples, permitting conservation 686
of computational resources). It does so by approximating the loss function gradient 687
around a point in parameter space, selecting two sampling locations at each iterate, 688
then moving the next iterate's point in parameter space "downhill" along this gradient. 689

Our loss function was a modified AIC; we weighted the residual sum of squares from 690
the PD-type mitotic probabilities by 1.5 in order to improve convergence, as the PD 691
data are the most informative part of the dataset regarding the critical phase 692
boundaries (PP mitoses occur in all 3 phases of the He model, DD occur in phases 2 693
and 3, PD mitoses occur only in phase 2). We also compared model induction count 694
output (that is, panels A-C in Fig 4) up to count values of 1000 to penalise parameters 695
producing very large lineage totals. 696

The values of the SPSA gain sequence constants were selected to be as large as 697
possible without resulting in instability and failure to converge; this ensured that the 698
algorithm stepped over the widest possible range of the parameter space in finding 699
optima. The $\alpha$ and $\gamma$ coefficients were initially set at the noise-tolerant suggested values 700
of .602 and .101 respectively [17]. 200 iterations were performed. Initially, each iterate 701
involved 250 seeds of each induction timepoint for the count data (panels A-C in Fig 4) 702

and 100 seeds of the entire 23-72 hr simulation time for mitotic mode rate data (panels D-F). At iterate 170, these were increased to 1000 and 250 seeds, respectively, decreasing RNG noise to a low level. For the last 10 iterations, RNG noise was reduced to close to nil by increasing the seed numbers to 5000 and 1250, respectively; $\alpha$ and $\gamma$ were set to the asymptotically optimal 1 and $\frac{1}{6}$ for these iterates, as well. The particular seeds used were kept constant throughout in order to take advantage of the improved convergence rate this affords [28].

Because the simple SPSA can assign any value to parameters, our implementation uses constraints, projecting nonsensical values for parameters (e.g. negative values, mitotic mode probabilities summing to $> 1$) back into legal parameter space. The use of constraints has no effect on the algorithm's ability to almost-surely converge within the given parameter space [29] As we felt that the deterministic mitotic mode models' "sister shift" should be relatively small in order to be biologically plausible, we also constrained this parameter such that 95% of sister shifts would be less than the mean length of the shortest phase.

The numerical results of the SPSA algorithm are available in `\python_fixtures\testoutput\SPSA\HeSPSAOutput`, alongside png images of output from each iterate, enabling the visualisation of the progress of the algorithm.

### Monte Carlo simulations

Subsequent to SPSA optimisation, the parameters derived from this procedure were used to perform Monte Carlo simulations of large numbers of individual lineages, using ranges of 10000 non-overlapping seeds for each of panels A, B, C, G, H, and I, and 1000 for each of panels D, E, and F of Fig 4 (using the SPSA-optimised parameter sets) and S1 Fig (using the original He et al. parameterisation). This was performed using `\python_fixtures\He_output_fixture.py`.

Similarly, 100 seeds were used in Monte Carlo fashion to simulate individual toroidal CMZ populations using `\python_fixtures\Wan_output_fixture.py`, which scripts the WanSimulator. Each of these simulations is initialised with an RPC population drawn from a normal distribution given the mean and standard deviation of the CMZ torus estimated from our 3dpf observations, as described above. Each RPC is given a TiL value randomly chosen from 0 (newly born from a stem cell) to 17hr (exiting the CMZ, "forced to differentiate", in the terms of Wan et al.). A further population of stem cells, sized at $\frac{1}{10}$ of the RPC population, is added using the `WanStemCellCycleModel`. This stem population divides asymmetrically except when it falls under its target population value determined by the model of lens growth described above; as long as this condition holds, the stem cells will divide symmetrically to keep up their numbers.

### Simulation data analysis

Not all of the numerical data used in He et al. and Wan et al. has been published. As such, we reconstructed the lineage data from He et al.'s Figure 5C, used by He et al. to obtain their mitotic mode probabilities; our reconstructed values are available in `\empirical_data\empirical_lineages\`. We also reconstructed the He et al. WT and Ath5 morpholino data (Fig 4 panels H, I) and Wan et al. lineage count probabilities (panel G) from the relevant figures. These values are entered into `\python_fixtures\figure_plots\He_output_plot.py`, where they are used in the AIC calculations described below.

In order to assess the performance of the models used in our simulations, we used Akiake's Information Criterion (AIC). This allows us to compare models with dissimilar numbers of parameters, as AIC trades off goodness-of-fit against parameterisation. The

values reported in Table 1 were produced by
`\python_fixtures\figure_plots\He_output_plot.py`.

The 95% CIs for the model output presented in this paper were estimated by repetitive sampling of the output data, using the same number of lineages observed empirically. 5000 such samples were performed. This provides a reasonable estimate of the confidence interval given the limited observational sample size.

# Supporting information

**S1 Fig.    He SSM original parameterisation model output** As Fig 4; model output uses the parameterisation given in He et al. [5]. Note significant deviation from observations in panel B, reflecting excess proliferative activity of modelled RPCs.

**S2 Fig.    SPSA-optimised He SSM model output** As Fig 4; stochastic model output displayed alone.

**S3 Fig.    SPSA-optimised deterministic mitotic mode model output** As Fig 4; deterministic model output displayed by itself.

**S4 Fig.    Cumulative EdU Labelling of 3dpf CMZ RPCs** Fraction of PCNA-labelled CMZ RPCs which bear EdU label over time, as determined by histochemical labelling of central coronal cryosections of 3dpf zebrafish retinas. Dots represent results from individual retinas. Line is the ordinary least squares fit $\pm$ 95% CI. Cell cycle length (Tc) and S-phase length (Ts) are estimated from this fit following the method of Nowakowski et al. [21]

**S1 File.    Code archive** Archive of all code and code output used in this paper.

**S1 Appendix.    Guide to replicating analyses** This file contains a guide to reproducing all of the figures and analyses presented in this paper, using the code supplied in the S1 File archive, also available on the

# References

1. Fagan MB. Collaborative Explanation and Biological Mechanisms. Studies in History and Philosophy of Science Part A. 2015;52:67–78. doi:10.1016/j.shpsa.2015.03.004.

2. Till JE, Mcculloch EA, Siminovitch L. A Stochastic Model of Stem Cell Proliferation, Based on the Growth of Spleen Colony-Forming Cells. Proceedings of the National Academy of Sciences of the United States of America. 1964;51:29–36.

3. Fagan MB. Philosophy of Stem Cell Biology: Knowledge in Flesh and Blood. New directions in the philosophy of science. Houndmills, Basingstoke, Hampshire: Palgrave Macmillan; 2013.

4. Jaynes ET, Bretthorst GL, EBSCO Publishing. Probability Theory: The Logic of Science. Cambridge: Cambridge University Press; 2003.

5. He J, Zhang G, Almeida AD, Cayouette M, Simons BD, Harris WA. How Variable Clones Build an Invariant Retina. Neuron. 2012;75(5):786–798. doi:10.1016/j.neuron.2012.06.033.

6. Wan Y, Almeida AD, Rulands S, Chalour N, Muresan L, Wu Y, et al. The Ciliary Marginal Zone of the Zebrafish Retina: Clonal and Time-Lapse Analysis of a Continuously Growing Tissue. Development. 2016;143(7):1099–1107. doi:10.1242/dev.133314.

7. Boije H, Rulands S, Dudczig S, Simons BD, Harris WA. The Independent Probabilistic Firing of Transcription Factors: A Paradigm for Clonal Variability in the Zebrafish Retina. Developmental Cell. 2015;34(5):532–543. doi:10.1016/j.devcel.2015.08.011.

8. Temple S, Raff MC. Clonal Analysis of Oligodendrocyte Development in Culture: Evidence for a Developmental Clock That Counts Cell Divisions. Cell. 1986;44(5):773–779. doi:10.1016/0092-8674(86)90843-3.

9. Burnham KP, Anderson DR. Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach. 2nd ed. New York: Springer; 2002.

10. Mirams GR, Arthurs CJ, Bernabeu MO, Bordas R, Cooper J, Corrias A, et al. Chaste: An Open Source C++ Library for Computational Physiology and Biology. PLoS Computational Biology. 2013;9(3):e1002970. doi:10.1371/journal.pcbi.1002970.

11. Gomes FLAF, Zhang G, Carbonell F, Correa JA, Harris WA, Simons BD, et al. Reconstruction of Rat Retinal Progenitor Cell Lineages in Vitro Reveals a Surprising Degree of Stochasticity in Cell Fate Decisions. Development (Cambridge, England). 2011;138(2):227–235. doi:10.1242/dev.059683.

12. Schaffner KF. Discovery and Explanation in Biology and Medicine. Science and its conceptual foundations. Chicago: University of Chicago Press; 1993.

13. Feyerabend P. Against Method. 3rd ed. London ; New York: Verso; 1993.

14. Holt CE, Bertsch TW, Ellis HM, Harris WA. Cellular Determination in the Xenopus Retina Is Independent of Lineage and Birth Date. Neuron. 1988;1(1):15–26.

15. Agathocleous M, Harris WA. From Progenitors to Differentiated Cells in the Vertebrate Retina. Annual Review of Cell and Developmental Biology. 2009;25:45–69. doi:10.1146/annurev.cellbio.042308.113259.

16. Cayouette M, Barres BA, Raff M. Importance of Intrinsic Mechanisms in Cell Fate Decisions in the Developing Rat Retina. Neuron. 2003;40(5):897–904.

17. Spall JC. Implementation of the Simultaneous Perturbation Algorithm for Stochastic Optimization. IEEE Transactions on Aerospace and Electronic Systems. 1998;34(3):817–823. doi:10.1109/7.705889.

18. Parichy DM, Elizondo MR, Mills MG, Gordon TN, Engeszer RE. Normal Table of Postembryonic Zebrafish Development: Staging by Externally Visible Anatomy of the Living Fish. Developmental Dynamics. 2009;238(12):2975–3015. doi:10.1002/dvdy.22113.

19. Tropepe V. Retinal Stem Cells in the Adult Mammalian Eye. Science. 2000;287(5460):2032–2036. doi:10.1126/science.287.5460.2032.

20. Harris WA, Perron M. Molecular Recapitulation: The Growth of the Vertebrate Retina. International Journal of Developmental Biology. 1998;42:299–304.

21. Nowakowski RS, Lewin SB, Miller MW. Bromodeoxyuridine Immunohistochemical Determination of the Lengths of the Cell Cycle and the DNA-Synthetic Phase for an Anatomically Defined Population. Journal of Neurocytology. 1989;18(3):311–318.

22. Centanin L, Ander JJ, Hoeckendorf B, Lust K, Kellner T, Kraemer I, et al. Exclusive Multipotency and Preferential Asymmetric Divisions in Post-Embryonic Neural Stem Cells of the Fish Retina. Development. 2014;141(18):3472–3482. doi:10.1242/dev.109892.

23. Engerer P, Suzuki SC, Yoshimatsu T, Chapouton P, Obeng N, Odermatt B, et al. Uncoupling of Neurogenesis and Differentiation during Retinal Development. The EMBO Journal. 2017;36(9):1134–1146. doi:10.15252/embj.201694230.

24. Raj A, van Oudenaarden A. Stochastic Gene Expression and Its Consequences. Cell. 2008;135(2):216–226. doi:10.1016/j.cell.2008.09.050.

25. Trimarchi JM, Stadler MB, Cepko CL. Individual Retinal Progenitor Cells Display Extensive Heterogeneity of Gene Expression. PLOS ONE. 2008;3(2):e1588. doi:10.1371/journal.pone.0001588.

26. Raser JM, O'Shea EK. Control of Stochasticity in Eukaryotic Gene Expression. Science; Washington. 2004;304(5678):1811–4.

27. Westerfield M. The Zebrafish Book : A Guide for the Laboratory Use of Zebrafish. http://zfinorg/zf_info/zfbook/zfbkhtml. 2000;.

28. Kleinman NL, Spall JC, Naiman DQ. Simulation-Based Optimization with Stochastic Approximation Using Common Random Numbers. Management Science. 1999;45(11):1570–1578.

29. Sadegh P. Constrained Optimization via Stochastic Approximation with a Simultaneous Perturbation Gradient Approximation. Automatica. 1997;33(5):889–892. doi:10.1016/S0005-1098(96)00230-0.