

Submitted to the *Annals of Applied Statistics*
arXiv: [arXiv:0000.0000](https://arxiv.org/abs/2000.0000)

1 **TRACKING HEMATOPOIETIC STEM CELL EVOLUTION**
2 **IN A WISKOTT-ALDRICH CLINICAL TRIAL**

3 BY DANILO PELLIN^{*}, LUCA BIASCO[†], SERENA SCALA[‡], CLELIA DI
4 SERIO[‡] AND ERNST C. WIT[§]

5 *Harvard Medical School^{*}, UCL GOS Institute of Child Health[†], Vita-Salute*
6 *San Raffaele University[‡] and Università della Svizzera italiana[§]*

7 Hematopoietic Stem Cells (HSC) are the cells that give rise to
8 all other blood cells and, as such, they are crucial in the healthy
9 development of individuals. Wiskott-Aldrich Syndrome (WAS) is a
10 severe disorder affecting the regulation of hematopoietic cells and is
11 caused by mutations in the WASP gene. We consider data from a
12 revolutionary gene therapy clinical trial, where HSC harvested from
13 3 WAS patients' bone marrow have been edited and corrected using
14 viral vectors. Upon re-infusion into the patient, the HSC multiply
15 and differentiate into other cell types. The aim is to unravel the cell
16 multiplication and cell differentiation process, which has until now
17 remained elusive.

18 This paper models the replenishment of blood lineages resulting
19 from corrected HSC via a multivariate, density-dependent Markov
20 process and develops an inferential procedure to estimate the dy-
21 namic parameters given a set of temporally sparsely observed tra-
22 jectories. Starting from the master equation, we derive a system of
23 non-linear differential equations for the evolution of the first- and
24 second-order moments over time. We use these moment equations in
25 a generalized method-of-moments framework to perform inference.
26 The performance of our proposal has been evaluated by consider-
27 ing different sampling scenarios and measurement errors of various
28 strengths using a simulation study. We also compared it to another
29 state-of-the-art approach and found that our method is statistically
30 more efficient.

31 By applying our method to the Wiskott-Aldrich Syndrome gene
32 therapy data we found strong evidence for a myeloid-based develop-
33 mental pathway of hematopoietic cells where fates of lymphoid and
34 myeloid cells remain coupled even after the loss of erythroid poten-
35 tial.

36 All code used in this manuscript can be found in the online Sup-
37 plement, and the latest version of the code is available at github.com/dp3111n/SLCDP_v1.0.
38

39 **1. Introduction.** Although mammalian organisms have more than a
40 hundred different cell types, many tissues are sustained by relatively few va-
41 rieties of multipotent stem and progenitor cells ([Weissman, 2000](#); [Blanpain,](#)

Keywords and phrases: gene therapy, clonal tracking, Wiskott-Aldrich Syndrome, Mul-
tivariate Markov process, Master equation,¹generalized method-of-moments, non-linear
differential equations

42 [Horsley and Fuchs, 2007](#); [Snippert and Clevers, 2011](#)). Given their impor-
43 tance, a comprehensive understanding of stem cells is crucial for advancing
44 the development of regenerative medicine. HSC represent a particular pool
45 of cells that resides mainly in the bone marrow and has the unique capa-
46 bility of self-renewal. Through a process of progressive specialization called
47 hematopoiesis, HSC can give rise and replenish all blood lineages in a human
48 being, lifelong. HSC are among the most clinically relevant cell population
49 and are used to treat many hematological malignancies and bone marrow
50 disorders. Despite being the focus of decades of research and clinical efforts,
51 many questions about HSC biology are still open and debated. For example,
52 it is well-established that a progressive loss of multi-lineage potential occurs
53 when descending the hematopoietic cell differentiation hierarchy from HSC
54 to committed cell types and then, finally, mature blood cells. However, it
55 is still unclear at what stage of the differentiation process the separation
56 between the three main cell lineage groups, lymphoid, myeloid, and ery-
57 throid, happens. Other essential aspects about the metabolism of human
58 blood cells, such as how duplication, death, and differentiation rates are or-
59 chestrated along the blood phylogeny to maintain the hematopoietic system
60 stable, are still unknown.

61 Gene therapy consists of delivering DNA or RNA fragments into cells of
62 patients as a drug to treat a disease. It has been mainly applied to inher-
63 ited monogenetic disease where deleterious mutations occurring in a specific
64 known gene lead to the synthesis of a dysfunctional protein causing the
65 symptoms. Under this setting, gene therapy offers a real opportunity and
66 can be used to provide cells with a correct copy of the gene, thereby produc-
67 ing a functional version of the protein. The treatment effect is tied to the
68 presence and activity of the therapeutic gene in specific cells or tissue, hence
69 for the long-term treatment of hematological disorders, HSC represent the
70 ideal target for gene therapy clinical trials ([Naldini, 2011](#); [Biffi et al., 2013](#);
71 [Aiuti et al., 2013](#)).

72 This paper will focus on a gene therapy clinical trial for Wiskott-Aldrich
73 Syndrome (WAS), an inherited immunodeficiency caused by mutations in
74 the gene encoding for WAS protein. The study was performed by the au-
75 thors of this paper and described in clinical detail in [Biasco et al. \(2016\)](#).
76 Briefly, HSC sorted from patient's bone marrow samples according to their
77 immunophenotyping characteristic — enrichment analysis for known pro-
78 tein on a cell's cellular membrane, such as CD34 molecules specifically for
79 HSC isolation — are distinctly labeled through the random incorporation of
80 the WASP gene into their genome, using a lentiviral vector. Importantly, all
81 progeny deriving from a marked HSC, through both duplication and differ-

82 entiation, will carry the corrected copy of the gene and the identical unique
83 markings defined by the original viral insertion site (IS). This procedure al-
84 lows not only to obtain a long-term and widespread expression of the WAS
85 protein among all blood lineages but also to perform in-vivo clonal tracking,
86 the longitudinal observation of multiple clones' evolution. It is crucial to
87 highlight that for ethical reasons, gene therapy is one of the few settings in
88 which scientists can collect information about human, *in-vivo* hematopoiesis
89 at clone level.

90 One of the first quantitative analysis of clonal tracking data was devel-
91 oped in the context of a non-human primate rhesus macaque study by [Wu](#)
92 [et al. \(2014\)](#). Using clustering methods on the multi-lineage clonal output of
93 barcoded HSC, authors demonstrated how the correlation among lineages
94 changes during reconstitution, with uni-lineage short-term progenitors being
95 supplanted over time by multi-lineage long-term clones. ([Biasco et al., 2016](#);
96 [Pellin et al., 2019](#)) model clones dynamics using local linear approximations.
97 Assuming linearity offers several advantages from a computational perspec-
98 tive, but also implies that cell type counts must eventually either go to zero
99 or infinity in the long term. This assumption is biologically unrealistic be-
100 cause the hematopoietic system evolves in a constrained environment with
101 limited resources and space available. At the same time, the replenishment
102 of blood cells lasts for the entire life span of a human being. To extrap-
103 olate insight from real data [Biasco et al. \(2016\)](#) and [Pellin et al. \(2019\)](#)
104 relied on a first-order local linear approximation of the dynamics: this is
105 efficient but not very accurate when the time between consecutive process
106 measurements is large, as it is in the case of gene therapy clinical trials. [Xu](#)
107 [et al. \(2019\)](#) re-analyzed the rhesus macaque data using a statistical frame-
108 work that models hematopoiesis as a multi-type Markov branching process,
109 similar to our set-up. In [Xu et al. \(2019\)](#), clone trajectories are considered
110 realizations from a stochastic process defined using a set of fundamental
111 cellular events with event-specific rates. The authors showed that it is pos-
112 sible to derive exact analytical formulation for the evolution of the moments
113 through a set of ordinary differential equations (ODEs), given the cell differ-
114 entiation tree configuration and assuming event rates to be linear in the cell
115 counts. The estimation of the cell differentiation dynamic is performed by
116 matching model-based correlation functions to empirical lineage temporal
117 correlations. An alternative approach, similar to ours, could a be Bayesian
118 implementation ([Wilkinson, 2006](#); [Golightly and Wilkinson, 2008](#)), which
119 can deal with temporal sampling an observational noise in a natural fash-
120 ion. We expect that implementation of those methods would yield similar
121 results to ours.

122 In section 2 we describe the clonal tracking data obtained from a gene
123 therapy clinical trial for Wiskott-Aldrich syndrome, for which the statistical
124 methodology in this paper has been developed. The stochastic cell differen-
125 tiation process and its characteristics are presented in section 3. In section 4
126 a non-linear generalized least squares estimation procedure for the param-
127 eters in the stochastic process is developed, both from a methodological and
128 computational point of view. section 5 is dedicated to simulation studies.
129 In section 5.1 the performance of our proposal is compared for different
130 sampling time intervals with a simpler polynomial generalized least squares
131 estimation procedure. Section 5.2 and section 5.3 are focused respectively
132 on the impact on inference performance of having (multiplicative) measure-
133 ment errors on cell count observations and the effect of potential model
134 misspecification. Section 5.4 compares our method to the correlation-based
135 moment estimator by Xu et al. (2019). In section 6 we return to the WAS
136 gene therapy clinical trial data and answer the main questions of this pa-
137 per, namely, estimate the coefficients driving HSC differentiation and verify
138 whether the WAS data support the classical dichotomy model or a recently
139 proposed myeloid-based model of hematopoietic stem cell differentiation.

140 **2. Hematopoietic stem cell gene therapy in Wiskott-Aldrich**
141 **Syndrome patients.** WAS syndrome is an X-linked primary immunod-
142 efficiency characterized by infections, micro-thrombocytopenia, eczema, au-
143 toimmunity, and lymphoid malignancies. The disorder is caused by muta-
144 tions in the WAS gene, which encodes for WASP, a protein that regulates
145 cytoskeleton conformation and is involved in proliferation, migration, and
146 immunological synapsis formation. For patients without a matched donor,
147 gene therapy based on the infusion of autologous gene-corrected HSC rep-
148 resent an alternative therapeutic strategy.

149 Three children with WAS, who did not have compatible allogeneic donors,
150 were enrolled in phase I/II clinical trial. Autologous BM-derived CD34+ cells
151 were collected, transduced with a lentiviral vector coding for human WASP
152 under the control of a 1.6-kb reconstituted WAS gene promoter (LV-w1.6W)
153 using an optimized protocol, and re-infused intravenously into the patients
154 three days after collection. Patients are given chemotherapy treatment before
155 receiving the engineered cell infusion to deplete the existing HSC compart-
156 ment and to facilitate the engraftment of corrected cells. This conditioning
157 procedure requires a fast replenishment of all blood lineages by corrected
158 HSC upon infusion until a homeostasis condition is met. All three WAS pa-
159 tients showed robust and multi-lineage engraftment of gene-corrected cells
160 in BM and PB up to the latest follow-up.

161 We collected IS from eight distinct Peripheral Blood (PB) and seven dis-
162 tinct Bone Marrow (BM) lineages at multiple time-points up to 36 months
163 after infusion of transduced HSCs using a combination of linear-amplification-
164 mediated (LAM)-PCR and next-generation sequencing (NGS) technologies
165 (Biasco et al., 2011).

166 After initial clonal fluctuations, we observed stable and polyclonal recon-
167 stitution in all hematopoietic lineages starting from 1 year after the infu-
168 sion of gene-corrected HSC. Importantly, no adverse event associated with
169 insertional mutagenesis was detected, allowing us to exploit IS to assess
170 hierarchical relationships among engineered blood cell types in humans.

171 A major distinction in three subgroups, named lymphoid, myeloid and
172 erythroid branches, can be made within the hematopoietic cell types. The
173 lymphoid branch, responsible for the adaptive immune system, can, in turn,
174 be subdivided into T-cells (CD3 in BM and CD4, CD8, CD3 in PB), B-
175 cell (CD19), and Natural Killer cells (NK-cells, CD56). Myeloid cell types
176 are involved in such diverse roles as innate immunity, adaptive immunity,
177 and blood clotting and are composed of monocytes (CD14), granulocytes
178 (CD15), and megakaryocytes (CD61). Erythrocytes are the oxygen-carrying
179 red blood cells (GLYCO).

180 Two different models of hematopoiesis are currently debated, shown in
181 Figure 1. The classical dichotomy model assumes that HSC first generate
182 a common myeloid-erythroid progenitor (CMEP) and a common lymphoid
183 progenitor (CLP). The CLP then produces only T-cells or B-cells. The al-
184 ternative myeloid-based model postulates that HSC first diverge into the
185 CMEP and a common myeloid-lymphoid progenitor (CMLP), which gener-
186 ates T- and B-cell progenitors through a bipotential myeloid-T progenitor
187 and a myeloid-B progenitor stage. The main difference is that according to
188 the second, all erythroid, T- and B-lineage branches retain the potential to
189 generate myeloid cells, even after the segregation of T- and B-cell lineages
190 (Kawamoto, Wada and Katsura, 2010).

191 This study aims to provide novel insights about human hematopoiesis
192 and the HSC differentiation process in-vivo. This crucial biological question
193 remained unresolved despite extensive efforts over the past years. Exploiting
194 clonal tracking data from WAS gene therapy clinical trial, in section 6 we
195 will investigate the hierarchical relationship among cell types and estimate
196 lineage-specific cell duplication and death rates.

197 **3. Stochastic logistic cell differentiation process.** We consider an
198 N -dimensional, continuous time counting process $\mathbf{X}_t = (X_{t1}, \dots, X_{tN})$,
199 where $t \in \mathbb{R}$ and $\mathbf{X}_t \in \mathbb{N}_0^N$. Each element of X_{ti} , corresponds to the number

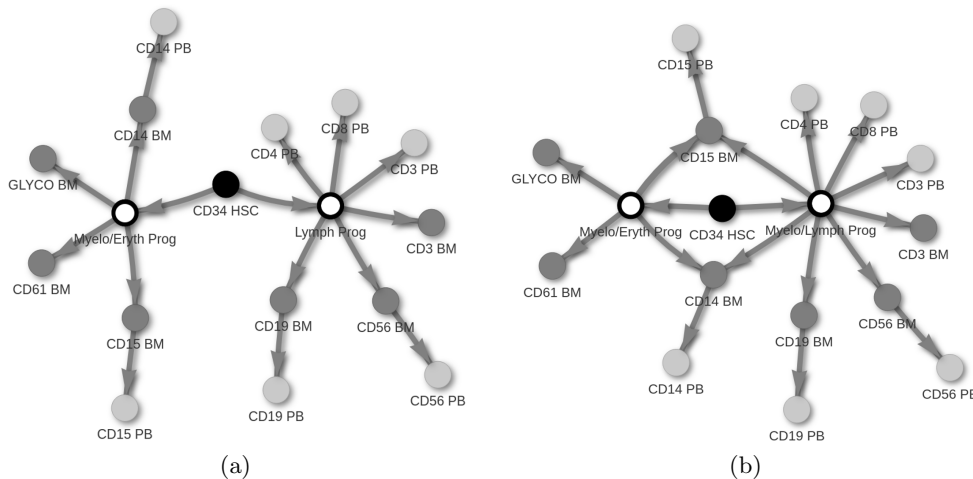


FIG 1. **Two competing hematopoiesis theories.** Filled nodes correspond to lineages analyzed in this manuscript. Black, dark gray and light gray nodes represent Hematopoietic Stem Cell, Bone Marrow, and Peripheral Blood lineages. Empty nodes are latent, unobserved cell types. The classical dichotomy model (a) assumes that HSC first generate a common myeloid-erythroid progenitor (CMEP) and a common lymphoid progenitor (CLP), whereas the alternative myeloid-based model (b) postulates that HSC first diverge into the CMEP and a common myelo-lymphoid progenitor (CMLP).

200 of cells of type C_i , ($i = 1, \dots, N$) present in the system at time t . X_1 refers
 201 to the HSC count, the most primitive and multi-potent cell type.

We assume that \mathbf{X} evolves according to a continuous-time Markov process. There are three event types in the process: cell duplication, cell death, and, importantly, cell differentiation. Individual cells are assumed to evolve independently from each other and cells belonging to the same cell type are assumed to obey the same laws. Event rates are assumed constant over time. The generic cell duplication rate $\alpha_i \geq 0$ is assumed to be a linear growth term, corresponding to the expected number of cell duplications per time unit per cell of type C_i , $i = 1, \dots, N$,

$$P(X_{t+\partial t, i} = x_{t, i} + 1, X_{t+\partial t, -i} = x_{t, -i} | \mathbf{X}_t = \mathbf{x}_t) \approx x_i \alpha_i \partial t.$$

Secondly, linear cell duplication is eventually overcome by quadratic cell death. This assumption results in a cell type specific logistic growth curve, represented by the following conditional transition probabilities for cell death of type C_i (for some $\delta_i \geq 0$, $i = 1, \dots, N$),

$$P(X_{t+\partial t, i} = x_{t, i} - 1, X_{t+\partial t, -i} = x_{t, -i} | \mathbf{X}_t = \mathbf{x}_t) \approx x_i^2 \delta_i \partial t.$$

Furthermore, it is assumed that cell differentiation from cell type i into cell type j is a process with constant rate $\lambda_{ij} \geq 0, i, j = 1, \dots, N, i \neq j$,

$$P(X_{t+\partial t, i} = x_{ti} - 1, X_{t+\partial t, j} = x_{tj} + 1, X_{t+\partial t, -ij} = x_{t, -ij} | \mathbf{X}_t = \mathbf{x}_t) \approx x_i \lambda_{ij} \partial t.$$

202 It is convenient to write the Markov process in a vectorized form. Each
203 cellular event $k \in \{1, \dots, r\}$ can be associated with an N -dimensional integer
204 vector \mathbf{v}_k , describing the net change in the state induced by event k . Given
205 the hazard $h_k(\mathbf{x}, \boldsymbol{\theta}) = (x_i \alpha_i, x_i^2 \delta_i, x_i \lambda_{ij})$ for $\boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\delta}, \boldsymbol{\lambda})$, we can write the
206 process generally as $P(X_{t+\partial t} = x_t + \mathbf{v}_k | X_t = x_t) \approx h_k(x_t; \boldsymbol{\theta}) \partial t$. The whole
207 process can be recast in matrix notation involving the net effect matrix, \mathbf{V} ,
208 corresponding to an $N \times r$ integer matrix, in which the columns correspond
209 to the vectors \mathbf{v}_k ($k = 1, \dots, r$). For simplicity we assume that the first
210 N columns of \mathbf{V} refer to cell duplications, the second N to cell deaths
211 and the remaining columns to differentiation events. The hazard $\mathbf{h}(\mathbf{X}, \boldsymbol{\theta}) =$
212 (h_1, \dots, h_r) , is the r -dimensional vector of the r individual event hazards.

213 We are here considering that cells can only divide symmetrically, gener-
214 ating two daughters cells of the same nature as the mother cell. Assuming
215 the alternative asymmetric division, such as in [Xu et al. \(2019\)](#), means that
216 division is always coupled with a differentiation event, resulting in the for-
217 mation of two cells with different properties and fate. Even though recent
218 literature based on in-vitro experiments supports the possibility for HSC to
219 undergo asymmetric division, little is known about the frequency of such
220 events in-vivo and whether other lineages also have this capability.

221 Logistic differential equation models are widely used in the study of
222 hematopoietic dynamics. Yet, it has not been applied in the context of clonal
223 tracking data. According to the transition probabilities specified in our cell
224 differentiation process, a clone will generate new cells purely based on its
225 current counts. When a given size is reached, scarcity of nutrients and space
226 in the niche makes cells die at a faster rate, preventing clone size from grow-
227 ing exponentially. Biologically, this is likely to be a too simplistic model of
228 steady-state maintenance. In-vivo, cell duplication and death are regulated
229 based on the current system needs using complex signaling mechanisms.
230 However, our assumption has the remarkable advantage of allowing infer-
231 ence on all parameters of the differentiation process, avoiding the necessity
232 to resort to literature data to set some coefficients, as proposed in [Xu et al.](#)
233 [\(2019\)](#), or to infer *net rates* (duplication minus death rates) as done in [Pellin](#)
234 [et al. \(2019\)](#).

3.1. *Moment equations.* For any stochastic process obeying the Markov property, given some initial condition \mathbf{X}_0 , it is possible to determine the

evolution of the probability distribution function associated with the system states over time, $P(\mathbf{X}; t)$, using the Chemical Master Equation (CME) (Bailey, 1964; Kampen, 1981; Risken, 1984; Gardiner, 1985). The CME is defined as a differential equation for the process transition probabilities and can be written as

$$(3.1) \quad \frac{dP(\mathbf{x}; t)}{dt} = \sum_{k=1}^r [h_k(\mathbf{x} - \mathbf{v}_{k,\cdot}; \boldsymbol{\theta})P(\mathbf{x} - \mathbf{v}_{k,\cdot}; t) - h_k(\mathbf{x}; \boldsymbol{\theta})P(\mathbf{x}; t)]$$

A solution and complete characterization of $P(\mathbf{x}; t)$ from (3.1) is unfeasible due to the large set of possible states configurations. However, important insights about cell differentiation dynamics and its parameters can be determined based on the time evolution of a few low-order statistical moments. Let $m_i(t)$ describe the time evolution of $E[\mathbf{X}_{it}] = \sum_{\mathbf{x}} \mathbf{x} P_{\mathbf{X}_i}(\mathbf{x}; t)$. Applying the derivate operator to both sides, we obtain the dynamics of the mean of $\mathbf{X}(t)$ can be summarized in the following ODE system,

$$(3.2) \quad \frac{dm_i(t)}{dt} = \sum_{k=1}^r v_{k,i} E[h_k(\mathbf{X}_t; \boldsymbol{\theta})]; \quad i = 1, \dots, N.$$

Similarly, let $m_{i,j}^2(t)$ be the time evolution for the symmetric second-order moments $E[X_{ti}X_{tj}]$ as

$$(3.3) \quad \frac{dm_{i,j}^2(t)}{dt} = \sum_{k=1}^r v_{k,j} E[X_{ti}h_k(\mathbf{X}_t; \boldsymbol{\theta})] + \sum_{k=1}^r v_{k,i} E[X_{tj}h_k(\mathbf{X}_t; \boldsymbol{\theta})] + \sum_{k=1}^r v_{k,i}v_{k,j} E[h_k(\mathbf{X}_t; \boldsymbol{\theta})].$$

235 A detailed derivation of (3.2) and (3.3) can be found in [Supplement A](#).
 236 With death rates $x_i^2\delta_i$ being polynomial of degree 2, the time evolution for
 237 the generic moment of order n depends on moments of order $n + 1$, leading
 238 to an infinite system of equations that can not be solved directly. There
 239 are different approaches to address this issue that consists of approximation
 240 methods. The most popular are the Chemical Langevin Equation, a diffu-
 241 sion approximation of the CME (Wilkinson, 2006; Golightly and Wilkinson,
 242 2005), the system size expansion (Kampen, 1981; Elf and Ehrenberg, 2003),
 243 the Linear Noise Approximation (Gardiner, 1985), and the moments closure
 244 approximation (Grima, 2012). Hematopoietic differentiation is a stochastic
 245 process with an output consisting of a relatively small amount of cells, that
 246 starts from an individual HSC. These are not ideal conditions to apply the
 247 CLE approximation (Schnoerr, Sanguinetti and Grima, 2017). In its funda-
 248 mental formulation, LNA requires the assumption that fluctuation and, as a

249 consequence the clone cell counts, have a multivariate Normal distribution.
250 This assumption, combined with the deterministic first-moment dynamics,
251 poses challenges for approximating systems with multimodal steady-state
252 behavior, as it is the cell differentiation process. We therefore approximate
253 the moments evolution using moment closure.

254 Several moment closure approaches have been proposed in the litera-
255 ture: (i) assuming a specific probability distribution for $P(\mathbf{X}; t)$ (Whittle,
256 1957; Nåsell, 2003a,b; Keeling, 2000) or (ii) a separable-derivative-matching
257 schema proposed in Singh and Hespánha (2007). The choice of the most
258 appropriate method depends on the application of interest and the nature
259 of the data analyzed. In this manuscript, we follow the indication provided
260 in Schnoerr, Sanguinetti and Grima (2017), where these methods have been
261 thoroughly tested and compared. Based on numerical evaluations, authors
262 conclude that moment closure based on a normal distribution assumption is
263 in general favorable for stability and precision. However, it is important to
264 notice that the approach presented here is in principle valid irrespectively
265 of the moment closure strategy adopted.

A Gaussian third-order moment approximation consists of setting the skewness equal to 0, leading to third-order moment definitions as follows,

$$\begin{aligned} \mathbb{E}[X_{ti}^3] &= 3 \mathbb{E}[X_{ti}] \mathbb{E}[X_{ti}^2] - 2 \mathbb{E}[X_{ti}]^3 \\ (3.4) \quad \mathbb{E}[X_{ti} X_{tj}^2] &= 2 \mathbb{E}[X_{tj}] \mathbb{E}[X_{ti} X_{tj}] + \mathbb{E}[X_{ti}] \mathbb{E}[X_{tj}^2] - 2 \mathbb{E}[X_{ti}] \mathbb{E}[X_{tj}]^2 \end{aligned}$$

266 Substituting these third-order moments in (3.3) with the appropriate non-
267 linear formulation in (3.4), we derive two coupled systems of ordinary differ-
268 ential equations for the first and second order moments for the stochastic cell
269 differentiation process. Based on this ODE system we will now propose an
270 inferential procedure able to obtain parameter estimates and to reconstruct
271 the cell differentiation structure.

272 **4. Inference.** The cell differentiation process is typically observed across
273 a discrete number of time points and some replicates. To simplify nota-
274 tion, we assume we have S equally Δt -spaced observations $\mathbf{X}_s, s = 1, \dots, S$
275 from one realization of an N -dimensional stochastic cell differentiation pro-
276 cess. It is computationally trivial to drop the equal spacing assumption. A
277 likelihood-based approach would need to integrate all possible states and in-
278 termediate time-points, effectively making closed-form inference impossible.
279 Instead, we will derive a methods-of-moments type estimator for inferring
280 the parameters of interest.

281 As mentioned in section 2, in an experimental setting, clone sizes are
282 estimated using NGS readouts. Despite several protocols, techniques and

estimators proposed in the literature (Berry et al., 2012; Calabria et al., 2014; Leonardelli et al., 2016), measurement error still plays an important role in the quantitative characterization of the progeny of an individual HSC. Therefore, we included in our model definition (4.1) a multiplicative noise term that can be adjusted using an intensity parameter to be set according to the protocol followed.

4.1. *Non-linear generalized method of moments.* We reformulate the process as a non-linear regression problem, i.e.,

$$(4.1) \quad \mathbf{X}_s = f(\mathbf{x}_{s-1}; \boldsymbol{\theta}) + \boldsymbol{\varepsilon}_s$$

where $f(\mathbf{x}_{s-1}, \boldsymbol{\theta}) = E[\mathbf{X}_s | \mathbf{x}_{s-1}; \boldsymbol{\theta}]$ is a known non-linear function of the process state at time step $s - 1$ and $\boldsymbol{\varepsilon}_s$ is an N -dimensional mismatch variable with $E[\boldsymbol{\varepsilon}_s] = \mathbf{0}_N$, $\text{Var}(\boldsymbol{\varepsilon}_s) = \mathbf{W}_s + \varphi \mathbf{N}_s$. $\mathbf{W}_s = \text{Cov}[X_i(s), X_j(s) | \mathbf{x}_{s-1}; \boldsymbol{\theta}]$ is a $N \times N$ matrix for some known non-linear function g modeling the stochastic process intrinsic covariance structure. The diagonal matrix $\varphi \mathbf{N}_s$ describes a multiplicative-like noise term that allows to include a measurement uncertainty on cell counts recordings. In particular, φ is a user-defined dispersion parameter that can be set by using a control experiment, as described in section 6, and $\mathbf{N}_s = \text{Diag}(\mathbf{x}_{s-1})$ is a $N \times N$ diagonal matrix with the cell counts on the diagonal. To avoid the usage of unnecessarily complicated notation in the description of our inference framework, throughout this section we will consider observations to be noise-free ($\varphi = 0$). However, the implemented method on the data does consider the dispersion parameter ($\varphi > 0$).

For each value of s the function $f(\mathbf{x}_{s-1}, \boldsymbol{\theta}) = \mathbf{m}(s)$ and matrix $\mathbf{W}_s = \mathbf{m}^2(s) - \mathbf{m}(s)\mathbf{m}(s)^t$ are defined through the solutions of the coupled ODE system (3.2) and (3.3) setting \mathbf{x}_{s-1} as initial conditions for $\mathbf{m}(s - 1)$ and $x_{s-1,i}x_{s-1,j}$ for $m_{ij}^2(s - 1)$. This projects the state and covariance matrix from one observed time-point to the next.

Applying this procedure to all observations available, we can perform parameter estimation by means of a generalized method of moments estimator with objective function,

$$(4.2) \quad \hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \geq \mathbf{0}_r} [\mathbf{x}_{1:S} - f(\mathbf{x}_{0:S-1}; \boldsymbol{\theta})]^T (\mathbf{W}_{1:S})^{-1} [\mathbf{x}_{1:S} - f(\mathbf{x}_{0:S-1}; \boldsymbol{\theta})]$$

where $\mathbf{x}_{1:S}$ and $f(\mathbf{x}_{0:S-1}; \boldsymbol{\theta})$ are $(N \times S)$ -dimensional column vectors and $\mathbf{W}_{1:S}$ is a $NS \times NS$ block diagonal matrix, in which blocks correspond to expected variance-covariance matrices \mathbf{W}_s within each time increment. In Supplement B all the elements introduced in this section are derived for a simple example involving 3 cell types.

312 To calculate the solution $\hat{\theta}$, we propose an iterative procedure in which
 313 moments estimation and parameter refinement alternate until a convergence
 314 criterion is met. The complete algorithm is described in section 4.2. It is
 315 worth noting that for the solution of (4.2), some initial values $\hat{\theta}^{(0)}$ for θ , must
 316 be provided as input in order to start the iterative optimization procedure.
 317 Given the amount the parameters involved in the model, especially if no
 318 or limited assumptions are made to limit possible cell differentiations (by
 319 setting $\lambda_{ij} = 0$ for some i, j), it is important to start the minimization of
 320 (4.2) from accurate starting values within the convex region surrounding
 321 the true, unknown θ . Supplement E presents a local linear approximation
 322 approach that can be used to obtain a sensible starting value (Pellin et al.,
 323 2019).

4.2. *Algorithm.* To find the solution to the minimization problem in (4.2), a modified implementation of the Gauss-Newton algorithm is proposed (Björck, 1996). Its pseudo-code is available in Algorithm 1. The procedure receives as input the initial cell counts, observations during the follow-up time, $\mathbf{x}_{0:S}$, and the system of ODEs for the first order, $\mathbf{m}(t)$, and second order, $\mathbf{m}^2(t)$. The algorithm starts with the initial estimate $\hat{\theta}^{(0)}$ that is then refined using an iterative procedure with the updating formula

$$\hat{\theta}^{(k+1)} = \hat{\theta}^{(k)} + \widehat{\Delta\theta}^{(k)},$$

324 where $\widehat{\Delta\theta}^{(k)}$ is the solution to the following constrained quadratic problem,

$$\begin{aligned} \widehat{\Delta\theta}^{(k)} = \arg \min_{\Delta\theta} [r(\hat{\theta}^{(k)}) - J(\hat{\theta}^{(k)})\Delta\theta]^\top [W(\hat{\theta}^{(k)})]^{-1} [r(\hat{\theta}^{(k)}) - J(\hat{\theta}^{(k)})\Delta\theta] \\ (4.3) \quad \text{such that } \Delta\theta \geq -\hat{\theta}^{(k)} \end{aligned}$$

in which $r(\hat{\theta}^{(k)}) = \mathbf{x}_{1:S} - f(\mathbf{x}_{0:S-1}; \theta^{(k)})$ is the residual NS -dimensional column vector and

$$J(\hat{\theta}^{(k)}) = \begin{bmatrix} \frac{df(\mathbf{x}_0; \hat{\theta}^{(k)})}{d\theta} & \frac{df(\mathbf{x}_1; \hat{\theta}^{(k)})}{d\theta} & \dots & \frac{df(\mathbf{x}_{S-1}; \hat{\theta}^{(k)})}{d\theta} \end{bmatrix}^\top$$

is the $NS \times r$ Jacobian matrix. Each $\frac{df(\mathbf{x}_s; \hat{\theta}^{(k)})}{d\theta}$ is a $N \times r$ matrix measuring the change in predicted evolution for the mean of each component of the process caused by a small displacement of parameter vector around $\hat{\theta}^{(k)}$. Finally,

$$W(\hat{\theta}^{(k)}) = \text{Diag} [W_1(\hat{\theta}^{(k)}) \quad W_2(\hat{\theta}^{(k)}) \quad \dots \quad W_S(\hat{\theta}^{(k)})]$$

Data: $x_{0:S}$: derive $dx_{1:S}$ and $M_{0:S-1}$ according to (7.5) and (7.6)

Result: Get parameters estimates $\hat{\theta}$

begin

Initialization: $\text{tol} = \epsilon$, $k = 0$;

$\hat{\theta}^{(0)} = \arg \min_{\theta} (dx_{1:S} - M_{0:S-1}\theta)^\top (dx_{1:S} - M_{0:S-1}\theta) \text{ s.t. } \theta \geq 0$;

while $(\|\widehat{\Delta\theta}^{(k)}\|_1) \geq \text{tol}$ **do**

Calculate $r(\hat{\theta}^{(k)})$, $J(\hat{\theta}^{(k)})$, $W_{0:S-1}(\hat{\theta}^{(k)})$;

$\widehat{\Delta\theta}^{(k)} =$

$\arg \min_{\Delta\theta} [r(\hat{\theta}^{(k)}) - J(\hat{\theta}^{(k)})\Delta\theta]^\top [W_{0:S-1}(\hat{\theta}^{(k)})]^{-1} [r(\hat{\theta}^{(k)}) - J(\hat{\theta}^{(k)})\Delta\theta]$

s.t. $\Delta\theta \geq -\hat{\theta}^{(k)}$;

$\hat{\theta}^{(k+1)} = \hat{\theta}^{(k)} + \widehat{\Delta\theta}^{(k)}$

$k = k + 1$;

end

$\hat{\theta} = \hat{\theta}^{(k)}$

end

Algorithm 1: Iterative procedure for the non-linear generalized method of moments based parameter estimation.

325 is the estimated $NS \times NS$ covariance matrix, setting the parameters vector
326 to current value $\hat{\theta}^{(k)}$.

327 For the local linear approximation method, some modifications to Algo-
328 rithm 1 have to be made. At each iteration, parameter refinement is not
329 performed by estimating increments vector $\widehat{\Delta\theta}$, but $\hat{\theta}^{(k)}$ directly by solving
330 the generalized (constrained) least square problem in (7.7) with covariance
331 matrix calculated using $\hat{\theta}^{(k-1)}$.

332 **5. Simulation study.** In this section, we present four simulation stud-
333 ies. In the first, we study the behavior of the non-linear inference procedure
334 simulating the data under that very model. In particular, we compare the
335 method to a linear alternative, known as the local linear approximation,
336 for several sampling intervals. For short sampling intervals, it is expected
337 that the local linear approximation will be a serious competitor, whereas for
338 longer sampling intervals the non-linearity will start to favor our non-linear
339 inference scheme. In the second simulation study we mimic an experimental
340 setting scenario by perturbing clones trajectory with multiplicative errors
341 before performing inference. Our goal here is to investigate how an addi-
342 tional and extrinsic source of variation affect parameter estimation. The
343 third simulation study focuses on how our model deals with model misspec-
344 ification. Although our model is a detailed and generative model of the cell
345 differentiation process, it is almost certain that this model is wrong — as all
346 models are (Wit, Heuvel and Romeijn, 2012). We report the performance

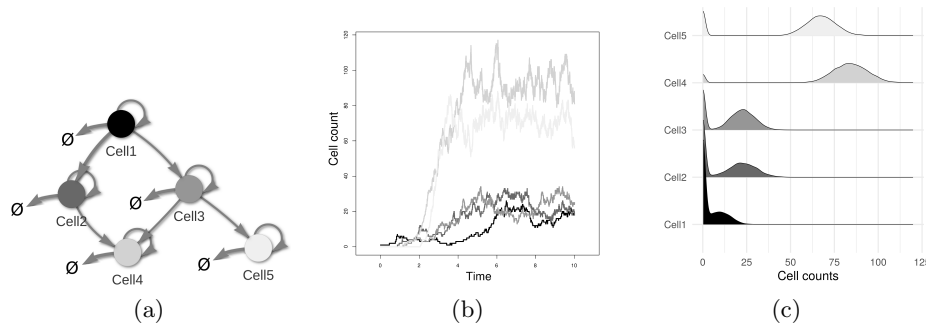


FIG 2. *Details on the cell differentiation process used in the simulation study.* A) Structure of the 5 cell types stochastic cell differentiation process. Cell types are represented by nodes. Self-connecting edges are duplication events. Death events are expressed by edges pointing to \emptyset . Edges connecting two nodes correspond to differentiation paths. B) An example of cell differentiation process trajectory (clone evolution) generated by means of Gillespie algorithm. C) Cell types multi-modal steady-state distribution calculated using 1000 trajectories.

347 of our model in recovering the differentiation process under an alternative
348 generative model. The fourth simulation study compares our proposal to an
349 alternative method-of-moments formulation proposed by Xu et al. (2019),
350 based on matching model-based and empirical correlations among cell types
351 dynamics.

352 5.1. *Improvement over local linear approximation approach.* The infer-
353 ence procedure presented in this paper requires one to calculate as many
354 solutions of the system of non-linear ODEs related to the first and second
355 moments of the process, as available observations. In Figure 2a the network
356 representation of the simulated system is shown. The precise parameter set-
357 tings are given in supplementary materials Supplement D.

358 The stochastic cell differentiation process implemented has been designed
359 with a low number of cell differentiations (5 out of 20) to reflect the expected
360 scenario of real biological systems. The simulation study aims to determine
361 whether our procedure is capable of correctly estimating the process param-
362 eters (both positive and zeros) and to investigate its performance for different
363 sampling intervals. Clone dynamics are simulated employing the Gillespie
364 algorithm (Gillespie, 1977), known to generate statistically correct trajec-
365 tories of the stochastic equation described in (3.1). An illustrative trajec-
366 tory is shown in Figure 2b, where it is possible to appreciate the logistic be-
367 havior generated by the model specification. Continuous-time trajectories are
368 then sampled at three different equally spaced time intervals $\Delta t = (0.1, 0.5, 1)$

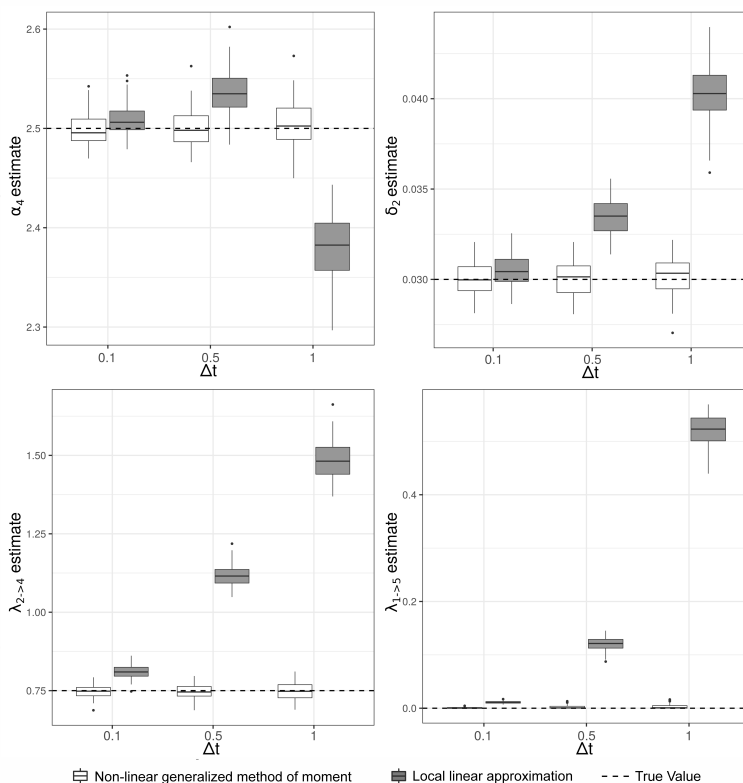


FIG 3. *Comparison between the non-linear generalized method of moments and local linear approximation for different Δt setting. Estimate distributions for the non-linear generalized method of moments and the local linear approximation are displayed using respectively white and gray boxplots. Dashed lines correspond to the true values. On top-left the performance of the methods for the estimation of the duplication rate α_4 (2.5) is shown. On top-right, death rate δ_2 (0.03) is reported. On the bottom the differentiation rates $\lambda_{1,3}$ (0.35, left) and $\lambda_{1,5}$ (0, right) are represented.*

369 until stopping time $t_{\text{end}} = 10$ is reached. Parameter estimates obtained by
 370 using the proposed algorithm and the local linear approximation approach
 371 are compared for 100 experiments, each composed of $n = 1000$ clones start-
 372 ing from initial conditions vector $\mathbf{x}_0 = (1, 0, 0, 0, 0)$. Having clone evolu-
 373 tions starting from a single cell makes steady-state behavior particularly sensi-
 374 tive to the initial (stochastic) sequence of cellular events. In Figure 2c the
 375 distribution at t_{end} , calculated based on 1000 clone trajectories, highlight
 376 the presence of a multi-modal steady-state configuration. On average, our
 377 algorithm converges in 2.8, 4.2 and 5.9 iterations, respectively, for Δt equal

378 to 0.1, 0.5 and 1. The local linear approximation approach converged on
379 average in 3.2, 6.2 and 7.2 iterations.

380 As shown in Figure 3, the local linear approximation based method suffers
381 in terms of accuracy in all settings. Due to the limited amount of cells present
382 in the system in the initial phases, the strong non-linearity component of the
383 dynamics is poorly approximated by the linear approach. As a consequence,
384 there is a considerable and fast decay of estimation precision as Δt increases.
385 For $\Delta t = 0.1$, the local linear approximation approach seems to be able to
386 recognize the underlying structure of the system, since almost all absent
387 differentiation paths are correctly estimated as very closed to zero. This
388 is not true for larger time gaps Δt , e.g., 0.5 or 1, where, in addition to a
389 considerable bias for all estimates, some of the absent links – for example,
390 $\lambda_{1,5}$ is shown in Figure 3 (bottom-right) – are systematically estimated as
391 greater than 0. The non-linear inference procedure, instead, shows unbiased
392 estimates for all Δt considered for all parameters.

5.2. *Performance introducing measurement errors.* To investigate how
measurement errors affect the performance of our proposal for inference, we
apply our algorithm to perturbed clone trajectories, $\tilde{\mathbf{x}}_s$. These trajectories
are generated by adding noise to the exact one, \mathbf{X}_s , as follows

$$(5.1) \quad \tilde{X}_{si} = \begin{cases} X_{si} + \tilde{\epsilon}_{si} & \text{if } X_{si} + \tilde{\epsilon}_{si} > 0 \\ 0 & \text{if } X_{si} + \tilde{\epsilon}_{si} \leq 0 \end{cases} \quad \text{and } \tilde{\epsilon}_{si} \sim \mathcal{N}(0, \varphi x_{s-1,i})$$

393 We considered the same system configuration and experiment setup as de-
394 scribed in section 5.1, using $\Delta t = 1$ and inspecting the impact of noise of
395 different strength by testing $\varphi = (0, 0.1, 0.5, 1)$.

396 In Figure 4 the performance in estimating a duplication rate (α_4), death
397 rate (δ_2), differentiation rate ($\lambda_{2,4}$) and an absent differentiation path ($\lambda_{1,5}$)
398 is shown. For all parameters, we observed an increase in the standard errors
399 as the value of φ increases. A shift in the parameter distribution is observed
400 for death and differentiation rates for the larger values of φ , but not for
401 the duplication coefficient. Most likely for large values of φ , as the states
402 are artificially truncated at 0, probably a bias is introduced. The higher
403 $\lambda_{1,5}$ average estimates we observed for φ values (0.5,1) is presumably due
404 to the increase of the estimator standard error. The vast majority of $\lambda_{1,5}$
405 estimates fall in the (0,6e-4) range, suggesting that the correct identification
406 of missing differentiation paths is robust to higher levels of observational
407 errors. To recover the underlying network structure and eliminate potential
408 spurious, low-intensity connections among lineages, in section 6 we propose

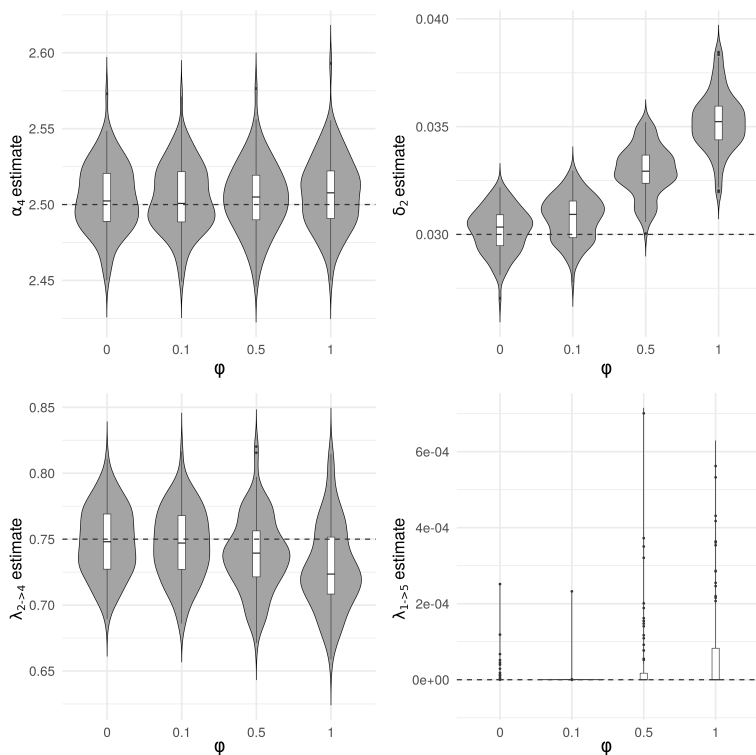


FIG 4. **Impact of measurement errors on inference.** On top-left the performance of methods for duplication rate α_4 (2.5) estimation. On top-right the performance of methods for death rate δ_2 (0.03) estimation. On bottom-left the performance of methods for differentiation rate $\lambda_{2,4}$ (0.75) estimation. On bottom-right the performance of methods for absent differentiation path $\lambda_{1,5}$ (0) estimation.

409 a model selection strategy based on backward stepwise selection and cross-
410 validation.

411 5.3. *Performance under model misspecification.* There are various dif-
412 ficulties associated with modeling biological processes, in particular when
413 dealing with questions related to the *in-vivo*, *in human*, investigation of
414 complex phenomena such as hematopoiesis. Many reasons limit sample size
415 and the type of experiments that can be performed, forcing the researcher in
416 making important assumptions about biological mechanisms based on evi-
417 dence gathered from *in-vitro* or animal studies, not always representative of
418 human dynamics. For these motivations, it is important to check how new
419 statistical procedures behave in case of model misspecification. In order to

420 test our proposal described in section 4 under this condition, we generated
421 clone trajectories using a corrupted version of the Gillespie algorithm. Differ-
422 entiation process structure has been kept as shown in Figure 2a. Parameters
423 have been set to the same values as reported in section 5.1, except for death
424 rates set to $\tilde{\delta} = (0.0, 0.3, 0.6, 2.0, 2.5)$. Individual clone evolution has been
425 simulated as described in the following steps:

- 426 1. Set initial state at $\mathbf{x}_0 = (1, 0, 0, 0, 0)$.
- 427 2. Generate time-to-next event, t_{s+1} , sampling from a Uniform distribu-
428 tion with parameters $\text{Unif}[0, (\frac{2}{\sum_{i=1}^N x_{s,i}})]$.
- 429 3. Select a cell type, $C_{s+1,i}$ sampling with probability proportional to cell
430 count among those cell type with $C_{s,i} \geq 1$.
- 431 4. Sample a cell event (duplication, death or differentiation) among those
432 available for the specific cell type $C_{s+1,i}$ with probability proportional
433 to event rates.
- 434 5. If total event time is less than 10, return to step 2.

435 It is worth noting that these modifications affect multiple aspects of the data
436 generating process, as visible from Figure 5a. Events frequency is much lower
437 throughout the simulation period and cell counts do not stabilize around a
438 cell type-specific value, as was the case for the original model shown in Fig-
439 ure 2b, but they rather exhibit exponential growth dynamics. In the correct
440 version of the Gillespie algorithm, the time-to-next-event is distributed as
441 an exponential with parameter $\exp(\sum_{k=1}^r h_k(\mathbf{X}_t; \boldsymbol{\theta}))$ and the same vector of
442 events hazard $h_k(\mathbf{X}_t; \boldsymbol{\theta})$ is rescaled to the unit sum in order to define events
443 sampling probabilities. Under the misspecification setting, the event times
444 are distributed uniformly, and the event probabilities are not directly linked
445 to the hazards.

446 Three different sample sizes have been tested: 30, 50, and 100 clones per
447 experiment. We evaluate our inference method for its capability to correctly
448 reconstruct the underlying differentiation structure, rather than for the pre-
449 cision in parameters estimation. Based on the data generated from a single
450 experiment, we test the null hypothesis $H_0 : \lambda_{ij} = 0, i, j = 1, \dots, N, i \neq j$
451 as described in [Supplement C](#).

452 Each ROC curve in Figure 5b shows the average of 100 ROC curves ob-
453 tained from independent replicates of the simulation experiments by varying
454 the significance threshold on differentiation rates. Our generalized method
455 of moments approach shows surprising accuracy in learning the true net-
456 work configuration for 30, 50, and 100 clone trajectories for a wide range of
457 significance threshold values.

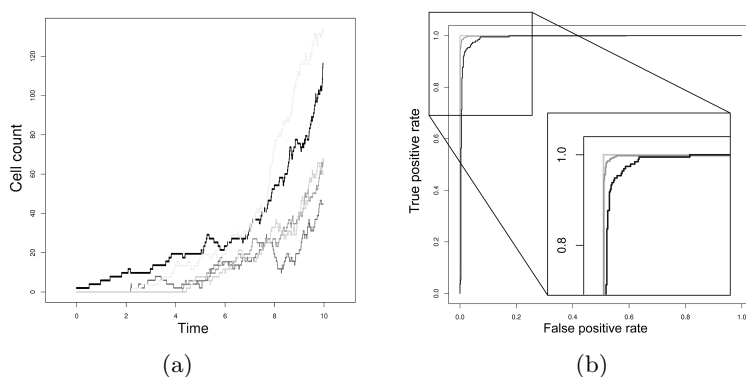


FIG 5. *Misspecified cell differentiation process.* A) Cell differentiation process trajectory generated by means of the misspecified generative model. B) Average ROC curves obtained from 100 experiments replicates each containing 30 (black), 50 (dark gray) and 100 (light gray) clones.

458 5.4. *Comparison with correlation-based M-estimator by Xu et al. (2019).*
459 Although the general model specification of the clonal expansion dynamics
460 by means of stochastic differential equations is similar to that described in
461 Xu et al. (2019), there are several differences between the data-handling and
462 estimation approaches. Wu et al. (2014) only have clone size measurements in
463 5 mature blood lineages and no information on the progenitors. To estimate
464 the hidden relationships among stem and progenitors cells, Xu et al. (2019)
465 resort to comparing known tree-like differentiation configurations by means
466 of cross-validation. Furthermore, in order to obtain an analytical solution for
467 the moments evolutions, they assume event hazards to be linear in process
468 states. This is probably the only sensible workable assumption, but it does
469 imply either exponential extinction and growth dynamics of the clones. On
470 the other hand, the gene therapy study motivating our method consists of 15
471 cell types from both BM and PB, providing a much more detailed description
472 of the complete hematopoietic process. Given this motivation, we designed
473 a modeling approach that assumes all lineages of interest to be observed.

474 In order to compare the two methods, we modified our methodology to
475 consider, as in Xu et al. (2019), asymmetric division (differentiation is cou-
476 pled to cell division) rather than symmetric division, whereby cell duplica-
477 tion is followed by a differentiation event. Furthermore, to match the two
478 stochastic processes we assumed that the dynamics does not involve satu-
479 ration by assuming linear ODEs. To make a reasonable comparison among
480 the two methods under the fully observed scenario, we extended the calcu-

481 lation of correlation-based M-estimator proposed by Xu et al. (2019) to all
482 correlations among lineages, including stem cells, progenitors and mature
483 cell types.

484 We set up a simulation study resembling the one described in Figure 2c
485 in Xu et al. (2019) (reproduced here in Figure 6a) both in terms of the
486 differentiation tree structure and the rate parameters. The process consists
487 of 8 cell types, starting from a *HSC* that duplicates with rate $\lambda = 0.285$ and
488 differentiates in progenitor cells, *Prog A* and *Prog B*, with rates $\nu_a = 0.14$ and
489 $\nu_b = 0.07$, respectively. Progenitor cell-types A die with rate $\mu_a = 0.14$ and
490 differentiate into three mature cell types with rates $\nu_1 = 36$, $\nu_2 = 18$ and
491 $\nu_3 = 10$, respectively. Progenitor cell-types B have two connected mature
492 lineages into which it differentiates with rates $\nu_4 = 20$ and $\nu_5 = 12$. As done
493 in Xu et al. (2019), we considered mature cells death rates known and equals
494 to $\mu_1 = 0.26$, $\mu_2 = 0.13$, $\mu_3 = 0.11$, $\mu_4 = 0.16$ and $\mu_5 = 0.09$. All trajectories
495 start with a single HSC at time $t_{\text{start}} = 0$. Each simulation experiment is
496 composed of 1000 clones, observed at intervals $\Delta t = 1$ unit apart, from
497 $t_{\text{start}} = 0$ up to the final time-point set at $t_{\text{end}} = 10$. The results of the
498 simulation study and the distributions of the coefficient estimates across
499 100 simulations are shown in Figure 6b.

500 Our proposal outperformed the method of Xu et al. (2019) in several as-
501 pects. The precision of our estimates is an order of magnitude better, and
502 the bias of our method is negligible, whereas their estimation of μ_a, μ_b, ν_a
503 and ν_b clearly suffers from bias. Furthermore, our computational algorithm
504 converged in 4.3 iterations on average, whereas the correlation-matching al-
505 gorithm converges on average in 60.6 iterations. The reason why our method
506 outperforms the method proposed by Xu et al. (2019) is that latter based
507 on second moment matching, whereas our method is based on first moment
508 matching, which is more stable, unbiased and computationally more efficient.
509 On the other hand, the main advantage of the method proposed by Xu et al.
510 (2019) is that their method can deal efficiently with missing progenitor and
511 HSC data. In certain experimental settings this can be crucial.

512 **6. Gene therapy study for Wiskott-Aldrich Syndrome.** In this
513 section, we return to the previously described clinical trial treating patients
514 suffering from Wiskott-Aldrich Syndrome with their stem cells, genetically
515 modified ex vivo, and then reinfused to the patient. We traced $N = 15$ cell
516 types over time in the three patients up to 36 months after GT. In Figure 7a
517 the differentiation trajectories observed for two clones are shown. The 15
518 distinct cell types can be organized in a three levels hierarchy, corresponding
519 to the original *HSC level*, i.e., CD34 stem cells, the *bone marrow (BM)*

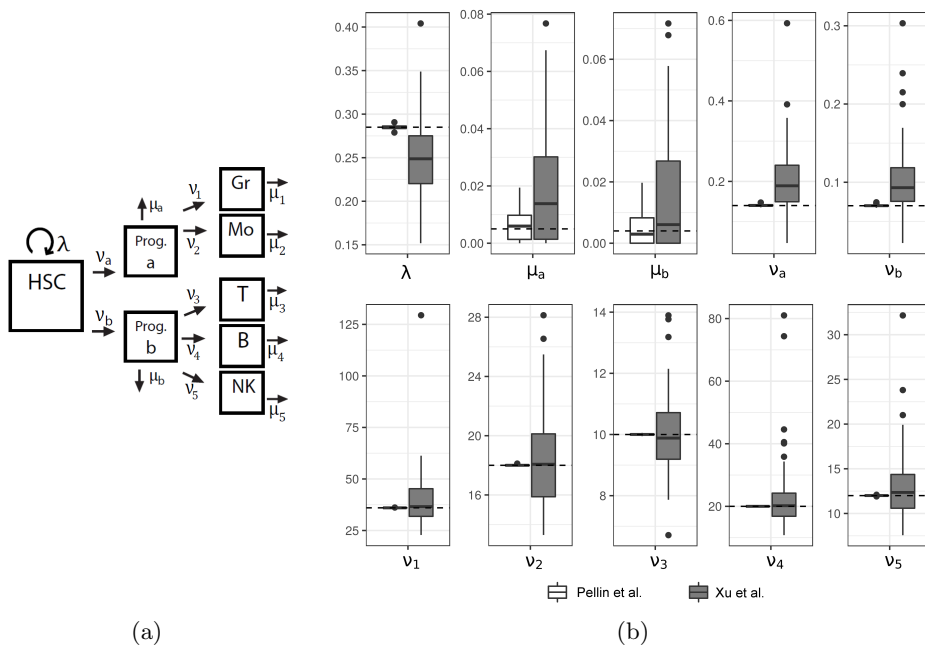


FIG 6. *Comparison with Xu et al. (2019) correlation-based M-estimator.* Considering the cell differentiation process shown (a), the boxplots in (b) show that our method is unbiased and more efficient than the M-estimator proposed in Xu et al. (2019). Boxplots show the distribution of estimates obtained using the method proposed in this manuscript (white) and the extended version of the correlation-based estimator (dark gray) for the 10 unknown rates, whose true value is indicated by the horizontal red dashed line.

520 level, corresponding to CD3, CD14, CD15, CD19, CD56, CD61 and GLYCO
 521 precursor cells and finally the *peripheral blood (PB) level*, i.e., CD3, CD4,
 522 CD8, CD14, CD15, CD19 and CD56 mature cells. Based on the available
 523 biological knowledge, the following assumptions are made,

- 524 • the HSC type can differentiate in any cell type in the BM level;
- 525 • cell types at the BM level can differentiate in any cell type in the PB
 526 level;
- 527 • cell types at the PB level can not differentiate.

528 These assumptions are graphically summarized in Figure 7b and incorpo-
 529 rated in the stochastic cell differentiation model and inferential algorithm
 530 by setting the corresponding λ_{ij} to zero.

531 From a practical perspective, the re-infusion of corrected HSC cells in a
 532 patient's body is considered as starting time $t = 0$. Initial conditions vector
 533 \mathbf{X}_0 consists of a 15-dimensional vector, with the count corresponding to

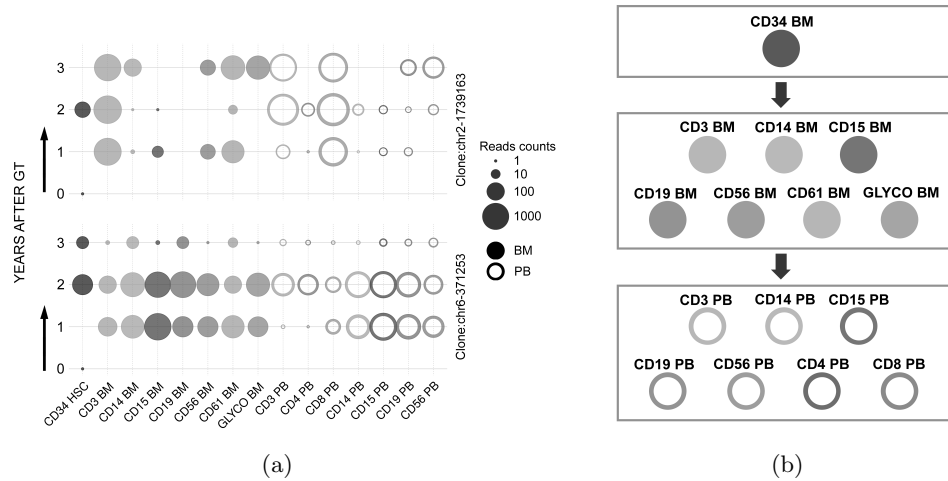


FIG 7. *Observed clones dynamics and schematic representation of hierarchical assumptions.* A) Reads counts trajectory over the 3 years follow-up for 2 clones. By assumption, all clones start with a count of 1 in CD34 HSC cell type at time 0. Cell count and are represented by circle of size proportional to their abundance. B) Schema reporting the biologically inspired three levels hierarchy used as a backbone for differentiation structure reconstruction. Arrows show directionality for potential differentiation paths.

534 CD34, the HSC, equal to 1 and the rest to zero. During the follow-up pe-
 535 riod, $S = 3$ samples from patient's HSC, BM, and PB cells are taken after 1,
 536 2, and 3 years. After excluding all clones detected only once throughout the
 537 study period, in total we obtain 17,195 unique chromosomal positions: 5,299
 538 from period 1, 5,300 from period 2, and 6,596 from period 3. The amount of
 539 cells, within each lineage, generated by individual labeled, re-infused HSC, is
 540 counted through an insertion site analysis technique described in [Aiuti et al.](#)
 541 (2013). For estimating the measurement error scaling coefficient associated
 542 with the protocol used in the processing of patients' samples, we took ad-
 543 vantage of the three independent experiments in which a pool of HSC cells
 544 have been sequenced 1-day after transduction. Given the low proliferative
 545 rate of HSC in culture conditions, all clones are expected to have a size of 1
 546 at time of sequencing. Based on these data we estimated $\hat{\phi} = 0.08$.

6.1. *Cell differentiation reconstruction.* Clonal tracking studies typically score and compare alternative but fixed models of hematopoiesis using experimental data. In this work, we opted for data-driven learning of the differentiation process structure. To recover the actual underlying data generating

process and eliminate differentiation paths caused by sampling issues and observational errors, we proceed as following described. We estimated the full model, m_0 , by solving the optimization problem 4.2 using the WAS data and $\hat{\varphi}$. We then iteratively eliminate the differentiation connection (λ_{ij}) with the least impact on the following Mahalanobis distance:

(6.1)

$$D_M = [\mathbf{x}_{1:S} - f(\mathbf{x}_{0:S-1}; \hat{\boldsymbol{\theta}}^k)]^\top (\mathbf{W}_{1:S} + \hat{\varphi} \mathbf{N}_{1:S})^{-1} [\mathbf{x}_{1:S} - f(\mathbf{x}_{0:S-1}; \hat{\boldsymbol{\theta}}^k)]$$

547 This method leads to a sequence of models, $m_k, k = 1 \dots 56$ with decreasing
548 complexity. To select the optimal model \tilde{m} among the set m_k , we used a
549 5-fold cross-validation strategy. We split the input dataset into five subsets
550 of equal size and used four subsets to estimate the process parameters and
551 the remaining as a validation subset on which the Mahalanobis distance
552 (6.1) has been calculated. The procedure has been repeated five times for
553 each model configuration, considering each subset for validation once. The
554 results are reported in Figure 8. We selected model m_{35} as optimal based
555 on its minimum median Mahalanobis distance across folds.

556 We then imposed the differentiation structure encoded in model m_{35} and
557 estimates the cell differentiation process parameters using all WAS data avail-
558 able. A graphical representation of the differentiation network is shown in
559 Figure 9a. Duplication and death have been omitted in the plot for clarity,
560 but all final parameters are available in supplementary materials [Supple-](#)
561 [ment F](#). In Figure 9b a trajectory of the HSC differentiation process esti-
562 mated using WAS gene therapy data is shown, generated using the Gillespie
563 algorithm.

564 Initialization with the local linear approximation aims at starting the
565 optimization procedure in the proximity of the objective function global
566 optima and reducing the number of iterations (m_{35} converges in 5 iterations)
567 required to meet the convergence criteria. We verified that the parameters
568 estimate in Appendix [Supplement G](#) are stable to random initialization by
569 sampling candidate values for $\hat{\boldsymbol{\theta}}^{(0)}$ from a Normal distribution $\mathcal{N}(0.1, 0.1)$
570 for duplication and differentiation rates and $\mathcal{N}(0.01, 0.01)$ for death rates.
571 We performed 100 random restarts showing that our estimates are robust.

572 *6.2. Relevance of the results.* CD34 HSC resulted in being the lineage
573 with the highest duplication rate. According to our estimate, a CD34 HSC
574 cell is expected to duplicate approximately every 6.51 weeks ($\alpha_{CD34_HSC} =$
575 $8.006e + 00$), a significantly higher rate than the 40 weeks (range, 25-50
576 weeks) previously reported ([Catlin et al., 2011](#)). The difference is attributable
577 to the following considerations. First, the patients enrolled in a GT clinical
578 trial receive a conditioning regimen before treatment. Upon reinfusion,

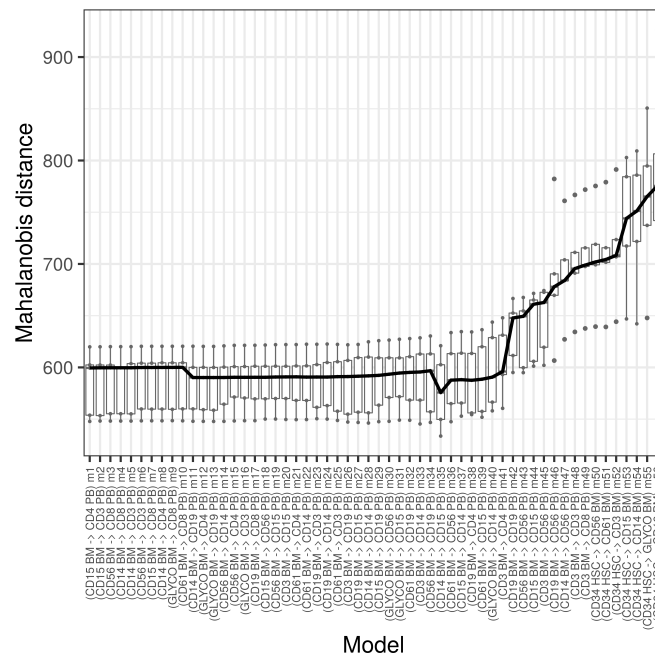


FIG 8. *Cross-validation and Mahalanobis distance based selection of the optimal model.* Models m_k are ordered from the most complex, m_1 (one differentiation path removed from the full model m_0) to the least complex, m_{56} (no connection among lineages). The additional differentiation path that is removed at each iteration is reported alongside model number. The distributions of the Mahalanobis distances calculated on the 5 validation subsets are represented with boxplots for each model configuration, m_k . Solid black line connects the median distances across models. The minimum median is observed for m_{35} that is therefore selected as the optimal model, \hat{m} .

579 the transduced cells are subjected to high proliferative stress because they
580 must replenish the depleted hematopoietic system. The estimate reported in
581 [Catlin et al. \(2011\)](#) instead is referred to a healthy, native, steady-state con-
582 dition and does not consider potential selective advantages that engineered
583 cells might have in disease settings. Second, the CD34 marker used in the
584 WAS study to isolate HSC from patients' BM samples is known to select
585 for a broader cell population that includes hematopoietic progenitor cells
586 in addition to stem cells, which are characterized by a higher proliferative
587 output and shorter half-lives compared to pure hematopoietic stem cells.

588 BM lymphoid lineages CD3 and CD19 show higher duplication coefficients
589 than myeloid cell types (CD14 and CD15). This result supports the idea of
590 the presence of long-lived lymphoid progenitors and the dependence of the

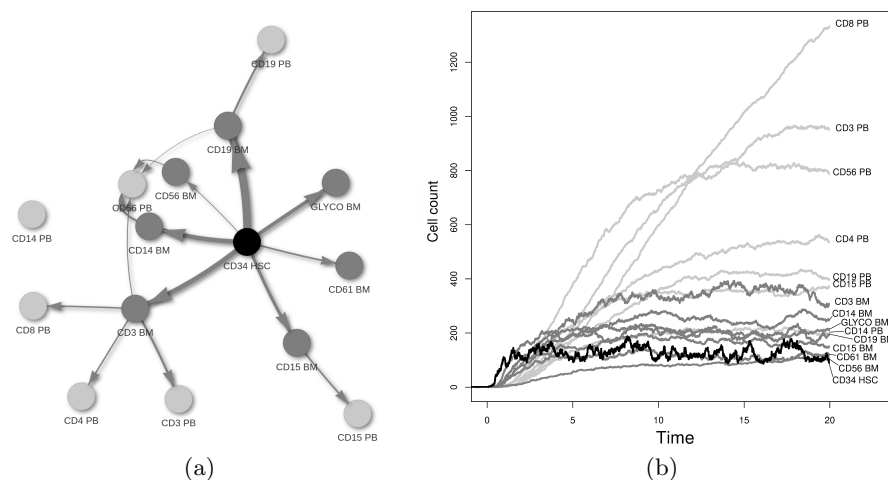


FIG 9. **HSC differentiation process.** (a) Network representation. Black, dark gray and light gray nodes represent CD34 HSC, BM, and PB cell types respectively. Edge thickness is proportional to the corresponding $\hat{\lambda}_{ij}$. Edges estimated are those included in the optimal model \tilde{m} . (b) HSC differentiation process trajectory simulated using the Gillespie algorithm, assuming model \tilde{m} and coefficient estimated using WAS data.

591 myeloid compartment from the continuous support of cells coming from the
 592 upstream CD34 HSC population (see supplementary materials [Supplement](#)
 593 [G](#)). CD61 BM cells are estimated to have a significant duplication rate. The
 594 distinct behavior of the megakaryocyte (CD61 BM) population is not sur-
 595 prising since megakaryo/erythrocyte-restricted progenitor, responsible for
 596 the production of platelets and red blood cells (erythrocytes), have been
 597 reported and validated in several studies, mostly based on gene expression
 598 data. Steady-state cell counts for individual lineages are not deterministic
 599 but depend on the specific evolution of each clone (see Figure 2c). However,
 600 in Figure 9b it is possible to appreciate how the combination of duplication
 601 and death rates estimate leads to a biologically meaningful differentiation
 602 process in which PB lineages are the most abundant, followed by BM and
 603 CD34 HSC.

604 In the optimal model configuration determined by our model selection
 605 strategy (Figure 9a), all BM lineages result directly connected to the HSC
 606 compartment. Surprisingly, HSC to B-Cell precursor ($\lambda_{CD34_HSC \rightarrow CD19_BM} =$
 607 1.453) differentiation rate is higher than HSC to myeloid cells (CD15 BM,
 608 CD14 BM), which are among the cell type with the fastest turnover in hu-
 609 mans ([Sender and Milo, 2021](#)). This finding agrees with the conclusion of
 610 [Meyer-Bahlburg et al. \(2008\)](#) who, using mouse models of WAS, highlighted

611 that upon transplantation, corrected B cells exhibit a marked selective ad-
612 vantage at both the precursor and mature stage.

613 The biology behind the maturation and migration of BM cells in the PB
614 stream is much better understood, and commitment paths are well charac-
615 terized. The consistency of our inferred structure at the BM and PB interface
616 with the biological expectation is remarkable, even though a limited set of
617 constraints to the network configurations has been provided. The separation
618 between lymphoid and myeloid branches is clear, with significant differen-
619 tiation parameters connecting CD3 at BM level to CD3-CD4-CD8 (T-cell)
620 and CD56 (NK) in the PB. Among the myeloid subpopulations, CD15 BM
621 is linked to CD15 PB as expected, but the differentiation from CD14 BM
622 to CD14 PB is missing. The isolation of CD14 PB from all BM lineages is
623 most likely a sampling issue since monocyte (CD14) account, on average,
624 for only 5% of the cells in a PB sample.

625 Our results support the myeloid-based model over the classical dichotomy
626 model. Mature NK cells (CD56 PB) are sustained by a cellular influx from
627 NK cells residing in the BM (CD56 BM), as expected, but also from CD14
628 BM (myeloid), CD19 BM, and CD3 BM (lymphoid lineages). Although it
629 is biologically challenging to conclude that all these cell populations can
630 directly give rise to CD56 PB cells, this pattern is compatible with the
631 presence of a common, unobserved progenitor cell type capable of generating
632 both myeloid and lymphoid mature cells.

633 Due to the poor approximation provided by the local linear method, as
634 also shown in our simulation study, [Biasco et al. \(2016\)](#) identified many
635 more low-intensity, most likely spurious, differentiation rates. For this rea-
636 son, the authors preferred to limit the inferential goal to calculate and com-
637 pare the likelihoods of only two known and competing tree configurations
638 using information-based criteria. Instead, the method presented in this pa-
639 per allows us to perform network and coefficients estimation simultaneously.
640 It requires only limited prior knowledge and is essentially data-driven. Nev-
641 ertheless, it also offers the flexibility to trade exploratory power for biological
642 interpretability by changing the settings of the differentiation rates fixed at
643 zero according to the scientific question.

644 Finally, to resolve the conundrum regarding *in-vivo* stem cell evolution
645 and hematopoietic differentiation structure, a more refined sorting strategy
646 for HSC (CD34 BM) is needed. Through additional known surface markers,
647 indicators of stem/progenitor cells priming towards specific lineages would
648 be possible to disentangle the complexity observed at the BM level.

649 **7. Conclusion.** To improve our knowledge about the cell differentia-
650 tion process, which in many contexts such as gene therapy might be fun-
651 damental for providing biological and therapeutic new insights, we have
652 devised and implemented a flexible statistical framework for the analysis of
653 clonal tracking data. The underlying stochastic process is assumed to be
654 a multidimensional Markov process and this allows a representation of the
655 moment dynamics by means of a system of non-linear ODEs. The partic-
656 ular definition of the transition probabilities induces a logistic behavior of
657 sub-population growth curves. The model and the proposed iterative infer-
658 ential procedure exhibit stability in terms of parameter estimation, structure
659 recognition, and convergence rate. The model can easily be extended to in-
660 corporate time-dependent individual cell rates, different feedback regulation
661 mechanisms, or random effects on specific parameters.

662 Applying the modeling and inference framework to a Wiskott-Aldrich
663 Syndrome gene therapy study, we have obtained insight into the underlying
664 stem cell differentiation dynamics. We found a high degree of agreement
665 between our results and the recently proposed myeloid-based model for hu-
666 man hematopoiesis *over* the predominant classical dichotomy model of cell
667 evolution.

668 References.

- 669 AHNERT, K. and MULANSKY, M. (2011). Odeint - Solving ordinary differential equations
670 in C++. *CoRR* **abs/1110.3397**.
- 671 AIUTI, A., BIASCO, L., SCARAMUZZA, S., FERRUA, F. . . and NALDINI, L. (2013). Lentivi-
672 ral Hematopoietic Stem Cell Gene Therapy in Patients with Wiskott-Aldrich Syndrome.
673 *Science* **341**.
- 674 BAILEY, N. (1964). *T., J.: The Elements of Stochastic Processes*. John Wiley.
- 675 BATES, D. and MAECHLER, M. (2015). Matrix: Sparse and Dense Matrix Classes and
676 Methods R package version 1.2-2.
- 677 BERRY, C. C., GILLET, N. A., MELAMED, A., GORMLEY, N., BANGHAM, C. R. and
678 BUSHMAN, F. D. (2012). Estimating abundances of retroviral insertion sites from DNA
679 fragment length data. *Bioinformatics* **28** 755–762.
- 680 BIASCO, L., AMBROSI, A., PELLIN, D., BARTHOLOMAE, C., BRIGIDA, I., RONCAR-
681 OLO, M. G., DI SERIO, C., VON KALLE, C., SCHMIDT, M. and AIUTI, A. (2011).
682 Integration profile of retroviral vector in gene therapy treated patients is cell-specific
683 according to gene expression and chromatin conformation of target cell. *EMBO Molec-*
684 *ular Medicine* **2** 1757-4684.
- 685 BIASCO, L., PELLIN, D., SCALA, S., DIONISIO, F., BASSO-RICCI, L., LEONARDELLI, L.,
686 SCARAMUZZA, S., BARICORDI, C., FERRUA, F., CICALESE, M. P. et al. (2016). In vivo
687 tracking of human hematopoiesis reveals patterns of clonal dynamics during early and
688 steady-state reconstitution phases. *Cell Stem Cell* **19** 107–119.
- 689 BIFFI, A., MONTINI, E., LORIOLO, L., CESANI, M., FUMAGALLI, F., PLATI, T., BAL-
690 DOLI, C., MARTINO, S., CALABRIA, A., CANALE, S. et al. (2013). Lentiviral hematoi-

- 691 etic stem cell gene therapy benefits metachromatic leukodystrophy. *Science* **341**
692 1233158.
- 693 BJÖRCK, A. (1996). *Numerical methods for least squares problems*. Siam.
- 694 BLANPAIN, C., HORSLEY, V. and FUCHS, E. (2007). Epithelial stem cells: turning over
695 new leaves. *Cell* 445–458.
- 696 CALABRIA, A., LEO, S., BENEDECENTI, F., CESANA, D., SPINOZZI, G., ORSINI, M.,
697 MERELLA, S., STUPKA, E., ZANETTI, G. and MONTINI, E. (2014). VISPA: a com-
698 putational pipeline for the identification and analysis of genomic vector integration
699 sites. *Genome medicine* **6** 1–12.
- 700 CATLIN, S. N., BUSQUE, L., GALE, R. E., GUTTORP, P. and ABKOWITZ, J. L. (2011).
701 The replication rate of human hematopoietic stem cells in vivo. *Blood, The Journal of*
702 *the American Society of Hematology* **117** 4460–4466.
- 703 ELF, J. and EHRENBERG, M. (2003). Fast evaluation of fluctuations in biochemical net-
704 works with the linear noise approximation. *Genome research* **13** 2475–2484.
- 705 GARDINER, C. W. (1985). *Handbook of Stochastic Methods*. Springer.
- 706 GILLESPIE, D. T. (1977). Exact stochastic simulation of coupled chemical reactions. *Jour-*
707 *nal of Physical Chemistry* **81** 2340–2361.
- 708 GOLIGHTLY, A. and WILKINSON, D. J. (2005). Bayesian inference for stochastic kinetic
709 models using a diffusion approximation. *Biometrics* **61** 781–788.
- 710 GOLIGHTLY, A. and WILKINSON, D. J. (2008). Bayesian inference for nonlinear multivari-
711 ate diffusion models observed with error. *Computational Statistics and Data Analysis*
712 **52** 1674–1693.
- 713 GRIMA, R. (2012). A study of the accuracy of moment-closure approximations for stochas-
714 tic chemical kinetics. *The Journal of chemical physics* **136** 04B616.
- 715 GUENNEBAUD, G., JACOB, B. et al. (2010). Eigen v3. <http://eigen.tuxfamily.org>.
- 716 IBM (2010). User’s Manual for CPLEX IBM ILOG CPLEX V12.1.
- 717 KAMPEN, N. G. V. (1981). *Stochastic Processes in Physics and Chemistry*.
718 Amsterdam:North-Holland.
- 719 KAWAMOTO, H., WADA, H. and KATSURA, Y. (2010). A revised scheme for developmental
720 pathways of hematopoietic cells: the myeloid-based model. *International immunology*
721 **22** 65–70.
- 722 KEELING, M. J. (2000). Multiplicative moments and measures of persistence in ecology.
723 *Journal of Theoretical Biology* **205** 269–281.
- 724 KILIAN, H., BARTKOWIAK, D., KAUFMANN, D. and KEMKEMER, R. (2008). The general
725 growth logistics of cell populations. *Cell biochemistry and biophysics* **51** 51–66.
- 726 LEONARDELLI, L., PELLIN, D., SCALA, S., DIONISIO, F., RICCI, L. B., CITTARO, D.,
727 DI SERIO, C., AIUTI, A. and BIASCO, L. (2016). 531. Computational Pipeline for the
728 Identification of Integration Sites and Novel Method for the Quantification of Clone
729 Sizes in Clonal Tracking Studies. *Molecular Therapy* **24** S212–S213.
- 730 MEYER-BAHLBURG, A., BECKER-HERMAN, S., HUMBLET-BARON, S., KHIM, S., WE-
731 BER, M., BOUMA, G., THRASHER, A. J., BATISTA, F. D. and RAWLINGS, D. J. (2008).
732 Wiskott-Aldrich syndrome protein deficiency in B cells results in impaired peripheral
733 homeostasis. *Blood, The Journal of the American Society of Hematology* **112** 4158–
734 4169.
- 735 NÅSELL, I. (2003a). An extension of the moment closure method. *Theoretical population*
736 *biology* **64** 233–239.
- 737 NÅSELL, I. (2003b). Moment closure and the stochastic logistic model. *Theoretical popu-*
738 *lation biology* **63** 159–168.
- 739 NAKARIAKOV, S. (2013). *The Boost C++ Libraries: Generic Programming*. CreateSpace
740 Independent Publishing Platform.

- 741 NALDINI, L. (2011). Ex vivo gene transfer and correction for cell-based therapies. *Nat Rev*
742 *Genet.* **12** 301-15.
- 743 PELLIN, D., BIASCO, L., AIUTI, A., DI SERIO, M. C. and WIT, E. C. (2019). Penalized
744 inference of the hematopoietic cell differentiation network via high-dimensional clonal
745 tracking. *Applied Network Science* **4** 115.
- 746 RISKEN, H. (1984). *The Fokker-Planck Equation*. Springer.
- 747 SCHNOERR, D., SANGUINETTI, G. and GRIMA, R. (2017). Approximation and inference
748 methods for stochastic biochemical kinetics—a tutorial review. *Journal of Physics A:*
749 *Mathematical and Theoretical* **50** 093001.
- 750 SENDER, R. and MILO, R. (2021). The distribution of cellular turnover in the human
751 body. *Nature Medicine* **27** 45–48.
- 752 SINGH, A. and HESPANHA, J. P. (2007). A derivative matching approach to moment
753 closure for the stochastic logistic model. *Bulletin of mathematical biology* **69** 1909.
- 754 SNIPPERT, H. J. and CLEVERS, H. (2011). Tracking adult stem cells. *EMBO Rep.* **12**
755 113–122.
- 756 STROUSTRUP, B. (1997). *The C++ Programming Language.*, third edition ed. Addison-
757 Wesley.
- 758 R CORE TEAM (2015). R: A Language and Environment for Statistical Computing R
759 Foundation for Statistical Computing, Vienna, Austria.
- 760 WEISSMAN, I. L. (2000). Stem cells: units of development, units of regeneration, and units
761 in evolution. *Cell* 157-168.
- 762 WHITTLE, P. (1957). On the use of the normal approximation in the treatment of stochas-
763 tic processes. *Journal of the Royal Statistical Society. Series B* **19** 268–281.
- 764 WILKINSON, D. J. (2006). *Stochastic Modelling for Systems Biology*. Chapman and Hall.
- 765 WIT, E., HEUVEL, E. v. D. and ROMELJN, J.-W. (2012). ‘All models are wrong...’: an
766 introduction to model uncertainty. *Statistica Neerlandica* **66** 217–236.
- 767 WU, C., LI, B., LU, R., KOELLE, S. J., YANG, Y., JARES, A., KROUSE, A. E., MET-
768 ZGER, M., LIANG, F., LORE, K. et al. (2014). Clonal tracking of rhesus macaque
769 hematopoiesis highlights a distinct lineage origin for natural killer cells. *Cell Stem Cell*
770 **14** 486–499.
- 771 XU, J., KOELLE, S., GUTTORP, P., WU, C., DUNBAR, C., ABKOWITZ, J. L., MININ, V. N.
772 et al. (2019). Statistical inference for partially observed branching processes with appli-
773 cation to cell lineage tracking of in vivo hematopoiesis. *The Annals of Applied Statistics*
774 **13** 2091–2119.

775 **Acknowledgements.** CdS and EW would like to acknowledge network-
776 ing support by the COST Action COSTNET (CA15109). EW acknowledges
777 funding from Swiss National Science Foundation (SNF grant: 188534).

SUPPLEMENTARY MATERIAL

778 “Tracking hematopoietic stem cell evolution in a Wiskott-Aldrich
 779 clinical trial”
 780
 781

Supplement A: Derivation of moments equations

(<http://www.e-publications.org/ims/support/download/imsart-ims.zip>). By means of the summation operator, $\sum_{\mathbf{x} \in \tilde{\mathbf{x}}}$, over the whole set of possible states for the process $\mathbf{X}(t)$, $\tilde{\mathbf{x}} = \mathbb{N}_0^N$, it is possible to derived a functional connection between the evolution for the expected population size of each process component and the dynamics of the process probability distribution $P(\mathbf{X}; t)$,

$$\begin{aligned} \frac{dm_i(t)}{dt} &= \frac{d \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i P(\mathbf{X} = \mathbf{x}; t)}{dt} \\ &= \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i \frac{dP(\mathbf{X} = \mathbf{x}; t)}{dt} \end{aligned}$$

The evolution of $P(\mathbf{X}; t)$ can be expressed by means of the master equation introduced in (3.1),

$$\frac{dm_i(t)}{dt} = \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i \sum_{k=1}^r [h_k(\mathbf{x} - \mathbf{V}_{k,\cdot}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x} - \mathbf{V}_{k,\cdot}; t) - h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t)]$$

Due to the fact that the summation operator $\sum_{\mathbf{x} \in \tilde{\mathbf{x}}}$ spans over all possible state configurations, the order of summation operators in the RHS can be inverted,

$$\begin{aligned} \frac{dm_i(t)}{dt} &= \sum_{k=1}^r \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i [h_k(\mathbf{x} - \mathbf{V}_{k,\cdot}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x} - \mathbf{V}_{k,\cdot}; t) - h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t)] \\ &= \sum_{k=1}^r \left[\sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i h_k(\mathbf{x} - \mathbf{V}_{k,\cdot}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x} - \mathbf{V}_{k,\cdot}; t) - \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) \right] \end{aligned}$$

Now, the summation variable in the first term of the right-end-side can be modified, without affecting the sum domain, since it cover all possible state

configurations,

$$\begin{aligned} \frac{dm_i(t)}{dt} &= \sum_{k=1}^r \left\{ \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} (x_i + v_{k,i}) h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) - \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) \right\} \\ &= \sum_{k=1}^r \left\{ \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) + v_{k,i} h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) - \right. \\ &\quad \left. \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) \right\} \\ &= \sum_{k=1}^r \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} v_{k,i} h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) \end{aligned}$$

Given the known property for expected value of function $f(x)$ of a r.v. x with probability distribution $P(x)$, $E[f(x)] = \sum_{\mathbf{x}} f(\mathbf{x}) P(\mathbf{x})$,

$$\frac{dm_i(t)}{dt} = \sum_{k=1}^r E[v_{k,i} h_k(\mathbf{X}_t; \boldsymbol{\theta})]$$

Finally, by linearity of expectation,

$$\frac{dm_i(t)}{dt} = \sum_{k=1}^r v_{k,i} E[h_k(\mathbf{X}_t; \boldsymbol{\theta})]$$

782 A similar approach can be extended to define a system of ODEs for the time
783 evolution for second order moments of $\mathbf{X}(t)$,

$$\begin{aligned} \frac{dm_{i,j}^2}{dt} &= \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i x_j \frac{dP(\mathbf{X} = \mathbf{x}; t)}{dt} \\ &= \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} x_i x_j \sum_{k=1}^r \{ h_k(\mathbf{x} - \mathbf{V}_{k,\cdot}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x} - \mathbf{V}_{k,\cdot}; t) - h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) \} \\ &= \sum_{k=1}^r \left\{ \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} v_{k,j} x_i h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) + \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} v_{k,i} x_j h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) \right. \\ &\quad \left. + \sum_{\mathbf{x} \in \tilde{\mathbf{x}}} v_{k,i} v_{k,j} h_k(\mathbf{x}; \boldsymbol{\theta}) P(\mathbf{X} = \mathbf{x}; t) \right\} \\ &= \sum_{k=1}^r v_{k,j} E[X_{ti} h_k(\mathbf{X}_t; \boldsymbol{\theta})] + \sum_{k=1}^r v_{k,i} E[X_{tj} h_k(\mathbf{X}_t; \boldsymbol{\theta})] + \sum_{k=1}^r v_{k,i} v_{k,j} E[h_k(\mathbf{X}_t; \boldsymbol{\theta})] \end{aligned}$$

785 **Supplement B: Example with $N=3$ cell types**

786 (<http://www.e-publications.org/ims/support/download/imsart-ims.zip>). In this
 787 section the most relevant elements defined in section 3 and section 4 are de-
 788 rived, to allow parameters inference for an illustrative hypothetical $N = 3$
 789 stochastic cell differentiation model. We define the parameters governing
 790 stochastic cell differentiation process as

- Individual cell duplication rates vector

$$\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \alpha_3);$$

- Individual cell death rates vector:

$$\boldsymbol{\delta} = (\delta_1, \delta_2, \delta_3);$$

- Individual cell differentiation rates:

$$\boldsymbol{\lambda} = \begin{bmatrix} 0 & \lambda_{12} & \lambda_{13} \\ \lambda_{21} & 0 & \lambda_{23} \\ \lambda_{31} & \lambda_{32} & 0 \end{bmatrix}.$$

According to the ordering rule described in section 3, the $r = 12$ distinct cellular events are associated with a vector of events rates, $\mathbf{h}(\mathbf{X}, \boldsymbol{\theta})$,

$$\mathbf{h}(\mathbf{X}, \boldsymbol{\theta}) = (\alpha_1 X_1, \alpha_2 X_2, \alpha_3 X_3, \delta_1 X_1^2, \delta_2 X_2^2, \delta_3 X_3^2, \lambda_{21} X_2, \lambda_{31} X_3, \lambda_{12} X_1, \lambda_{3,2} X_3, \lambda_{13} X_1, \lambda_{23} X_2);$$

and a net effect matrix \mathbf{V} ,

$$\mathbf{V} = \begin{bmatrix} 1 & 0 & 0 & -1 & 0 & 0 & 1 & 1 & -1 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & -1 & 0 & -1 & 0 & 1 & 1 & 0 & -1 \\ 0 & 0 & 1 & 0 & 0 & -1 & 0 & -1 & 0 & -1 & 1 & 1 \end{bmatrix}.$$

Within the local linear approximation framework described in section Supplement E, the diagonal matrix $D(\mathbf{X})$ corresponds to

$$D(\mathbf{X}) = \text{Diag}(X_1, X_2, X_3, X_1^2, X_2^2, X_3^2, X_2, X_3, X_1, X_3, X_1, X_2)$$

The ODEs systems for time evolutions of process first-order moments is given by

$$\begin{cases} \frac{dm_1(t)}{dt} = \alpha_1 m_1(t) - \delta_1 m_{11}^2(t) + \lambda_{21} m_2(t) + \lambda_{31} m_3(t) - \lambda_{12} m_1(t) - \lambda_{13} m_1(t); \\ \frac{dm_2(t)}{dt} = \alpha_2 m_2(t) - \delta_2 m_{22}^2(t) - \lambda_{21} m_2(t) + \lambda_{12} m_1(t) + \lambda_{3,2} m_3(t) - \lambda_{23} m_2(t); \\ \frac{dm_3(t)}{dt} = \alpha_3 m_3(t) - \delta_3 m_{33}^2(t) - \lambda_{31} m_3(t) - \lambda_{3,2} m_3(t) + \lambda_{13} m_1(t) + \lambda_{23} m_2(t); \end{cases}$$

and for second-order moments

$$\left\{ \begin{array}{l} \frac{dm_{11}^2(t)}{dt} = (\alpha_1 m_1(t) + \delta_1 m_{11}^2(t) + \lambda_{21} m_2(t) + \lambda_{31} m_3(t) + \lambda_{12} m_1(t) + \\ \lambda_{13} m_1(t)) + 2(\alpha_1 m_{11}^2(t) - \delta_1 E[X_1^3] + \lambda_{21} m_{12}^2(t) + \\ \lambda_{31} m_{13}^2(t) - \lambda_{12} m_{11}^2(t) - \lambda_{13} m_{11}^2(t)); \\ \\ \frac{dm_{12}^2(t)}{dt} = (-\lambda_{21} m_2(t) - \lambda_{12} m_1(t) + (\alpha_1 m_{12}^2(t) - \delta_1 E[X_1^2 X_2] + \\ \lambda_{21} m_{22}^2(t) + \lambda_{31} m_{23}^2(t) - \lambda_{12} m_{12}^2(t) - \lambda_{13} m_{12}^2(t)) + \\ (\alpha_2 m_{12}^2(t) - \delta_2 E[X_1 X_2^2] - \lambda_{21} m_{12}^2(t) + \lambda_{12} m_{11}^2(t) + \\ \lambda_{3,2} m_{13}^2(t) - \lambda_{23} m_{12}^2(t)); \\ \\ \frac{dm_{13}^2(t)}{dt} = (-\lambda_{31} m_3(t) - \lambda_{13} m_1(t) + (\alpha_1 m_{13}^2(t) - \delta_1 E[X_1^2 X_3] + \\ \lambda_{21} m_{23}^2(t) + \lambda_{31} m_{33}^2(t) - \lambda_{12} m_{13}^2(t) - \lambda_{13} m_{13}^2(t)) + \\ (\alpha_3 m_{13}^2(t) - \delta_3 E[X_1 X_3^2] - \lambda_{31} m_{13}^2(t) - \lambda_{3,2} m_{13}^2(t) + \\ \lambda_{13} m_{11}^2(t) + \lambda_{23} m_{12}^2(t)); \\ \\ \frac{dm_{22}^2(t)}{dt} = (\alpha_2 m_2(t) + \delta_2 m_{22}^2(t) + \lambda_{21} m_2(t) + \lambda_{12} m_1(t) + \lambda_{3,2} m_3(t) + \\ \lambda_{23} m_2(t)) + 2(\alpha_2 m_{22}^2(t) - \delta_2 E[X_2^3] - \lambda_{21} m_{22}^2(t) + \\ \lambda_{12} m_{12}^2(t) + \lambda_{3,2} m_{23}^2(t) - \lambda_{23} m_{22}^2(t)); \\ \\ \frac{dm_{23}^2(t)}{dt} = (-\lambda_{3,2} m_3(t) - \lambda_{23} m_2(t) + (\alpha_2 m_{23}^2(t) - \delta_2 E[X_2^2 X_3] - \\ \lambda_{21} m_{23}^2(t) + \lambda_{12} m_{13}^2(t) + \lambda_{3,2} m_{33}^2(t) - \lambda_{23} m_{23}^2(t)) + \\ (\alpha_3 m_{23}^2(t) - \delta_3 E[X_2 X_3^2] - \lambda_{31} m_{23}^2(t) - \lambda_{3,2} m_{23}^2(t) + \\ \lambda_{13} m_{12}^2(t) + \lambda_{23} m_{22}^2(t)); \\ \\ \frac{dm_{33}^2(t)}{dt} = (\alpha_3 m_3(t) + \delta_3 m_{33}^2(t) + \lambda_{31} m_3(t) + \lambda_{3,2} m_3(t) + \lambda_{13} m_1(t) + \\ \lambda_{23} m_2(t)) + 2(\alpha_3 m_{33}^2(t) - \delta_3 E[X_3^3] - \lambda_{31} m_{33}^2(t) - \\ \lambda_{3,2} m_{33}^2(t) + \lambda_{13} m_{13}^2(t) + \lambda_{23} m_{23}^2(t)); \end{array} \right.$$

791 To remove the dependence of second-order moments on higher-order mo-
 792 ments, is possible to apply the moment closure schema introduced in sec-
 793 tion 3.1 and formulated in (3.4).

Supplement C: Reconstructing cell differentiation network
 (<http://www.e-publications.org/ims/support/download/imsart-ims.zip>). In order to investigate the structure of the differentiation tree, differentiation parameters λ are tested by means of the following asymptotic approximation

derived from the generalized method of moments theory (?),

$$(7.1) \quad \hat{\boldsymbol{\theta}} \sim \mathcal{N}_r(\boldsymbol{\theta}, \boldsymbol{\Sigma})$$

where $\hat{\boldsymbol{\theta}}$ is the final vector estimates returned by Algorithm 1 and the asymptotic covariance matrix $\boldsymbol{\Sigma}$ is a $r \times r$ matrix, estimated by means of

$$(7.2) \quad \hat{\boldsymbol{\Sigma}} = [\mathbf{J}(\hat{\boldsymbol{\theta}})^\top \mathbf{W}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{J}(\hat{\boldsymbol{\theta}})]^{-1}.$$

794 These distributional consideration are used to define Wald-type tests for the
795 differentiation parameters,

$$(7.3) \quad H_0 : \lambda_{ij} = 0$$

$$(7.4) \quad H_1 : \lambda_{ij} \neq 0.$$

796 In general, we reject H_0 and conclude that cell type i can differentiate into
797 cell type j , if $\hat{\lambda}_{ij} / \sqrt{\hat{\Sigma}_{\lambda_{ij}}} \geq z_\alpha$. To take into account the positivity constraint,
798 we consider a truncated normal distribution under H_0 as asymptotic distri-
799 bution, with mean zero and variance equal to the corresponding diagonal
800 element of $\hat{\boldsymbol{\Sigma}}$ and domain restricted to $[0, +\infty)$.

801 **Supplement D: Simulation study with 5 cell types.**

802 (<http://www.e-publications.org/ims/support/download/imsart-ims.zip>). In this
803 supplement, we describe the parameter setting used in the simulation study
804 of section 5.1 and shown in Figure 2a. We consider a cell differentiation net-
805 work with 5 cell types, and therefore 5 cell duplication parameters $\boldsymbol{\alpha}$, 5 cell
806 death parameters $\boldsymbol{\delta}$, as well as 5 cell differentiation parameters $\boldsymbol{\lambda}$:

$$\begin{aligned} \boldsymbol{\alpha} &= (1.0, 1.5, 1.8, 2.5, 2.8) \\ \boldsymbol{\delta} &= (0.033, 0.03, 0.045, 0.0312, 0.043) \\ \boldsymbol{\lambda} &= \begin{bmatrix} 0 & 0.2 & 0.35 & 0 & 0 \\ 0 & 0 & 0 & 0.75 & 0 \\ 0 & 0 & 0 & 0.25 & 0.5 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{aligned}$$

807 The Gillespie algorithm is implemented in C++ (Stroustrup, 1997) with the
808 support of **Eigen** library (Guennebaud et al., 2010). Our inferential proce-
809 dure, described in Algorithm 1, is implemented in **R** (R Core Team, 2015) by
810 means of custom scripts requiring Matrix packages for efficient dense and
811 sparse matrices manipulations Bates and Maechler (2015) and integrated
812 with C++ scripts calling **ODEint** (Ahnert and Mulansky, 2011) routines

813 that are available in the **Boost** library (Nakariakov, 2013). The quadratic
 814 programming problem is solved by means of **IBM ILOG CPLEX Op-**
 815 **imizer**, freely available under IBM Academic Initiative program (IBM,
 816 2010). All code used in this manuscript can be found in the online Sup-
 817 plement, and the latest version of the code is available at [github.com/](https://github.com/dp3111n/SLCDP_v1.0)
 818 [dp3111n/SLCDP_v1.0](https://github.com/dp3111n/SLCDP_v1.0).

Supplement E: Local linear approximation

(<http://www.e-publications.org/ims/support/download/imsart-ims.zip>). In this supplement, we describe a linear approximation of (4.1), which provides quick estimates for the parameters θ . This linear estimate is used in this paper in two different situations. First and foremost, it provides reasonable initial values for the exact non-linear algorithm described in Section 4.1. Secondly, it serves as a comparison in the evaluation of the proposed inference procedure for different sampling intervals. The linear approximation consists of calculating a computationally efficient, albeit approximate, solution for $m_i(s)$ and $m_{i,j}^2(s)$ in (??) by Euler's method,

$$\begin{aligned}
 m_i(s) &\simeq x_{i,s-1} + \sum_{k=1}^r v_{k,i} h_k(\mathbf{x}_{s-1}; \theta) \Delta t \\
 m_{i,j}^2(s) &\simeq x_{s-1,i} x_{s-1,j} + \sum_{k=1}^r v_{k,j} x_{s-1,i} h_k(\mathbf{x}_{s-1}; \theta) \Delta t \\
 &+ \sum_{k=1}^r v_{k,i} x_{s-1,j} h_k(\mathbf{x}_{s-1}; \theta) \Delta t + \sum_{k=1}^r v_{k,i} v_{k,j} h_k(\mathbf{x}_{s-1}; \theta) \Delta t
 \end{aligned}
 \tag{7.5}$$

Since (7.5) is linear in θ , the regression model in (4.1) can be conveniently reformulated as

$$d\mathbf{x}_{1:S} = \mathbf{M}_{0:S-1} \theta + \epsilon_{1:S}
 \tag{7.6}$$

where $d\mathbf{x}_{1:S} = \mathbf{x}_{1:S} - \mathbf{x}_{0:S-1}$ is column vector with observed cells counts differences between consecutive time points, $\mathbf{M}_{0:S-1} \theta = \mathbf{V}^\top D(\mathbf{x}_{0:S-1}) \Delta t \theta$ is a compact matrix equivalent of (7.5) with $D(\mathbf{x}_s)$ an $r \times r$ diagonal matrix with the appropriate polynomial of \mathbf{x}_s and $\text{Var}(\epsilon_s)$ component is estimated using $\Omega_{0:S-1} = \mathbf{V}^\top D(\mathbf{x}_{0:S-1}) \Delta t \text{Diag}(\theta) \mathbf{V}$. Analogously to (4.2) the local linear estimate $\tilde{\theta}$ are derived by means of an iterative procedure, in which the following constraint least squares problem is solved,

$$\tilde{\theta} = \arg \min_{\theta} (d\mathbf{x}_{1:S} - \mathbf{M}_{0:S-1} \theta)^\top (\Omega_{0:S-1})^{-1} (d\mathbf{x}_{1:S} - \mathbf{M}_{0:S-1} \theta) \text{ s. t. } \theta \geq \mathbf{0}_r.
 \tag{7.7}$$

819 The first estimate $\hat{\theta}^{(0)}$, used also as a starting point for the non-linear
 820 procedure, is calculated assuming homoscedastic and uncorrelated errors
 821 $\Omega_s = I_N$.

822 **Supplement F: Human hematopoiesis parameter estimates**
 823 (<http://www.e-publications.org/ims/support/download/imsart-ims.zip>). In the
 824 main paper, we compare a large number of models. For the selected model,
 we provide here the parameter estimates.

Parameter	Estimate	Parameter	Estimate
$\alpha_{CD34.HSC}$	8.006e+00	$\lambda_{CD34.HSC \rightarrow CD3.BM}$	0.721
$\alpha_{CD3.BM}$	7.930e-01	$\lambda_{CD34.HSC \rightarrow CD14.BM}$	0.867
$\alpha_{CD14.BM}$	9.187e-11	$\lambda_{CD34.HSC \rightarrow CD15.BM}$	0.591
$\alpha_{CD15.BM}$	3.900e-10	$\lambda_{CD34.HSC \rightarrow CD19.BM}$	1.453
$\alpha_{CD19.BM}$	1.251e-01	$\lambda_{CD34.HSC \rightarrow CD56.BM}$	0.146
$\alpha_{CD56.BM}$	2.852e-10	$\lambda_{CD34.HSC \rightarrow CD61.BM}$	0.335
$\alpha_{CD61.BM}$	5.737e-01	$\lambda_{CD34.HSC \rightarrow GLYCO.BM}$	0.713
$\alpha_{GLYCO.BM}$	1.643e-10	$\lambda_{CD3.BM \rightarrow CD3.PB}$	0.386
$\alpha_{CD3.PB}$	6.690e-11	$\lambda_{CD3.BM \rightarrow CD4.PB}$	0.180
$\alpha_{CD4.PB}$	9.959e-11	$\lambda_{CD3.BM \rightarrow CD8.PB}$	0.276
$\alpha_{CD8.PB}$	4.176e-11	$\lambda_{CD3.BM \rightarrow CD56.PB}$	0.151
$\alpha_{CD14.PB}$	1.006e-10	$\lambda_{CD14.BM \rightarrow CD56.PB}$	0.384
$\alpha_{CD15.PB}$	7.344e-11	$\lambda_{CD15.BM \rightarrow CD14.PB}$	0.207
$\alpha_{CD19.PB}$	3.763e-11	$\lambda_{CD15.BM \rightarrow CD15.PB}$	0.223
$\alpha_{CD56.PB}$	1.770e-10	$\lambda_{CD19.BM \rightarrow CD4.PB}$	0.149
$\delta_{CD34.HSC}$	2.393e-02	$\lambda_{CD19.BM \rightarrow CD19.PB}$	0.372
$\delta_{CD3.BM}$	1.735e-04	$\lambda_{CD19.BM \rightarrow CD56.PB}$	0.054
$\delta_{CD14.BM}$	2.308e-04	$\lambda_{CD56.BM \rightarrow CD56.PB}$	0.153
$\delta_{CD15.BM}$	2.510e-04	$\lambda_{CD61.BM \rightarrow CD15.PB}$	0.281
$\delta_{CD19.BM}$	2.322e-03	$\lambda_{CD61.BM \rightarrow CD56.PB}$	0.085
$\delta_{CD56.BM}$	6.831e-04	$\lambda_{GLYCO.BM \rightarrow CD14.PB}$	0.153
$\delta_{CD61.BM}$	4.734e-03		
$\delta_{GLYCO.BM}$	1.333e-03		
$\delta_{CD3.PB}$	1.417e-04		
$\delta_{CD4.PB}$	3.366e-04		
$\delta_{CD8.PB}$	2.601e-05		
$\delta_{CD14.PB}$	1.512e-03		
$\delta_{CD15.PB}$	5.115e-04		
$\delta_{CD19.PB}$	4.390e-04		
$\delta_{CD56.PB}$	2.630e-04		

TABLE 1

Parameter estimates for hematopoiesis in human, in-vivo, based on gene therapy clinical trial data, assuming an underlying stochastic cell differentiation process.

825

826 **Supplement G: Parameter estimates sensitivity to random ini-**
 827 **tialization**

828 (<http://www.e-publications.org/ims/support/download/imsart-ims.zip>). Here
 829 we show the sensitivity of the estimates to random initializations.

830 HARVARD MEDICAL SCHOOL,
 GENE THERAPY PROGRAM,
 DANA-FARBER/BOSTON CHILDREN'S
 CANCER AND BLOOD DISORDERS CENTER,
 BOSTON, MA, USA
 E-MAIL: daniello.pellin@dfci.harvard.edu

GOS INSTITUTE OF CHILD HEALTH
 UNIVERSITY COLLEGE LONDON
 GOWER STREET
 WC1E 6BT LONDON, UK
 E-MAIL: l.biasco@ucl.ac.uk

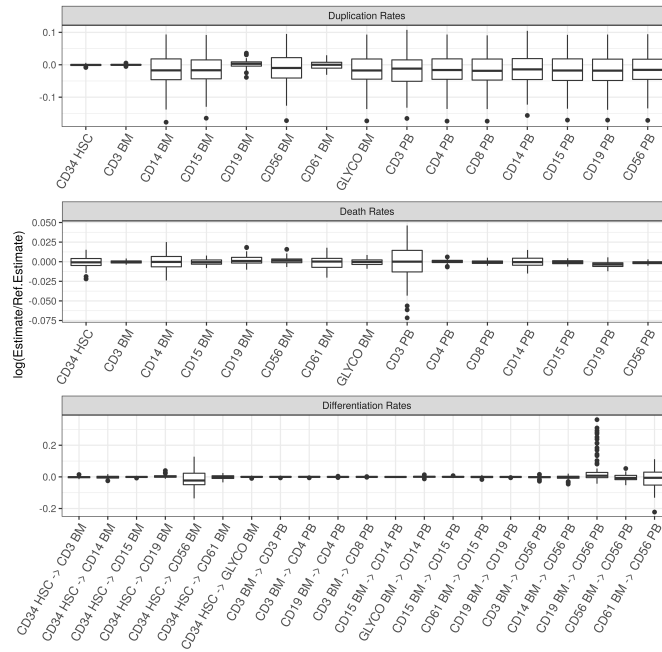


FIG 10. *Sensitivity analysis to random initialization.* Model \tilde{m} has been estimated starting from 100 different $\theta^{(0)}$ settings. Duplication and differentiation rates are sampled from a Normal distribution with $N(0.1, 0.1)$ and death rates from a $N(0.01, 0.01)$. Absolute value transformation was applied to avoid negative initial values. The distribution of logarithm of the ratio between the random restart estimates and the local linear initialization estimate (Ref. estimate, see values in Appendix Supplement F) is represented using boxplots.