



## CENTRE DE RECERCA MATEMÀTICA

This is a preprint of: *Oncometabolic nuclear reprogramming of cancer stemness*

Journal Information: *Stem Cell Reports*,

Author(s): J.A. Menendez, B. Corominas-Faja, E. Cuyas, M.G. García, S. Fernandez-Arroyo, A.F. Fernandez, J. Joven, M.F. Fraga, T. Alarcon.

Volume, pages: 6(3) 1-56, DOI:[10.1016/j.stemcr.2015.12.012]

## Oncometabolic Nuclear Reprogramming of Cancer Stemness

Javier A. Menendez,<sup>1,2,10,\*</sup> Bruna Corominas-Faja,<sup>2</sup> Elisabet Cuyàs,<sup>2</sup> María G. García,<sup>3</sup> Salvador Fernández-Arroyo,<sup>4</sup> Agustín F. Fernández,<sup>3</sup> Jorge Joven,<sup>4</sup> Mario F. Fraga,<sup>3,5</sup> and Tomás Alarcón<sup>6,7,8,9,11,\*</sup>

<sup>1</sup>ProCURE (Program Against Cancer Therapeutic Resistance), Metabolism and Cancer Group, Catalan Institute of Oncology, 17007 Girona, Catalonia, Spain

<sup>2</sup>Molecular Oncology Group, Girona Biomedical Research Institute (IDIBGI), 17190 Salt, Catalonia, Spain

<sup>3</sup>Cancer Epigenetics Laboratory, Instituto Universitario de Oncología del Principado de Asturias (IUOPA-HUCA), Universidad de Oviedo, 33006 Oviedo, Spain

<sup>4</sup>Unitat de Recerca Biomèdica, Hospital Universitari de Sant Joan, IISPV, Universitat Rovira i Virgili, Campus of International Excellence Southern Catalonia, 43201 Reus, Spain

<sup>5</sup>Nanomaterials and Nanotechnology Research Center (CINN-CSIC), 33940 San Martín del Rey Aurelio, Spain

<sup>6</sup>Institució Catalana d'Estudis i Recerca Avançats (ICREA), 08010 Barcelona, Spain

<sup>7</sup>Computational & Mathematical Biology Research Group, Centre de Recerca Matemàtica (CRM), 08193 Barcelona, Spain

<sup>8</sup>Departament de Matemàtiques, Universitat Autònoma de Barcelona, 08193 Barcelona, Spain

<sup>9</sup>Barcelona Graduate School of Mathematics (BGSMath), 08193 Barcelona, Spain

<sup>10</sup>Girona Biomedical Research Institute (IDIBGI), Parc Hospitalari Martí i Julià, Edifici M2, E-17190 Salt, Girona, Spain

<sup>11</sup>Centre de Recerca Matemàtica (CRM), Office 29 (C3b/140), Edifici C, Campus de Bellaterra, E-08193 Bellaterra, Barcelona, Spain

\*Correspondence: [jmenendez@iconcologia.net](mailto:jmenendez@iconcologia.net) (J.A.M.), [talarcon@crm.cat](mailto:talarcon@crm.cat) (T.A.)

<http://dx.doi.org/10.1016/j.stemcr.2015.12.012>

This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

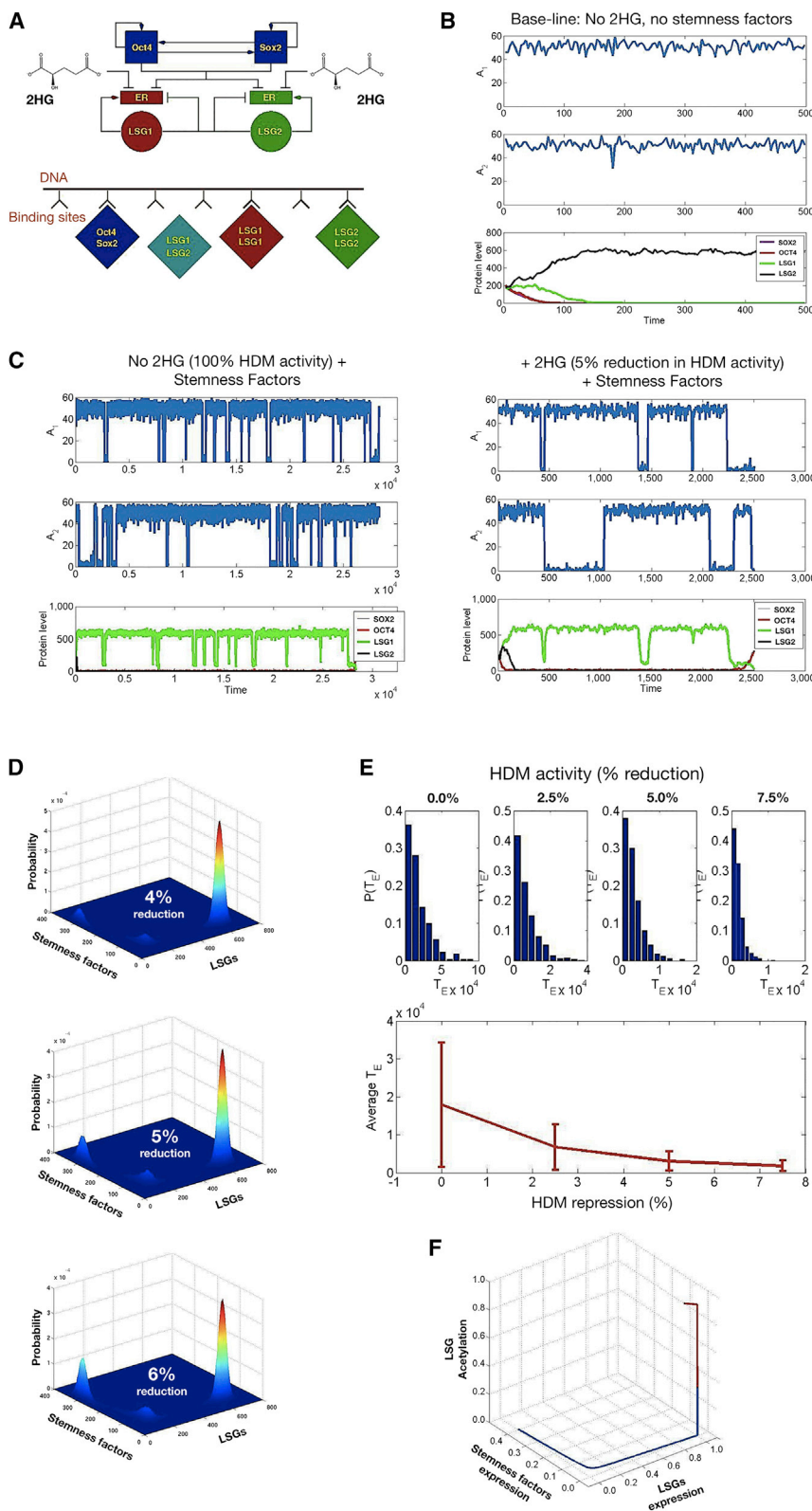
### SUMMARY

By impairing histone demethylation and locking cells into a reprogramming-prone state, oncometabolites can partially mimic the process of induced pluripotent stem cell generation. Using a systems biology approach, combining mathematical modeling, computation, and proof-of-concept studies with live cells, we found that an oncometabolite-driven pathological version of nuclear reprogramming increases the speed and efficiency of dedifferentiating committed epithelial cells into stem-like states with only a minimal core of stemness transcription factors. Our biomathematical model, which introduces nucleosome modification and epigenetic regulation of cell differentiation genes to account for the direct effects of oncometabolites on nuclear reprogramming, demonstrates that oncometabolites markedly lower the “energy barriers” separating non-stem and stem cell attractors, diminishes the average time of nuclear reprogramming, and increases the size of the basin of attraction of the macrostate occupied by stem cells. These findings establish the concept of oncometabolic nuclear reprogramming of stemness as a bona fide metabolo-epigenetic mechanism for generation of cancer stem-like cells.

### INTRODUCTION

The correct functioning of the epigenome ensures fidelity in the establishment of gene-expression programs that are compatible with specific cell identities. The need for tightly controlled epigenetic landscapes is of critical importance for stem cells, which are able to both self-renew and generate differentiated progeny (Barrero et al., 2010; Chen and Dent, 2014; Papp and Plath, 2013; Spivakov and Fisher, 2007). The inability to stabilize stem cell states and functions by maintaining epigenome integrity, a process in which DNA methylation plays a major role, can trigger pathological self-renewal processes that ultimately lead to cancer (Ohnishi et al., 2014; Suva et al., 2013; Tung and Knoepfler, 2015). Interestingly, remodeling of DNA methylation is a cancer-initiating event manifesting in the presence of particular types of cancer-driving metabolites, termed oncometabolites, and in the nuclear reprogramming process of transcription factor-generated induced pluripotent stem cell (iPSC) derivation.

The shared mechanism by which abnormal accumulation of the oncometabolites 2-hydroxyglutarate (2HG), succinate, and fumarate causes potential transformation to malignancy is the ability to promote DNA hypermethylation through suppression of histone demethylation, which, in turn, results in the repression of genes involved in the epigenetic rewiring of lineage-specific differentiation and in the promotion of stem cell-like transcriptional signatures (Chowdhury et al., 2011; Killian et al., 2013; Letouzé et al., 2013; Lu et al., 2012; Terunuma et al., 2014; Saha et al., 2014; Xiao et al., 2012; Xu et al., 2011; Yang et al., 2013). The transient expression of stemness-associated transcription factors, i.e., *OCT4*, *SOX2*, *KLF4*, and *c-MYC*, in vivo generates tumors consisting of undifferentiated dysplastic cells exhibiting global changes in DNA methylation (Ohnishi et al., 2014), suggesting that the epigenetic regulatory machinery associated with iPSC derivation might initiate cancer development in a manner that does not require mutational changes in the genomic sequence (Ben-David and Benvenisty, 2011; Ohnishi et al., 2014; Knoepfler, 2009; Tung and Knoepfler, 2015).



**Figure 1. Computation Simulation of Oncometabolic Nuclear Reprogramming Phenomena**

(A–C) A stochastic model of oncometabolic nuclear reprogramming. (A) Top: Schematic representation of the minimal gene regulatory network considered in our stochastic model, consisting of a coupled pluripotency module (self-activation of *Oct4* and *Sox2*) and a differentiation module (mutual antagonism between *LSGs*). Arrows denote activation and blunt-ended lines denote inhibitory interactions. Bottom: Schematic representation of the competitive binding model for activation/repression in the minimal gene regulatory network. (B) A realization path in which our stochastic model was run under baseline conditions (baseline HDM activity and lack of induction of stemness-related transcription factors, i.e.,  $h_i$ -values as per values given in Table S7 [Supplemental Appendix E] and  $p_1 = p_2 = 0$ ). Since the system is symmetric with respect to *LSG1* and *LSG2*, a state where  $O = 0$ ,  $S = 0$ , and  $L_2 = 0$ , whereas  $L_1 > 0$ , is also an absorbing state. (C) A realization path in which our stochastic model was run under induction of stemness-related transcription factors (parameter values  $p_1 = p_2 = 1.85 \times 10^7$ ). At the onset of stemness factor induction, i.e., we let  $p_1 > 0$  and  $p_2 > 0$ , the absorbing states observed in the simulations shown in (B) are not absorbing any longer and, therefore, there is a positive probability for the system to go from the differentiated cell state to the stem cell state. Left: Normal-like metabolism, baseline HDM activity; right: 2HG-induced reduction of HDM activity by 5% with respect to the baseline scenario.

(D–F) Epigenetic landscapes and reprogramming performance in response to 2HG. (D) 2HG-induced inhibition of HDM activity affects the depth of the stem cell attractors by lowering the barriers of the epigenetic landscape. Figures show the joint probability of the random variables  $O + S$  (stemness factors) and  $L_1 + L_2$  (LSGs) for different values of the relative oncometabolic-induced reduction of HDM activity with respect to the baseline scenario. To obtain the epigenetic landscapes for different degrees of 2HG-induced reduction of HDM activity in shorter computational time, we considered the following parameter values:  $p_1 = p_2 = 5.55 \times 10^{-7}$  and  $\vartheta_o = \vartheta_s =$

(legend continued on next page)



Because oncometabolites partially mimic the process of iPSC generation, a metabolically driven pathological version of nuclear reprogramming might represent an underappreciated epigenetic mechanism of enrichment for cellular states with increased tumor-initiating capacities and aberrant self-renewal potential (Goding et al., 2014; Menendez and Alarcón, 2014; Menendez et al., 2014b), often termed cancer stem cells. However, although a role for oncometabolite-driven changes in the epigenetic landscape is mechanistically attractive (Lu and Thompson, 2012; Yun et al., 2012; Johnson et al., 2015), the existence of bona fide oncometabolic reprogramming of differentiated cells into cancer stem-like states has never been demonstrated. In an attempt to resolve this issue, we have used a systems biology approach that combined mathematical modeling, computation, and proof-of-concept experimental validation of stochastic predictions in vitro.

## RESULTS

We initially developed methods and procedures for the mathematical modeling of oncometabolo-epigenetic regulatory networks involved in the acquisition of stemness (see [Supplemental Information](#)). Our stochastic model considers the interactions between a minimal core of stemness-associated transcription factors (*OCT4* and *SOX2*) and two generic lineage-specific genes (*LSG1* and *LSG2*) (Shu et al., 2013) ([Figure 1A](#)). The basic effector mechanism of the coupling between metabolism and the epi-transcriptional reprogramming system relies on histone- and nucleosome-modifying enzymes (Dodd et al., 2007). In particular, we considered a metabolo-epigenetic link in which the oncometabolite 2HG drastically inhibits the activity of DNA histone demethylases (HDMs), thus restricting the methylation plasticity that is required for the transition between stem cells and differentiated cells (Lu and Thompson, 2012; Lu et al., 2012).

The first consistency check we performed was that the “normal metabolism” scenario, defined as baseline HDM

activity and lack of induction of stemness transcription factors, should lead to cell differentiation. [Figure 1B](#) shows that, after an initial transient regime, the system settles down to a steady state whereby the protein levels of *OCT4*, *SOX2*, and *LSG1* decay to zero, whereas *LSG2* protein climbs to its stationary positive value. In the absence of induction, this cellular state is an absorbing state, i.e., once reached by the system it is not possible to exit. [Figure 1C](#) (left panel) shows a particular sample path where reprogramming of stemness is accomplished by merely adding *OCT4* and *SOX2* to the baseline scenario. Importantly, under baseline conditions both *LSG1* and *LSG2* are predominantly acetylated and, therefore, their promoters remain accessible to transcription factors, with short-lived journeys into the methylated state ([Supplemental Appendix D](#)). Despite the fact that during episodes of transient methylation the expression levels of the *LSGs* become downregulated, this does not necessarily lead to reprogramming since the stochastic dynamics of the gene regulatory network has to pass through an unstable saddle point ([Supplemental Appendix E](#)). Therefore, several episodes of transient methylation should occur before a period of transient methylation of sufficient duration allows the gene-regulation system to successfully pass through the bottleneck.

We then introduced the ability of 2HG to reduce HDM activity. [Figure 1C](#) (right panel) shows a realization of the stochastic model whereby reprogramming is achieved with a 5% reduction of HDM activity with respect to the baseline scenario ([Figure 1C](#), left panel). Oncometabolic reduction of HDM activity increases the characteristic duration of the transient episodes of methylation ([Supplemental Appendix D](#)) which, in turn, increases the likelihood of one such period of sufficient duration for the gene regulatory system to overcome the bottleneck. In this scenario, the system must transit from the differentiated state to the stem cell state. [Figure 1D](#) illustrates how the relative height of the peak corresponding to the stem cell state increases in relation to the peak corresponding to the differentiated cell state, thus implying that the

remaining parameter values are given in [Table S7](#) ([Supplemental Appendix E](#)). (E) 2HG-induced inhibition of HDM activity affects the kinetic efficiency of the reprogramming process. The panel shows statistics of the average reprogramming time,  $T_E$ , as well as its probability density,  $P(T_E)$ , as a function of the 2HG-induced reduction of HDM activity. The top panels illustrate that the predicted probability distribution of  $T_E$ ,  $P(T_E)$  is approximately exponential. The lower panel shows the average and SD (error bars) of the predicted reprogramming time, illustrating that the reprogramming rate increases exponentially with the 2HG-induced reduction of HDM activity. HDM activity reductions of 0%, 2.5%, 5%, and 7.5% correspond to  $h_2 = 1.00$ ,  $h_2 = 0.975$ ,  $h_2 = 0.95$ , and  $h_2 = 0.925$ , respectively. (F) 2HG-induced inhibition of HDM activity affects the size of the basin of attraction of the induced stem cell state. The graphic shows a solution of the semiclassical QSSA approximation (see Equations 29–32 and 40–44 in [Supplemental Information](#)) for the baseline scenario with no HDM inhibition ( $h_2 = 1.00$ , red line), and for the case with a 2HG-induced 5% reduction in HDM activity ( $h_2 = 0.95$ , blue line). The uninhibited scenario converges to the differentiated cell state (red line), whereas the inhibited scenario converges to the stem cell state (blue line). Parameter values  $\rho_1 = \rho_2 = 1.85 \times 10^7$ ,  $c_{E1} = c_{E2} = 2$ ,  $C_0 = C_5 = C_1 = C_2 = 1$ , and  $\vartheta_0 = \vartheta_5 = 0.2$ . The remaining parameter values are given in [Table S7](#) ([Supplemental Appendix E](#)).



epigenetic barriers are significantly lowered in response to 2HG-induced reduction of HDM activity. There is a “third peak” that arises from the fact that, while attempting reprogramming, the gene regulatory system spends a long time in the vicinity of the saddle point, trying to overcome the bottleneck. The height of this third peak appears to be rather insensitive to the 2HG-regulated activity of HDM, as it depends on the kinetic parameters of the gene regulatory system alone. [Figure 1E](#) shows the immediate and significant consequences for the kinetic efficiency of nuclear reprogramming; specifically, the reduction in average reprogramming time varies exponentially with the 2HG-induced reduction of HDM activity. Therefore, even modest 2HG-driven reductions of HDM activity are predicted to drive a considerable increase in the reprogramming efficiency.

We finally interrogated our stochastic model to examine whether the basins of attractions of each of these cell states, i.e., the set of developmental states which are attracted to each of them, are also altered in response to 2HG-induced reduction of HDM activity. We carried out a semiclassical quasi-steady-state approximation (QSSA) ([Supplemental Appendix B](#)) to analyze the existence of initial conditions that, in the absence of HDM inhibition, converge to a differentiated cell state and, upon inhibition, converge to the stem cell state. We found that such initial conditions exist, i.e., a portion of the basin of attraction of the differentiated cell state is transferred to the stem cell state upon HDM inhibition. [Figure 1F](#) illustrates how oncometabolic-induced repression of HDM activity actually enlarges the basin of attraction of the stem cell state.

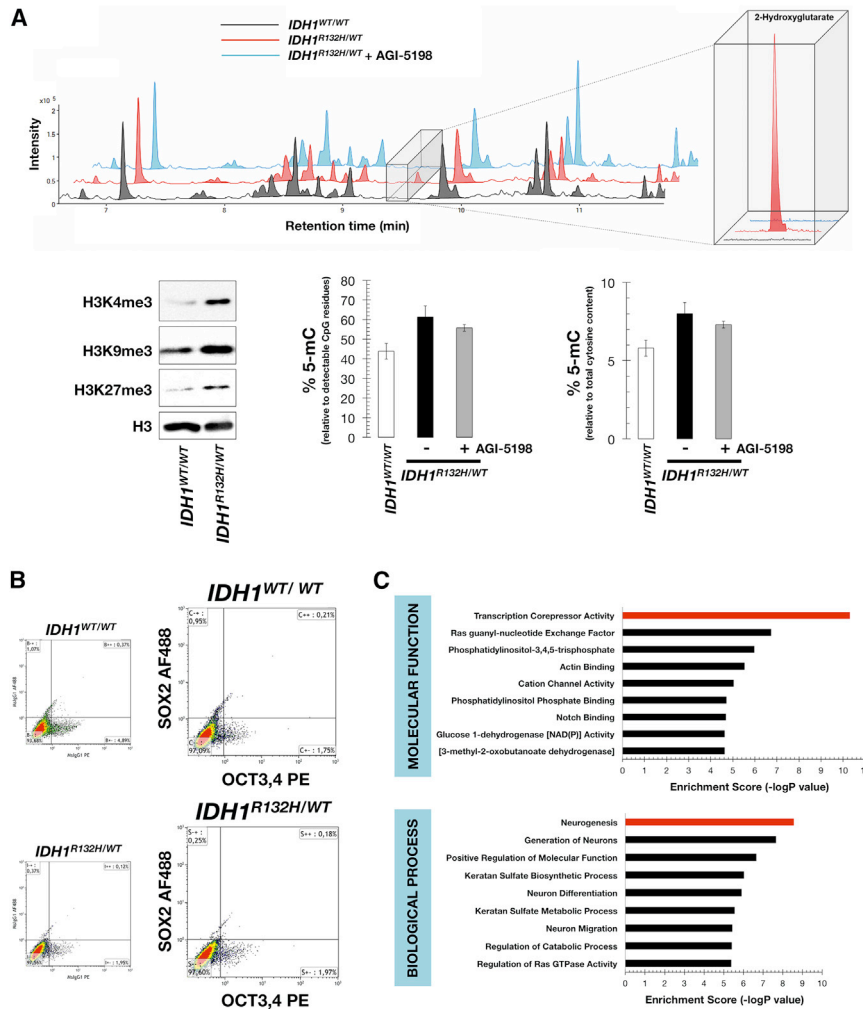
[Nishi et al. \(2014a, 2014b\)](#) have developed a method of inducing cancer stem-like cells (CSCs) through the reprogramming and partial differentiation of the immortalized but otherwise normal MCF10A human mammary epithelial cell line. To experimentally test our computational model of oncometabolic nuclear reprogramming, we similarly employed the MCF10A cell line and an isogenic derivative endogenously heterozygous for the *R132H* mutation of isocitrate dehydrogenase 1 (*IDH1*) gene, generating 2HG ([Grassian et al., 2012](#)). Quantification of intracellular 2HG showed that the levels of the oncometabolite were more than 30-fold higher in cell lysates from the knockin MCF10A *IDH1*<sup>R132H/WT</sup> cells, confirming neomorphic *IDH1*<sup>R132H</sup> enzymatic activity ([Figure 2A](#)). To corroborate that the sole accumulation of 2HG in an otherwise isogenic background was sufficient to significantly impair histone demethylation, we examined the pattern of histone lysine methylation in *IDH1*<sup>R132H/WT</sup> knockin and parental cells. Western blot analysis showed that the global levels of H3K4me<sub>3</sub>, H3K9me<sub>3</sub>, and H3K27me<sub>3</sub> were increased in 2HG-overproducing MCF10A knockin cells compared with *IDH1*<sup>WT/WT</sup> parental cells ([Figure 2A](#)). These results

were consistent with 2HG-induced broad inhibition of histone demethylation, showing agreement with previous models overexpressing *IDH* mutants ([Duncan et al., 2012; Lu et al., 2012](#)). We then employed a commercially available ELISA-based global DNA methylation able to indirectly provide a global measurement of 5-methylcytosine (5-mC) levels in genomic DNA obtained from *IDH1*<sup>WT/WT</sup> and *IDH1*<sup>R132H/WT</sup> knockin cells. Measurement of 5-mC levels of long interspersed nucleotide element 1 (*LINE-1*) confirmed that *LINE-1* methylation significantly increased in 2HG-overproducing *IDH1*<sup>R132H/WT</sup> cells when compared with *IDH1*<sup>WT/WT</sup> parental cells. Remarkably, a 2-day treatment with the selective R132H-*IDH1* inhibitor AGI-5198, which fully suppressed 2HG to background levels ([Figure 2A](#)), partially reverted *LINE-1* hypermethylation in *IDH1*<sup>R132H/WT</sup> cells.

Because it could be argued that 2HG-induced chromatin reorganization might promote the pluripotency-associated genes transition from an inactive to an active stage, we assessed whether the overproduction of 2HG promoted the expression of pluripotency regulators in normal breast epithelial cells. Flow cytometry analyses confirmed that the baseline expression of the core transcription factors OCT4 and SOX2 remained essentially unaltered in *IDH1*<sup>R132H/WT</sup> knockin cells compared with parental *IDH1*<sup>WT/WT</sup> cells ([Figure 2B](#)). A preliminary evaluation of the top ten most significant Gene Ontology (GO) “molecular function” and “biological process” term annotations overrepresented in the 290 differentially hypermethylated CpG sites that were identified in *IDH1*<sup>R132H/WT</sup> knockin cells ([Figure 2C](#)) strongly suggested that 2HG had a significant impact on the transcriptional repression of differentiation programs. Intracellular accumulation of the oncometabolite 2HG due to the heterozygous expression of the *IDH1*<sup>R132H</sup> allele is therefore sufficient to notably alter global histone lysine methylation without varying the baseline expression of the most critical reprogramming factors (i.e., OCT4 and SOX2) but altering expression of differentiation genes, thus providing an idoneous experimental model to validate the stochastic predictions of our biomathematical model in vitro.

MCF10A *IDH1*<sup>R132H/WT</sup> and MCF10A *IDH1*<sup>WT/WT</sup> cells were then transduced with OCT4 and SOX2 (hereafter called OS) to examine whether endogenously produced 2HG could substitute for combinations of stemness factors (i.e., *KLF4* and *c-MYC*) to reprogram MCF10A mammary cells into iPS-like (iPSL-10A) cells. At day 15 after infection, MCF10A *IDH1*<sup>R132H/WT</sup> cells growing on feeder layers showed a >10-fold increase in reprogramming efficiency relative to MCF10A *IDH1*<sup>WT/WT</sup> cells, as assessed by counting the number of alkaline phosphatase (AP)-positive (AP<sup>+</sup>) colonies ([Figure 3A](#)). We found that the colonies identified by the highly AP<sup>+</sup> criterion were also positive for strong





**Figure 2. Effects of 2HG on Histone Demethylation, Activation of Pluripotency Genes, and Genome-wide DNA Methylation**

(A) Top: Base-peak chromatograms of extracts from *IDH1*<sup>WT/WT</sup> cells (black line), *IDH1*<sup>R132H/WT</sup> cells (red line), and *IDH1*<sup>R132H/WT</sup> cells treated with 40  $\mu$ mol/l AGI-5198 for 2 days (blue line). A combined mass spectrum of the region where 2HG was eluted ( $m/z$  349.1317) is shown in the inset (two technical replicates per  $n$ ;  $n = 3$  biological replicates). Bottom: 2HG promotes broad inhibition of histone demethylation and LINE-1 global methylation. Left panel, western blots for total H3K4me3, H3K9me3, and H3K27me3 histone modifications in parental *IDH1*<sup>WT/WT</sup> and *IDH1*<sup>R132H/WT</sup> knockin cells. Also shown are total H3 controls (two technical replicates per  $n$ ;  $n = 2$  biological replicates). Middle and right panels, % of 5-mC relative to either detectable CpG residues or total cytosine content in *IDH1*<sup>WT/WT</sup> and *IDH1*<sup>R132H/WT</sup> cells, the latter being cultured in the absence or presence of 40  $\mu$ mol/l AGI-5198 for 2 days. The data are presented as the mean  $\pm$  SD (error bars); three technical replicates per  $n$ ;  $n = 2$  biological replicates.

(B) Flow cytometry analysis of OCT4/SOX2 expression in *IDH1*<sup>WT/WT</sup> and *IDH1*<sup>R132H/WT</sup> knockin cells. Representative dot plots showing the distribution of *IDH1*<sup>WT/WT</sup> and *IDH1*<sup>R132H/WT</sup> cells along the signal obtained with the isotype-specific control antibodies or with the OCT3,4 and SOX2 direct conjugated antibodies (two technical replicates per  $n$ ;  $n = 2$  biological replicates).

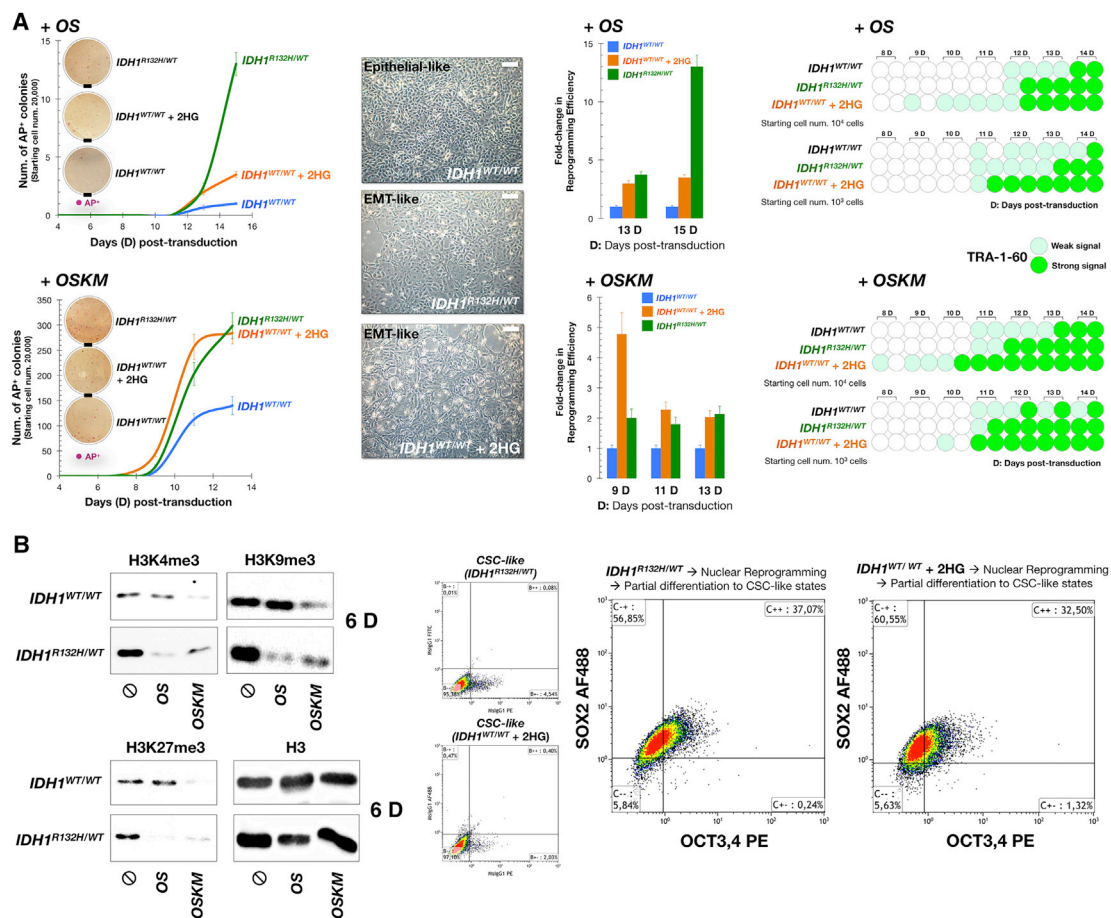
(C) Ten top-ranked GO molecular functions and biological processes associated with hypermethylated genes in *IDH1*<sup>R132H/WT</sup> cells (accession number GEO: GSE76263). x axis, negative logarithm (-lg) of the p value; y axis, GO category.

endogenous expression of NANOG (data not shown), which was considered a characteristic of bona fide iPSL-10A cells.

For live-cell imaging, we established a 96-well plate-based screening assay to assess expression of the pluripotency-associated surface marker TRA-1-60, a more reliable and specific marker for predicting successful reprogramming and iPS cell derivation than other markers including AP (Mali et al., 2010), during reprogramming. Cell clusters were scored as null, weak, or strong depending on the expression of TRA-1-60. Cell clusters positive for TRA-1-60 were detected in MCF10A *IDH1*<sup>R132H/WT</sup> cells 2–3 days earlier relative to MCF10A *IDH1*<sup>WT/WT</sup> cells (Figure 3A). When the standard reprogramming protocol was followed with four reprogramming factors, OCT4, SOX2, KLF4, and

*c-MYC* (hereafter OSKM), a significantly greater number of AP<sup>+</sup> colonies (~300) were observed in OSKM-transduced *IDH1*<sup>R132H/WT</sup> cells compared with *IDH1*<sup>WT/WT</sup> parental cells (~140) at day 13. TRA-1-60<sup>+</sup> clusters were also detected at an earlier stage in OSKM-transduced MCF10A *IDH1*<sup>R132H/WT</sup> cells than in MCF10A *IDH1*<sup>WT/WT</sup> control cells (Figure 3A).

Because the stimulatory effect of 2HG on the nuclear reprogramming efficiency was more striking in the absence of *KLF4* and *c-MYC* transgenes (>10-fold with OS) than in their presence (2-fold with OSKM), we preliminarily explored the time frame during which 2HG might exert its positive reprogramming effects. Exogenous supplementation of *IDH1*<sup>WT/WT</sup> parental cells with 1 mmol/l D-2HG octylester, a concentration of a cell membrane-permeable



**Figure 3. Effects of 2HG on the Nuclear Reprogramming of Breast Epithelial Cells into CSC-like States**

(A) Left panels: Kinetics of reprogramming in the absence or presence of 2HG. MCF10 *IDH1*<sup>WT/WT</sup> control cells and 2HG-overproducing MCF10A *IDH1*<sup>R132H/WT</sup> isogenic derivatives were reprogrammed by the retroviral delivery of OS (top) or OSKM (bottom) transcription factors. Alternatively, octyl-2HG, a cell-permeable esterified form of 2HG, was added at a final concentration of 1 mmol/l to the culture medium immediately after transduction of *IDH1*<sup>WT/WT</sup> parental cells with OS and OSKM, and was maintained for 4 days. The total number of highly AP<sup>+</sup> colonies for each condition was counted at different days until day 15 after transduction under feeder conditions. The data are presented as the mean ± SD (error bars); n = 3 biological replicates. Representative microphotographs of AP<sup>+</sup> colonies are also shown (scale bar, 5 mm). Middle panels: The reprogramming efficiencies of various conditions were compared with that obtained without octyl-2HG treatment in *IDH1*<sup>WT/WT</sup> parental cells, and are presented as relative fold changes (mean [columns] ± SD [error bars]). Insets show microscopy images of the representative cell morphology of *IDH1*<sup>WT/WT</sup>, *IDH1*<sup>R132H/WT</sup>, and *IDH1*<sup>WT/WT</sup> cells growing in the presence of octyl-2HG (scale bar, 10 μm). Right panels: Temporal activation of stemness during reprogramming was analyzed by live-cell staining with an antibody against TRA-1-60 (n = 3 biological replicates).

(B) Western blots for total H3K4me3, H3K9me3, and H3K27me3 histone modifications in parental *IDH1*<sup>WT/WT</sup> and *IDH1*<sup>R132H/WT</sup> knockin cells at day 6 post-OS or post-OSKM transduction. Also shown are total H3 controls (two technical replicates per n; n = 2 biological replicates). Right panels: Flow cytometry analysis of OCT4/SOX2 expression in CSC-like derivatives obtained from partial differentiation of reprogrammed *IDH1*<sup>WT/WT</sup> and *IDH1*<sup>R132H/WT</sup> knockin cells. Representative dot plots showing the distribution of *IDH1*<sup>WT/WT</sup> and *IDH1*<sup>R132H/WT</sup> cells along the signal obtained with the isotype-specific control antibodies or with the OCT3,4 and SOX2 direct conjugated antibodies (two technical replicates per n; n = 2 biological replicates).

form of 2HG that has previously been shown to mimic 2HG levels in tumors with aberrant 2HG accumulation by promoting a >100-fold increased intracellular concentration of 2HG (Lu et al., 2012; Xu et al., 2011; Terunuma et al., 2014), beginning soon after OS and OSKM transduc-

tion, for 4 days, resulted in a significantly increased number of reprogrammed colonies (Figure 3A). The fact that an early short-term supplementation with exogenous octyl-2HG, which caused prominent epithelial-to-mesenchymal (EMT)-like changes in cell fate (Figure 3A), was



sufficient to promote a pro-reprogramming effect similar to that of continued endogenous exposure in 2HG-overproducing EMT-like MCF10A *IDH1*<sup>R132H/WT</sup> cells (Grassian et al., 2012), was consistent with the notion that 2HG might contribute to the time-sensitive activation of a mesenchymal stage required during the initiation period of successful reprogramming (Liu et al., 2013; O'Malley et al., 2013).

Furthermore, increased numbers of AP<sup>+</sup> colonies and TRA-1-60<sup>+</sup> clusters were detected within a shorter period of time in octyl-2HG-treated *IDH1*<sup>WT/WT</sup> parental cells (Figure 3A). Interestingly, when monitoring the effect of 2HG on the appearance of tightly packed colonies morphologically resembling human embryonic stem cells (hESCs), 2HG was found to increase the ratio of AP<sup>+</sup> iPSL-10A colonies to total hESC-like colonies, i.e., the 2HG-driven increase in the number of AP<sup>+</sup> iPSL-10A colonies was not accompanied by changes in the total colony number, thus implying that 2HG enhances the destination of reprogrammed cells to the stem cell fate (data not shown). To confirm that the pre-existing status of histone modifications might differentially regulate the global chromatin environment controlling reprogramming toward a pluripotent state, we reexamined the global histone modification variation that occurred in early stages of reprogramming (i.e., 6 days after OS and OSKM transduction) before genuine AP<sup>+</sup> iPS-like patches become apparent in the cultures (Figure 3B). Interestingly, H3K9 methylation, which has been defined as the primary epigenetic determinant for the intermediate pre-iPS state as its removal leads to fully reprogrammed iPS cells (Chen et al., 2013), was notably suppressed by early OS induction in 2HG-overproducing *IDH1*<sup>R132H/WT</sup> cells but not in 2HG-negative *IDH1*<sup>WT/WT</sup> cells. Moreover, a reduction of H3K27me3 and H3K4me3, a phenomenon that has been associated with the acquisition of a transient open/primed chromatin state during the early transcriptional events of nuclear reprogramming (Hussein et al., 2014), was apparently observed upon early OS induction in 2HG-overproducing *IDH1*<sup>R132H/WT</sup> cells but not in 2HG-negative *IDH1*<sup>WT/WT</sup> cells. Thus, whereas the baseline levels of functionally opposing histone methylation marks were increased in 2HG-overproducing *IDH1*<sup>R132H/WT</sup> cells, consistent with a broad inhibition of histone demethylation, 2HG-induced histone modifications provided a collaborative rewiring of the chromatin organization that responded more rapidly and efficiently to the core stemness factors (OS) at the start of reprogramming.

Because the process of dedifferentiation through the addition of Yamanaka factors is extremely inefficient, the actual contribution of de novo generated “CSC states” via oncometabolic reprogramming to cancer evolution might be a matter of conjecture. To evaluate such a situation, we

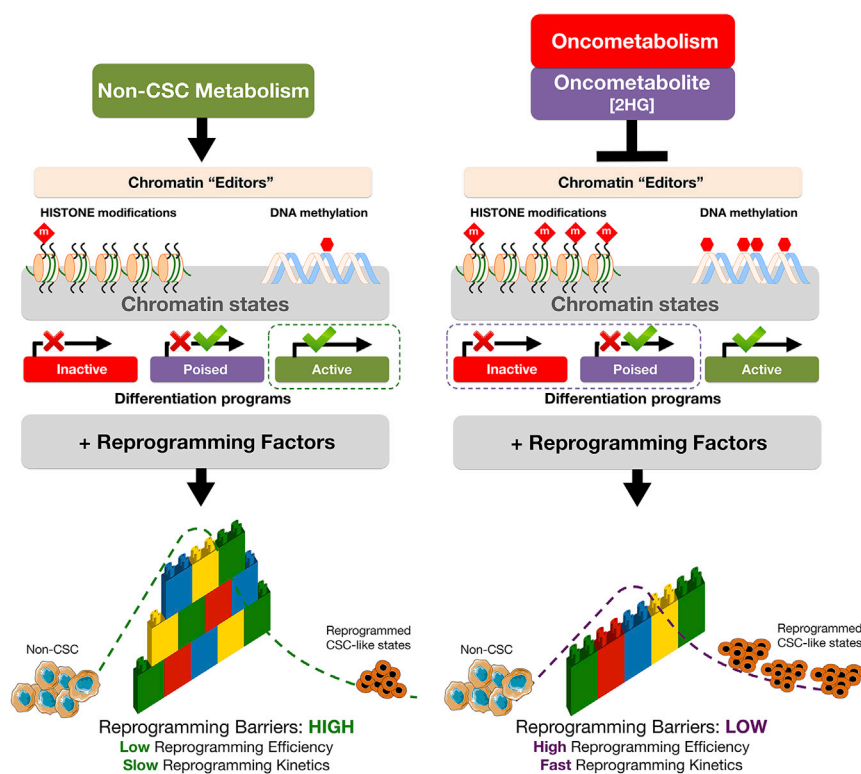
first confirmed that, upon the introduction of defined reprogramming factors and subsequent partial differentiation, proliferating CSC-like cells stably overexpressing OCT4 and SOX2 with tumor-initiating capacity (Nishi et al., 2014b) likewise arise from non-CSC, OCT4/SOX2-negative MCF10A *IDH1*<sup>R132H/WT</sup> cells (Figure 3B). We then designed a mathematical model to investigate the expected dynamics of tumor progression when the presence of oncometabolic signals can favor differentiated cells to revert to a multipotent CSC-like state. When a native, “resident” population sustained by normal stem cells competes with “invader” clones of CSC-like cells generated de novo as the result of nuclear reprogramming, our mathematical model predicts that the chances of prolonged survival increase exponentially with the size of the reprogrammed clones (Supplemental Appendix F). By solely affecting epigenetic events involving histone methylation, oncometabolites such as 2HG can functionally replace stemness transcription factors (e.g., *KLF4* and *c-MYC*) and accelerate the dedifferentiation rates to efficiently drive the de novo generation of reprogrammed CSC-like states, thus confirming that the possibility that metabolically driven nuclear reprogramming-like phenomena contribute to cancer initiation, and that progression cannot be neglected in terms of cancer prognosis and therapeutic planning (Brooks et al., 2015; Leder et al., 2010; Martin-Castillo et al., 2015; Menendez et al., 2014a).

The experimental approach using no-2HG versus 2HG-overproducing cellular models in an identical non-transformed genomic background, functionally confirms the predictions of our stochastic model, demonstrating that an oncometabolite markedly lowers the “energy barriers” separating non-stem and stem cell attractors, diminishes the average time of reprogramming, and increases the size of the basin of attraction of the macrostate occupied by stem cells (Figure 4). For 2HG to improve nuclear reprogramming performance, it is sufficient to be present only during the first few days of reprogramming, when it appears to exert partial functional redundancy with other reprogramming factors that ensure the supply of chromatin-modifying enzymes with metabolic intermediates for the epigenetic activation of stemness-related gene networks (Goding et al., 2014; Gut and Verdin, 2013; Menendez and Alarcón, 2014).

## DISCUSSION

One of the most challenging issues in the field of cancer research is understanding how cellular metabolism influences chromatin structure and the epigenome to drive tumor formation (Johnson et al., 2015; Menendez and Alarcón, 2014; Lu and Thompson, 2012; Yun et al.,





**Figure 4. Oncometabolite-Driven Nuclear Reprogramming of Cancer Stemness: A Framework Proposal**

HDMs, such as Jumonji histone demethylases (JHDM) and ten-eleven translocation (TET) family members, remove repressive histone methylation marks and activate the expression of differentiation-related genes by protecting promoters from aberrant DNA methylation. Oncometabolites such as 2HG inhibit the epigenetic “editors” HDMs and TETs, which leads to histone modifications (e.g., increased H3K9me3, H3K27me3, and H3K4me3) and DNA hypermethylation. Oncometabolites reprogram chromatin state to promote the downregulation of genes involved in differentiation as well as bias in developmental gene-expression patterns. This metabolo-epigenetic modification of inactive/poised states of lineage-specific genes is sufficient to significantly alter the efficiency and speed of nuclear reprogramming by lowering the “reprogramming barriers” of the epigenetic landscape and increasing the size of the stem cell state basin of attraction, which results in the acceleration (i.e., higher efficiency and faster kinetics) of the nuclear reprogramming

process. Oncometabolites such as 2HG permissively alleviate the unfavorable developmental process of “jumping” from differentiated cell states to CSC-like attractors while concomitantly stabilizing the ground-state self-maintaining character of CSC states. This conceptual figure represents cells stabilized in an initial non-CSC attractor and how nuclear reprogramming can make cells exceed the “reprogramming barrier,” represented as a wall of interlocking bricks, harder or easier in the absence or presence of the oncometabolite 2HG, respectively, and fall down in a final CSC attractor. The cellular reprogramming process is represented as a dashed line from the initial to the final cellular state.

2012). To date, however, there have been no attempts to delineate predictive mathematical platforms that operatively integrate the required contribution of certain metabolites for the extensive remodeling of the epigenetic landscape that drives nuclear reprogramming (Morris et al., 2014). From a mathematical standpoint, here we introduce nucleosome modification and epigenetic regulation of lineage-specific genes as an essential element of stochastic modeling that successfully integrate the recognized ability of oncometabolites to competitively inhibit epigenetic regulation of cell differentiation with the process whereby the stemness regulatory circuitry is established during nuclear reprogramming (Ben-David et al., 2013; Shu et al., 2013). By combining mathematical modeling and computation simulation with wet-lab in vitro experiments in an isogenic model, we demonstrate the existence of bona fide oncometabolic nuclear reprogramming phenomena able to efficiently generate CSC-like states (Figure 4). Our model provides a stochastic tool as well as a conceptual framework that should be extremely useful in helping to understand and investigate

the underexplored link between cellular metabolism and cancer-driving alterations in the epigenome. Beyond the numerous “common” metabolites that are used as substrates and cofactors for reactions that coordinate epigenetic status (Locasale, 2013; Johnson et al., 2015; Yun et al., 2012), a recent systems approach predicted more than 40 compounds and substructures of potential “oncometabolites” that could result from the loss-of-function and gain-of-function mutations of metabolic enzymes (Nam et al., 2014). In this context, our model can be a starting point for future studies on the processes by which cellular metabolism influences chromatin structure and epi-transcriptional circuits to causally drive stemness in cancer tissues.

## EXPERIMENTAL PROCEDURES

### Stochastic Model

A detailed mathematical formulation of the stochastic model of oncometabolic nuclear reprogramming can be found in the [Supplemental Information](#).



## Reagents

The *octyl ester* derivative of [2R]-2-hydroxyglutaric acid was purchased from US Biologicals Life Sciences (cat. #01386; Deltaclon). AGI-5198, a highly potent and selective inhibitor of IDH1 R132H/R132C mutants, was purchased from Selleck Chemicals (cat. #S7185).

## Cell Lines

MCF10A cells with heterozygous knockin of *IDH1* dominant-negative (*R132H*) point mutation and MCF10A isogenic parental cells were obtained from Horizon Discovery (cat. #HD 101-013 and #HD PAR-058, respectively). *IDH1* mutational status was verified by sequencing (Grassian et al., 2012).

## Reprogramming of Human Breast Epithelial Cells and Cell Infection

MCF10A *IDH1*<sup>R132H/WT</sup> and MCF10A *IDH1*<sup>WT/WT</sup> cells were transduced with retroviral vectors encoding nuclear reprogramming factors as previously described (Nishi et al., 2014a, 2014b) (Figure S1). The pMX vectors containing human cDNA for *OCT4*, *SOX2*, *KLF4*, and *c-MYC* were obtained from Addgene (<http://www.addgene.org>).

## Alkaline Phosphatase Staining

AP staining was performed using the Leukocyte Alkaline Phosphatase kit (cat. #86-R; Sigma-Aldrich) according to the manufacturer's protocol in OS- and OSKM-transduced cells reseeded onto mouse embryo fibroblast feeder layers.

## Live Staining by the TRA-1-60 Antibody

A mouse anti-human StainAlive TRA-1-60 antibody (DyLight 488; cat. #09-0068) was employed according to the manufacturer's protocol to identify and track the appearance of iPS-like colonies in OS- and OSKM-transduced cells reseeded onto Matrigel-coated 96-well plates.

## Histone Extraction and Western Analysis

Histones were acid-extracted following a modified version of the original protocol published by Sarg et al. (2002). For western blot analyses of H3K4me3, H3K9me3, H3K27me3, and total H3, 12 µg of the histone lysates were electrophoresed on 17% SDS-PAGE gel, transferred to a 0.45-mm polyvinylidene fluoride membrane, and incubated with antibodies against histone H3 (Abcam; cat. #ab1791), H3K4me3 (Abcam; cat. #ab8580), H3K9me3 (Millipore; cat. #CS200604), and H3K27me3 (Upstate-Millipore; cat. #07-449, lot #DAM1421462), followed by horseradish peroxidase-conjugated secondary and chemiluminescence detection.

## Targeted Metabolomics

Quantitative measurements of 2HG were performed by employing a method based on gas chromatography coupled to a quadrupole time-of-flight mass spectrometer and an electron impact interface (GC-EI-QTOF-MS). A detailed description of this procedure is given in Cuyàs et al. (2015) and Riera-Borrull et al. (2016).

## Global DNA Methylation

DNA was extracted and purified with a DNeasy Blood & Tissue kit (Qiagen; cat. #69504 or #69506) according to the manufacturer's instructions. Global DNA methylation levels were determined using the Global DNA Methylation LINE-1 kit (Active Motif; cat. #55017) according to the manufacturer's instructions.

## Genome-wide DNA Methylation: DNA Methylation Microarrays and Data Analysis

Microarray-based DNA methylation profiling was performed using Illumina Infinium HumanMethylation450 BeadChip Array (Bibikova et al., 2011). Methylation levels (beta values) were obtained using Illumina's GenomeStudio Software. The beta value represents a quantitative measure of the DNA methylation level of specific CpG sites and ranges from 0 (completely unmethylated) to 1 (completely methylated). Before analyzing the methylation data (accession number GEO: GSE76263), we excluded possible sources of technical biases that could alter the results. We excluded probes with a detection p value of  $\geq 0.01$  and removed the probes containing a SNP at the CpG interrogation site. Because there were only two samples in the experimental design, we used stringent statistical criteria to define differential methylated probes. Thus, we defined a probe to be hypermethylated or hypomethylated if the differences between beta values was larger than 0.5. The HOMER (Hypergeometric Optimization of Motif EnRichment) suite of tools (<http://homer.salk.edu/homer/>) was used to determine the enrichment of individual ontology terms and create GO maps in the groups of differentially methylated genes.

## Multivariate Permeabilized Cell Flow Cytometry

OCT3,4 and SOX2 protein levels were analyzed by intracellular staining using the Fix & Perm Cell Permeabilization kit (Invitrogen, cat. #GAS004), and flow cytometry using the anti-hOct4-PE (Becton Dickinson, cat. #560186) and anti-hSox2-AF488 (BD, cat. #561593) primary antibodies. Corresponding isotype antibodies MsIgG1 PE (BD, cat. #556650) and MsIgG1 AF488 (BD, cat. #551954) were used as controls. Plots show the fluorescence intensity distribution and percentages of cells above or below the thresholds determined by the staining with the isotype antibodies.

## ACCESSION NUMBERS

The accession number for the DNA methylation data reported in this paper is GEO: GSE76263.

## SUPPLEMENTAL INFORMATION

Supplemental Information includes Supplemental Experimental Procedures, Supplemental data, Supplemental appendices, three figures, and eight tables and can be found with this article online at <http://dx.doi.org/10.1016/j.stemcr.2015.12.012>.

## AUTHOR CONTRIBUTIONS

J.A.M. and T.A. conceived the idea for this project and wrote the manuscript. T.A. developed the stochastic mathematical model. J.A.M. designed and analyzed the experiments on living cells. B.C.F. and E.C. conducted reprogramming experiments and data



analysis. M.G.G. conducted the western blot experiments and data analysis. A.E.F. and M.F.F. conducted DNA methylation arrays and data analysis.

## ACKNOWLEDGMENTS

This work was supported by grants from the Ministerio de Ciencia e Innovación (Grant SAF2012-38914), Plan Nacional de I + D + I, Spain and the Agència de Gestió d'Ajuts Universitaris i de Recerca (AGAUR) (Grant 2014 SGR229), Departament d'Economia i Coneixement, Catalonia, Spain to J.A.M. T.A. acknowledges financial support from the Ministerio de Ciencia e Innovación (MICINN) under grant MTM2011-29342 and Agència de Gestió d'Ajuts Universitaris i de Recerca (AGAUR) under grant 2014SGR1307. We are thankful to José M.A. Vaquero (Center of Regenerative Medicine, CMR[B] Core Facility-Cytometry Unit, Barcelona, Spain) for excellent assistance with flow cytometry.

Received: May 18, 2015

Revised: December 30, 2015

Accepted: December 31, 2015

Published: February 11, 2016

## REFERENCES

- Barrero, M.J., Boué, S., and Izpisua Belmonte, J.C. (2010). Epigenetic mechanisms that regulate cell identity. *Cell Stem Cell* 7, 565–570.
- Ben-David, U., and Benvenisty, N. (2011). The tumorigenicity of human embryonic and induced pluripotent stem cells. *Nat. Rev. Cancer* 11, 268–277.
- Ben-David, U., Nissenbaum, J., and Benvenisty, N. (2013). New balance in pluripotency: reprogramming with lineage specifiers. *Cell* 153, 939–940.
- Bibikova, M., Barnes, B., Tsan, C., Ho, V., Klotzle, B., Le, J.M., Delano, D., Zhang, L., Schroth, G.P., Gunderson, K.L., et al. (2011). High density DNA methylation array with single CpG site resolution. *Genomics* 98, 288–295.
- Brooks, M.D., Burness, M.L., and Wicha, M.S. (2015). Therapeutic implications of cellular heterogeneity and plasticity in breast cancer. *Cell Stem Cell* 17, 260–271.
- Chen, T., and Dent, S.Y. (2014). Chromatin modifiers and remodelers: regulators of cellular differentiation. *Nat. Rev. Genet.* 15, 93–106.
- Chen, J., Liu, H., Liu, J., Qi, J., Wei, B., Yang, J., Liang, H., Chen, Y., Chen, J., Wu, Y., et al. (2013). H3K9 methylation is a barrier during somatic cell reprogramming into iPSCs. *Nat. Genet.* 45, 34–42.
- Chowdhury, R., Yeoh, K.K., Tian, Y.M., Hillringhaus, L., Bagg, E.A., Rose, N.R., Leung, I.K., Li, X.S., Woon, E.C., Yang, M., et al. (2011). The oncometabolite 2-hydroxyglutarate inhibits histone lysine demethylases. *EMBO Rep.* 12, 463–469.
- Cuyàs, E., Fernández-Arroyo, S., Corominas-Faja, B., Rodríguez-Gallego, E., Bosch-Barrera, J., Martín-Castillo, B., De Llorens, R., Joven, J., and Menendez, J.A. (2015). Oncometabolic mutation IDH1 R132H confers a metformin-hypersensitive phenotype. *Oncotarget* 6, 12279–12296.
- Dodd, I.B., Micheelsen, M.A., Sneppen, K., and Thon, G. (2007). Theoretical analysis of epigenetic cell memory by nucleosome modification. *Cell* 129, 813–822.
- Duncan, C.G., Barwick, B.G., Jin, G., Rago, C., Kapoor-Vazirani, P., Powell, D.R., Chi, J.T., Bigner, D.D., Vertino, P.M., and Yan, H. (2012). A heterozygous IDH1R132H/WT mutation induces genome-wide alterations in DNA methylation. *Genome Res.* 22, 2339–2355.
- Goding, C.R., Pei, D., and Lu, X. (2014). Cancer: pathological nuclear reprogramming? *Nat. Rev. Cancer* 14, 568–573.
- Grassian, A.R., Lin, F., Barrett, R., Liu, Y., Jiang, W., Korpai, M., Astley, H., Gitterman, D., Henley, T., Howes, R., et al. (2012). Isocitrate dehydrogenase (IDH) mutations promote a reversible ZEB1/microRNA (miR)-200-dependent epithelial-mesenchymal transition (EMT). *J. Biol. Chem.* 287, 42180–42194.
- Gut, P., and Verdin, E. (2013). The nexus of chromatin regulation and intermediary metabolism. *Nature* 502, 489–498.
- Hussein, S.M., Puri, M.C., Tonge, P.D., Benevento, M., Corso, A.J., Clancy, J.L., Mosbergen, R., Li, M., Lee, D.S., Cloonan, N., et al. (2014). Genome-wide characterization of the routes to pluripotency. *Nature* 516, 198–206.
- Johnson, C., Warmoes, M.O., Shen, X., and Locasale, J.W. (2015). Epigenetics and cancer metabolism. *Cancer Lett.* 256, 309–314.
- Killian, J.K., Kim, S.Y., Miettinen, M., Smith, C., Merino, M., Tsoikos, M., Quezada, M., Smith, W.L., Jr., Jahromi, M.S., Xekouki, P., et al. (2013). Succinate dehydrogenase mutation underlies global epigenomic divergence in gastrointestinal stromal tumor. *Cancer Discov.* 3, 648–657.
- Knoepfler, P.S. (2009). Deconstructing stem cell tumorigenicity: a roadmap to safe regenerative medicine. *Stem Cells* 27, 1050–1056.
- Leder, K., Holland, E.C., and Michor, F. (2010). The therapeutic implications of plasticity of the cancer stem cell phenotype. *PLoS One* 5, e14366.
- Letouzé, E., Martinelli, C., Lorient, C., Burnichon, N., Abermil, N., Ottolenghi, C., Janin, M., Menara, M., Nguyen, A.T., Benit, P., et al. (2013). SDH mutations establish a hypermethylator phenotype in paraganglioma. *Cancer Cell* 23, 739–752.
- Liu, X., Sun, H., Qi, J., Wang, L., He, S., Liu, J., Feng, C., Chen, C., Li, W., Guo, Y., et al. (2013). Sequential introduction of reprogramming factors reveals a time-sensitive requirement for individual factors and a sequential EMT-MET mechanism for optimal reprogramming. *Nat. Cell Biol.* 15, 829–838.
- Locasale, J.W. (2013). Serine, glycine and one-carbon units: cancer metabolism in full circle. *Nat. Rev. Cancer* 13, 572–583.
- Lu, C., and Thompson, C.B. (2012). Metabolic regulation of epigenetics. *Cell Metab.* 16, 9–17.
- Lu, C., Ward, P.S., Kapoor, G.S., Rohle, D., Turcan, S., Abdel-Wahab, O., Edwards, C.R., Khanin, R., Figueroa, M.E., Melnick, A., et al. (2012). IDH mutation impairs histone demethylation and results in a block to cell differentiation. *Nature* 483, 474–478.
- Mali, P., Ye, Z., Chou, B.K., Yen, J., and Cheng, L. (2010). An improved method for generating and identifying human induced pluripotent stem cells. *Methods Mol. Biol.* 636, 191–205.



- Martin-Castillo, B., Lopez-Bonet, E., Cuyàs, E., Viñas, G., Pernas, S., Dorca, J., and Menendez, J.A. (2015). Cancer stem cell-driven efficacy of trastuzumab (Herceptin): towards a reclassification of clinically HER2-positive breast carcinomas. *Oncotarget* 6, 32317–32338.
- Menendez, J.A., and Alarcón, T. (2014). Metabostemness: a new cancer hallmark. *Front. Oncol.* 4, 262.
- Menendez, J.A., Alarcón, T., Corominas-Faja, B., Cuyàs, E., López-Bonet, E., Martin, A.G., and Vellon, L. (2014a). Xenopatients 2.0: reprogramming the epigenetic landscapes of patient-derived cancer genomes. *Cell Cycle* 13, 358–370.
- Menendez, J.A., Corominas-Faja, B., Cuyàs, E., and Alarcón, T. (2014b). Metabostemness: metaboloepigenetic reprogramming of cancer stem-cell functions. *Oncoscience* 1, 803–806.
- Morris, R., Sancho-Martinez, I., Sharpee, T.O., and Izpisua Belmonte, J.C. (2014). Mathematical approaches to modeling development and reprogramming. *Proc. Natl. Acad. Sci. USA* 111, 5076–5082.
- Nam, H., Campodonico, M., Bordbar, A., Hyduke, D.R., Kim, S., Zielinski, D.C., and Palsson, B.O. (2014). A systems approach to predict oncometabolites via context-specific genome-scale metabolic networks. *PLoS Comput. Biol.* 10, e1003837.
- Nishi, M., Akutsu, H., Kudoh, A., Kimura, H., Yamamoto, N., Umezawa, A., Lee, S.W., and Ryo, A. (2014a). Induced cancer stem-like cells as a model for biological screening and discovery of agents targeting phenotypic traits of cancer stem cell. *Oncotarget* 5, 8665–8680.
- Nishi, M., Sakai, Y., Akutsu, H., Nagashima, Y., Quinn, G., Masui, S., Kimura, H., Perrem, K., Umezawa, A., Yamamoto, N., et al. (2014b). Induction of cells with cancer stem cell properties from nontumorigenic human mammary epithelial cells by defined reprogramming factors. *Oncogene* 33, 643–652.
- Ohnishi, K., Semi, K., Yamamoto, T., Shimizu, M., Tanaka, A., Mitsunaga, K., Okita, K., Osafune, K., Arioka, Y., Maeda, T., et al. (2014). Premature termination of reprogramming in vivo leads to cancer development through altered epigenetic regulation. *Cell* 156, 663–677.
- O'Malley, J., Skylaki, S., Iwabuchi, K.A., Chantzoura, E., Ruetz, T., Johnsson, A., Tomlinson, S.R., Linnarsson, S., and Kaji, K. (2013). High-resolution analysis with novel cell-surface markers identifies routes to iPS cells. *Nature* 499, 88–91.
- Papp, B., and Plath, K. (2013). Epigenetics of reprogramming to induced pluripotency. *Cell* 152, 1324–1343.
- Riera-Borrull, M., Rodríguez-Gallego, E., Hernández-Aguilera, A., Luciano, F., Ras, R., Cuyàs, E., Camps, J., Segura-Carretero, A., Menendez, J.A., Joven, J., and Fernández-Arroyo, S. (2016). Exploring the process of energy generation in pathophysiology by targeted metabolomics: performance of a simple and quantitative method. *J. Am. Soc. Mass Spectrom.* 27, 168–177.
- Saha, S.K., Parachoniak, C.A., Ghanta, K.S., Fitamant, J., Ross, K.N., Najem, M.S., Gurumurthy, S., Akbay, E.A., Sia, D., Cornella, H., et al. (2014). Mutant IDH inhibits HNF-4 $\alpha$  to block hepatocyte differentiation and promote biliary cancer. *Nature* 513, 110–114.
- Sarg, B., Koutzamani, E., Helliger, W., Rundquist, I., and Lindner, H.H. (2002). Postsynthetic trimethylation of histone H4 at lysine 20 in mammalian tissues is associated with aging. *J. Biol. Chem.* 277, 39195–39201.
- Shu, J., Wu, C., Wu, Y., Li, Z., Shao, S., Zhao, W., Tang, X., Yang, H., Shen, L., Zuo, X., et al. (2013). Induction of pluripotency in mouse somatic cells with lineage specifiers. *Cell* 153, 963–975.
- Spivakov, M., and Fisher, A.G. (2007). Epigenetic signatures of stem-cell identity. *Nat. Rev. Genet.* 8, 263–271.
- Suva, M.L., Riggi, N., and Bernstein, B.E. (2013). Epigenetic reprogramming in cancer. *Science* 339, 1567–1570.
- Terunuma, A., Putluri, N., Mishra, P., Mathé, E.A., Dorsey, T.H., Yi, M., Wallace, T.A., Issaq, H.J., Zhou, M., Killian, J.K., et al. (2014). MYC-driven accumulation of 2-hydroxyglutarate is associated with breast cancer prognosis. *J. Clin. Invest.* 124, 398–412.
- Tung, P.Y., and Knoepfler, P.S. (2015). Epigenetic mechanisms of tumorigenicity manifesting in stem cells. *Oncogene* 34, 2288–2296.
- Xiao, M., Yang, H., Xu, W., Ma, S., Lin, H., Zhu, H., Liu, L., Liu, Y., Yang, C., Xu, Y., et al. (2012). Inhibition of  $\alpha$ -KG-dependent histone and DNA demethylases by fumarate and succinate that are accumulated in mutations of FH and SDH tumor suppressors. *Genes Dev.* 26, 1326–1338.
- Xu, W., Yang, H., Liu, Y., Yang, Y., Wang, P., Kim, S.H., Ito, S., Yang, C., Wang, P., Xiao, M.T., et al. (2011). Oncometabolite 2-hydroxyglutarate is a competitive inhibitor of  $\alpha$ -ketoglutarate-dependent dioxygenases. *Cancer Cell* 19, 17–30.
- Yang, M., Soga, T., and Pollard, P.J. (2013). Oncometabolites: linking altered metabolism with cancer. *J. Clin. Invest.* 123, 3652–3658.
- Yun, J., Johnson, J.L., Hanigan, C.L., and Locasale, J.W. (2012). Interactions between epigenetics and metabolism in cancers. *Front. Oncol.* 2, 163.



**Stem Cell Reports, Volume 6**

## **Supplemental Information**

### **Oncometabolic Nuclear Reprogramming of Cancer Stemness**

**Javier A. Menendez, Bruna Corominas-Faja, Elisabet Cuyàs, María G. García, Salvador Fernández-Arroyo, Agustín F. Fernández, Jorge Joven, Mario F. Fraga, and Tomás Alarcón**

# Supplemental data

---

## Model formulation

Our model is based on a minimal gene regulatory network, which enables us to study the oncometabolic nuclear reprogramming of differentiated cells into pluripotent stem cells. This network considers the interactions between a minimal core of stemness-associated transcription factors (*OCT4* and *SOX2*) and two generic lineage-specific genes, referred to as *LSG1* and *LSG2*. The model schematically represented in **Fig. 1A** (*top panel*) is similar to a number of previous approaches (Cinquin and Demongeot, 2005; Cinquin and Page, 2007; MacArthur and Lemischka, 2013; MacArthur et al., 2008; Shu et al., 2013), particularly the one considered by Shu et al. (2013), which involves a coupled pluripotency module (i.e., self-activation of *OCT4* and *SOX2*) and a differentiation module (i.e., mutual antagonism between the *LSGs*). However, the model presented here originally introduces the epigenetic regulation of the *LSGs*, which is the essential element that enables us to account for the regulatory effects of certain metabolic features (i.e., oncometabolites) in the nuclear reprogramming process.

Our stochastic model of oncometabolic nuclear reprogramming is formulated in terms of a continuous-time Markov process governed by the corresponding master equation [1], which determines the temporal evolution of the probability density function  $P(X,t)$  for the random variable  $X(t)$ .

$$\frac{\partial P(X,t)}{\partial t} = \sum_{i=1}^R (W_i(X - r_i)P(X - r_i, t) - W_i(X)P(X, t)) \quad [1]$$

$X$  is a random vector whose components correspond to the number of cell states present in our system at time  $t$ .  $X(t)$  is therefore the state vector of the network whose components are the number of each element involved in the dynamics of the nuclear reprogramming gene network (i.e., the number of molecules of each transcription factor, the number of available binding sites left in the promoter of each gene, and the number of binding sites bound to each of three transcription factor dimers: OCT4-SOX2, LSG1-LSG1, and LSG2-LSG2).  $R$  is the number of reactions or events that the system can undergo, also referred to as channels.  $W_i(X)$  is the transition rate corresponding to channel  $i$  (i.e., the probability that the event associated with the channel  $i$  occurs in the time interval  $(t, t+\Delta t)$  is  $W_i(X(t))\Delta t$  for  $\Delta t \rightarrow 0$ ), and  $r_i$  is the change in the state vector  $X$  when channel  $i$  fires up. In mathematical terms, the probability that  $X(t+\Delta t) = X(t) + r_i$  conditioned to the system to be in state  $X(t)$  at time  $t$  is given by  $P(X(t+\Delta t) = X(t) + r_i | X(t)) = W_i(X(t))\Delta t$ .

In other words, the quantities  $W_i(X)$  and  $r_i$  are, respectively, the transition rates and the state vector change associated with the occurrence of the elementary reaction  $i$ . For each gene in the reprogramming network, our stochastic model considers the following set of elementary

reactions: (i) reversible binding of transcription factor dimers to free binding sites in the promoter region, (ii) uninduced protein synthesis, (iii) protein degradation. For the stemness-related transcription factors specifically, we consider a fourth elementary process, namely induced protein synthesis, whereby OCT4 and SOX2 proteins are synthesised independently of the binding of the OCT4-SOX2 heterodimer to their promoters. A detailed description and model of the different processes involved in our stochastic model of oncometabolic nuclear reprogramming follows.

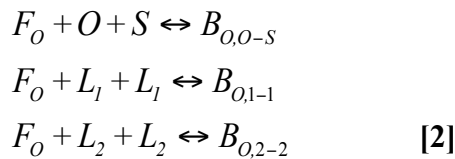
## 1. Model of the gene regulatory network

Models of gene regulatory networks controlling cell differentiation have generally been studied in terms of mean-field (i.e., deterministic) high-dimensional switches, particularly in the context of competitive heterodimerisation networks (Cinquin and Demongeot, 2005; Cinquin and Page, 2007). Instead, our stochastic model of oncometabolic nuclear reprogramming is based on premises regarding transcription factor binding to gene promoters that have been proposed by the MacArthur's group (MacArthur and Lemischka, 2013; MacArthur et al., 2008).

As mentioned above, our stochastic model considers the pluripotency genes *OCT4* and *SOX2* and two lineage-specific genes (*LSGs*). The stemness-associated transcription factors OCT4 and SOX2 upregulate each other and downregulate the expression of the *LSGs*. The OCT4 and SOX2 protein products form a dimer that binds to the gene promoters of the network; when the OCT4-SOX2 dimer binds to the promoters of *OCT4* and *SOX2*, the OCT4-SOX2 dimer upregulates the expression of *OCT4* and *SOX2*. If, by contrast, the OCT4-SOX2 dimer binds to the promoters of *LSGs*, the OCT4-SOX2 dimer downregulates the expression of *LSGs*. Our model also considers the induction of the expression of *OCT4* and *SOX2* independently of the binding of OCT4-SOX2 dimers to their promoter regions, which makes it possible to model the original Yamanaka mechanism (Takahashi and Yamanaka, 2006; Takahashi et al., 2007), where reprogramming is achieved by transfection with retroviruses encoding the RNA of *OCT4* and *SOX2*.

The stoichiometric equations for the processes involved in the regulation of the *OCT4* promoter region are as follows:

- Reversible binding of dimers to free binding sites in the promoter region of *OCT4*:



- Uninduced protein synthesis of OCT4:



- Degradation of the OCT4 protein:

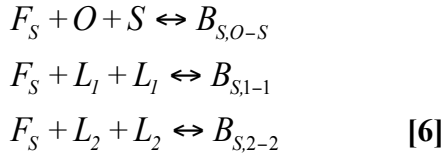


- Induced synthesis of the OCT4 protein:



Similarly, the stoichiometric equations for the processes involved in the regulation of the *SOX2* promoter region are as follows:

- Reversible binding of dimers to free binding sites in the promoter region of *SOX2*:



- Uninduced protein synthesis of *SOX2*:



- Degradation of the *SOX2* protein:



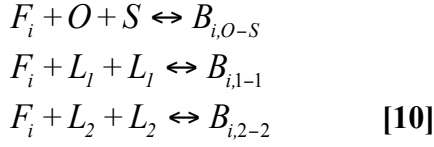
- Induced synthesis of the *SOX2* protein:



Each of the LSGs self-upregulate and downregulate the expression of the other LSGs as well as of the stemness transcription factors OCT4 and *SOX2*. The protein products of the LSGs form homodimers (LSG1-LSG1 and LSG2-LSG2) and heterodimers (LSG1-LSG2) that bind to the gene promoters; when the LSG homodimers binds the corresponding promoter, it self-upregulates the expression of the corresponding *LSG*. If the LSG homodimer binds to the promoter of the other *LSG*, the LSG homodimer downregulates the expression of the other *LSG*. The LSG heterodimer always represses the expression of the corresponding *LSG*, regardless of which promoter it binds to. All the possible LSG dimers repress the expression of the stemness factors. The stoichiometric equations corresponding to the dynamics of the LSGs are as follows:



- Reversible binding of dimers to free binding sites in the promoter regions of *LSGs*:



- Synthesis of LSG proteins:



- Degradation of LSG proteins



where  $i = 1, 2$ .

The details of how we model the dynamics produced by these stoichiometric reactions and, more importantly, how the models for gene and epigenetic regulation are coupled follow.

**1.1. Activation and downregulation of transcription: Competitive inhibition.** Here, we provide details of how we model the stoichiometric processes [2], [6] and [10]. Regulation of transcription is achieved using a model where the transcription factors of each gene in the minimal regulatory network compete for the binding sites in their promoter regions (**Fig. 1A, bottom panel**). Upon dimerisation, the protein products of the stemness-related transcription factors *OCT4* and *SOX2* and of each *LSG* bind to the promoter region of the different genes in the network. Dimers of *OCT4*-*SOX2* proteins that bind to the promoter regions of *OCT4* and *SOX2* upregulate their transcription, whereas dimers made out of combinations of *LSGs* repress the transcription of *OCT4* and *SOX2*. Similarly, *LSG* homodimers (*LSG1*-*LSG1* and *LSG2*-*LSG2*) bound to the promoters of the corresponding *LSGs* promote their self-expression. By contrast, binding of any other dimer represses the expression of *LSGs*.

To model the processes of binding and unbinding, we used the standard law of mass action kinetics (Gillespie, 1976). We made two simplifying assumptions, which we argue should have solely minor quantitative effects without altering the qualitative properties of the system. First, we did not explicitly consider the transcription factor dimerisation process. Instead, we assumed that the whole process of dimerisation and dimer binding to the corresponding promoter region can be subsumed under ternary reactions. Thus, under the usual fast kinetics assumption for the formation of the dimer, when the dimer and its components are assumed to be in equilibrium (MacArthur et al., 2008), this approximation should have only minor effects. Second, we did not consider the formation of *LSG1*-*LSG2* heterodimers. The

resulting transition rates are given in **Tables S1, S2** and **S3** for the promoter regions of *OCT4*, *SOX2* and the *LSGs*, respectively.

**Table S1**

Transition rate	$\Delta O$	$\Delta S$	$\Delta F_O$	$\Delta B_{O,O-S}$	$\Delta B_{O,1-1}$	$\Delta B_{O,2-2}$	Description
$W_{11} = k_{11}OSF_O$	-1	-1	-1	+1	0	0	<i>O-S</i> binding
$W_{12} = k_{12}B_{O,O-S}$	+1	+1	+1	-1	0	0	<i>O-S</i> unbinding
Transition rate	$\Delta L_1$	$\Delta L_2$	$\Delta F_O$	$\Delta B_{O,O-S}$	$\Delta B_{O,1-1}$	$\Delta B_{O,2-2}$	
$W_{13} = k_{13}L_1(L_1 - 1)F_O$	-2	0	-1	0	+1	0	$L_1$ dimer binding
$W_{14} = k_{14}B_{O,1-1}$	+2	0	+1	0	-1	0	$L_1$ dimer unbinding
$W_{15} = k_{15}L_2(L_2 - 1)F_O$	0	-2	-1	0	0	+1	$L_2$ dimer binding
$W_{16} = k_{16}B_{O,2-2}$	0	+2	+1	0	0	-1	$L_2$ dimer unbinding

**Transition rates corresponding to the stochastic model of competitive transcription factor binding to the promoter region of Oct4 (Equation [2]).** *O*, *S*,  $L_1$  and  $L_2$  refer to the numbers of OCT4, SOX2, LSG1 and LSG2 proteins.  $B_{O,O-S}$ ,  $B_{O,1-1}$  and  $B_{O,2-2}$  are the numbers of sites in the promoter region of *OCT4* bound to OCT4-SOX2 dimers, LSG1 dimers and LSG2 dimers, respectively.  $F_O$  is the number of free (i.e., unbound) sites in the *OCT4* promoter.  $\Delta O$ ,  $\Delta S$ ,  $\Delta L_1$ ,  $\Delta L_2$ ,  $\Delta F_O$ ,  $\Delta B_{O,O-S}$ ,  $\Delta B_{O,1-1}$  and  $\Delta B_{O,2-2}$  are the components of the state change vector corresponding to each channel. The remaining components of the state change vector have not been explicitly stated in this table, as they do not affect these reactions and are identically zero.

**Table S2**

Transition rate	$\Delta O$	$\Delta S$	$\Delta F_S$	$\Delta B_{S,O-S}$	$\Delta B_{S,1-1}$	$\Delta B_{S,2-2}$	Description
$W_{17} = k_{17}OSF_S$	-1	-1	-1	+1	0	0	<i>O-S</i> Binding
$W_{18} = k_{18}B_{S,O-S}$	+1	+1	+1	-1	0	0	<i>O-S</i> Unbinding
Transition rate	$\Delta L_1$	$\Delta L_2$	$\Delta F_S$	$\Delta B_{S,O-S}$	$\Delta B_{S,1-1}$	$\Delta B_{S,2-2}$	
$W_{19} = k_{19}L_1(L_1 - 1)F_S$	-2	0	-1	0	+1	0	$L_1$ dimer binding
$W_{20} = k_{20}B_{S,1-1}$	+2	0	+1	0	-1	0	$L_1$ dimer unbinding
$W_{21} = k_{21}L_2(L_2 - 1)F_S$	0	-2	-1	0	0	+1	$L_2$ dimer binding
$W_{22} = k_{22}B_{S,2-2}$	0	+2	+1	0	0	-1	$L_2$ dimer unbinding

**Transition rates corresponding to the stochastic model of competitive transcription factor binding to the promoter region of Oct4 (Equation [6]).** *O*, *S*,  $L_1$  and  $L_2$  refer to the

numbers of OCT4, SOX2, LSG1 and LSG2 proteins.  $B_{S,O-S}$ ,  $B_{S,1-1}$  and  $B_{S,2-2}$  are the numbers of sites in the promoter region of *SOX2* bound to OCT4-SOX2 dimers, LSG1 dimers and LSG2 dimers, respectively.  $F_S$  is the number of free (i.e., unbound) sites in the *OCT4* promoter.  $\Delta O$ ,  $\Delta S$ ,  $\Delta L_1$ ,  $\Delta L_2$ ,  $\Delta F_O$ ,  $\Delta B_{S,O-S}$ ,  $\Delta B_{S,1-1}$  and  $\Delta B_{S,2-2}$  are the components of the state change vector corresponding to each channel. The remaining components of the state change vector have not been explicitly stated in this table, as they do not affect these reactions and are identically zero.

**Table S3**

Transition rate	$\Delta O$	$\Delta S$	$\Delta F_i$	$\Delta B_{i,O-S}$	$\Delta B_{i,1-1}$	$\Delta B_{i,2-2}$	Description
$W_{23+6(i-1)} = k_{23+6(i-1)} O S F_i$	-1	-1	-1	+1	0	0	<i>O-S</i> binding
$W_{24+6(i-1)} = k_{24+6(i-1)} B_{i,O-S}$	+1	+1	+1	-1	0	0	<i>O-S</i> unbinding
Transition rate	$\Delta L_1$	$\Delta L_2$	$\Delta F_i$	$\Delta B_{i,O-S}$	$\Delta B_{i,1-1}$	$\Delta B_{i,2-2}$	
$W_{25+6(i-1)} = k_{25+6(i-1)} L_1 (L_1 - 1) F_i$	-2	0	-1	0	+1	0	$L_1$ dimer binding
$W_{26+6(i-1)} = k_{26+6(i-1)} B_{i,1-1}$	+2	0	+1	0	-1	0	$L_1$ dimer unbinding
$W_{27+6(i-1)} = k_{27+6(i-1)} L_2 (L_2 - 1) F_i$	0	-2	-1	0	0	+1	$L_2$ dimer binding
$W_{28+6(i-1)} = k_{28+6(i-1)} B_{i,2-2}$	0	+2	+1	0	0	-1	$L_2$ dimer unbinding

**Transition rates corresponding to the stochastic model of competitive transcription factor binding to the promoter region of the LSGs (Equation [10]).**  $O$ ,  $S$ ,  $L_1$  and  $L_2$  refer to the numbers of OCT4, SOX2, LSG1 and LSG2 proteins.  $B_{i,O-S}$ ,  $B_{i,1-1}$  and  $B_{i,2-2}$  are the numbers of sites in the promoter region of the *LSGs* bound to OCT4-SOX2 dimers, LSG1 dimers and LSG2 dimers, respectively.  $F_i$  is the number of free (i.e., unbound) sites in the *LSGs* promoters.  $\Delta O$ ,  $\Delta S$ ,  $\Delta L_1$ ,  $\Delta L_2$ ,  $\Delta F_i$ ,  $\Delta B_{i,O-S}$ ,  $\Delta B_{i,1-1}$  and  $\Delta B_{i,2-2}$  are the components of the state change vector corresponding to each channel. The remaining components of the state change vector have not been explicitly stated in this table, as they do not affect these reactions and are identically zero. The index  $i$  ( $i = 1, 2$ ) spans the set of the different LSGs considered.

The reactions [2], [6], and [10], whose transition rates are given in **Tables S1**, **S2**, and **S3**, respectively, are such that the total number of binding sites within the promoter region of each gene must be conserved, i.e., the following balance equations must be satisfied:

$$F_O + B_{O,O-S} + B_{O,1-1} + B_{O,2-2} = T_O \quad [13]$$

$$F_S + B_{S,O-S} + B_{S,1-1} + B_{S,2-2} = T_S \quad [14]$$

$$F_1 + B_{1,O-S} + B_{1,1-1} + B_{1,2-2} = T_1(A_1) \quad [15]$$

$$F_2 + B_{2,O-S} + B_{2,I-I} + B_{2,2-2} = T_2(A_2) \quad [16]$$

The above equations, particularly [15] and [16], are of crucial relevance for the formulation of our stochastic model, as they incorporate the coupling between genetic and epigenetic regulation. Specifically, we consider that the *LSGs* are under epigenetic regulation and that the number of binding sites within the promoter regions of *LSGs* that are accessible to transcription factors at every given time,  $T_1$  and  $T_2$ , are determined by the acetylation status. As a first approximation, we consider the case where  $T_1 = C_1 A_1$  and  $T_2 = C_2 A_2$ , where  $C_1$  and  $C_2$  are constant. By contrast, we consider the case where  $T_O$  and  $T_S$  are constant. The process of introducing the model for protein synthesis and degradation is described in **section 1.2** (see below).

**1.2. Uninduced and induced protein synthesis and degradation.** We proceed further by introducing the model for protein synthesis and degradation. We should first consider the stoichiometric processes described by equations [3], [4] and [5], which include protein synthesis and degradation of OCT4. The equation [3] describes the uninduced protein synthesis of OCT4. Because the transcription of *OCT4* is upregulated by the binding of OCT4-SOX2 dimers to the free sites within the *OCT4* promoter, we model the synthesis of the OCT4 protein by assuming that its probability rate is proportional to the number of OCT4-SOX2-bound sites ( $B_{O,O-S}$ ). It is relevant to note that uninduced synthesis of OCT4 requires the presence of OCT4 protein, i.e., in the absence of OCT4 protein,  $B_{O,O-S} = 0$ , and consequently, no synthesis of OCT4 occurs. By contrast, the induced synthesis of OCT4 occurs at a constant rate,  $R_1$ , independently of  $B_{O,O-S}$ , which means that there is a positive probability of induced synthesis of OCT4 at all times. These two processes give rise to a global rate of OCT4 synthesis given by  $W_1 = k_1 B_{O,O-S} + R_1$  (see **Table S4**). The degradation process is modelled by the standard first-order decay kinetics (as shown in the  $W_2$ -entry of **Table S4**). The principles for SOX2 synthesis and degradation are exactly the same as for OCT4. The dynamics of LSG synthesis and degradation are also determined by the same kinetics, except that we do not consider induced synthesis for the LSGs.

**Table S4**

Transition rate	$\Delta O$	$\Delta S$	$\Delta L_1$	$\Delta L_2$	Description
$W_1 = k_1 B_{O,O-S} + R_1$	+1	0	0	0	Oct4 synthesis
$W_2 = k_2 O$	-1	0	0	0	Oct4 degradation
$W_3 = k_3 B_{S,O-S} + R_1$	0	+1	0	0	Sox2 synthesis
$W_4 = k_4 S$	0	-1	0	0	Sox2 degradation
$W_5 = k_5 B_{1,1-1}$	0	0	+1	0	LSG1 synthesis
$W_6 = k_6 L_1$	0	0	-1	0	LSG1 degradation



$W_7 = k_7 B_{2,2-2}$	0	0	0	+1	LSG2 synthesis
$W_8 = k_8 L_2$	0	0	0	-1	LSG2 degradation

**Transition rates corresponding to the stochastic model of competitive transcription factor binding to the promoter region of the LSGs (Equation [10]).**  $O$ ,  $S$ ,  $L_1$  and  $L_2$  refer to the numbers of OCT4, SOX2, LSG1 and LSG2 proteins.  $B_{i,O-S}$ ,  $B_{i,1-1}$  and  $B_{i,2-2}$  are the numbers of sites in the promoter region of the LSGs bound to OCT4-SOX2 dimers, LSG1 dimers and LSG2 dimers, respectively.  $F_i$  is the number of free (i.e., unbound) sites in the LSGs promoters.  $\Delta O$ ,  $\Delta S$ ,  $\Delta L_1$  and  $\Delta L_2$  are the components of the state change vector corresponding to each channel. The remaining components of the state change vector have not been explicitly stated in this table, as they do not affect these reactions and are identically zero.

## 2. Epigenetic regulation of the lineage-specific genes (LSGs)

Our stochastic model of oncometabolic nuclear reprogramming is a generalisation of the epigenetic model proposed by Dodd et al. (2007), which considers nucleosome modification as the basic mechanism for epigenetic cell memory. Nucleosomes are assumed to be in one of three states, namely methylated ( $M$ ), unmodified ( $U$ ) and acetylated ( $A$ ), and the dynamics of the model is given in terms of the transition rates between the  $M$ ,  $U$ , and  $A$  nucleosome states.

As originally postulated by Dodd et al. (2007), our model considers direct transitions between  $M$  and  $A$  to be highly unlikely. Instead, our models assume that transitions occur in a linear sequence [17] in which  $M$  nucleosomes can only undergo loss of the corresponding methyl group to become  $U$  nucleosomes, which can then, through the intervention of the corresponding histone-modifying enzyme, acquire an acetyl group to become  $A$  nucleosomes, and vice versa.

$$M_i \leftrightarrow U_i \leftrightarrow A_i \quad [17]$$

The subindex  $i$  in [17] spans the set of LSGs considered (in the current model  $i = 1, 2$ ).

Nucleosome modifications are of two types, namely recruited and unrecruited. A recruited modification refers to a positive feedback mechanism where change in the modification status of the nucleosome is facilitated by the presence of other modified nucleosomes (i.e., by the presence of other methylated or acetylated nucleosomes). Mathematically, a recruited modification is expressed through a non-linear dependence on the number of  $M$ -nucleosomes and  $A$ -nucleosomes of the corresponding transition rates (see **Table S5**). An unrecruited modification refers to nucleosome modifications whose probability is independent of the modification status of the other nucleosomes. The corresponding transition rates are given in **Table S5**.

**Table S5**

Transition rate	$\Delta M_i$	$\Delta U_i$	$\Delta A_i$	Description
$W_{35+8(i-1)} = k_{35+8(i-1)} h_2 M_i A_i$	-1	+1	0	Recruited nucleosome demethylation
$W_{36+8(i-1)} = k_{36+8(i-1)} h_1 M_i U_i$	+1	-1	0	Recruited nucleosome methylation
$W_{37+8(i-1)} = k_{37+8(i-1)} h_4 M_i A_i$	0	+1	-1	Recruited nucleosome deacetylation
$W_{38+8(i-1)} = k_{38+8(i-1)} h_3 M_i U_i$	0	-1	+1	Recruited nucleosome acetylation
$W_{39+8(i-1)} = k_{39+8(i-1)} M_i$	-1	+1	0	Unrecruited nucleosome demethylation
$W_{41+8(i-1)} = k_{41+8(i-1)} h_1 U_i$	+1	-1	0	Unrecruited nucleosome methylation
$W_{40+8(i-1)} = k_{40+8(i-1)} A_i$	0	+1	-1	Unrecruited nucleosome deacetylation
$W_{42+8(i-1)} = k_{42+8(i-1)} h_3 U_i$	0	-1	+1	Unrecruited nucleosome acetylation

**Transition rates corresponding to the stochastic model of epigenetic regulation.**  $M_i$ ,  $U_i$  and  $A_i$  with  $i=1,2$  are the number of methylated, unmodified, and acetylated nucleosomes corresponding to the LSG  $i$ .  $\Delta M_i$ ,  $\Delta U_i$  and  $\Delta A_i$  are the components of the state change vector corresponding to each reaction channel. The remaining components of the state change vector have not been explicitly stated in this table, as they do not affect these reactions and are identically zero.

We added an additional element to the model by Dodd et al. (2007). Nucleosome modifications are mediated by four types of enzymes, namely histone methyltransferases (HMT) and histone acetyl transferases (HAT), which mediate methylation and acetylation, and histone demethylases (HDM) and histone deacetylases (HDAC), which catalyse demethylation and deacetylation. The activities of HAT, HMT, HDM and HDAC are accounted for by the parameters  $h_1$ ,  $h_3$ ,  $h_2$  and  $h_4$ , respectively (see **Table S5**). These enzymes play a pivotal role in our stochastic model of oncometabolic nuclear reprogramming, as their activity status is the coupling effector between a given metabolic feature (i.e., oncometabolites) and the regulatory differentiation/reprogramming system. In particular, we considered the metabolo-epigenetic connection originally proposed by Thompson's group (Lu and Thompson, 2012; Lu et al., 2012). They were pioneers in showing that, whereas the wild-type form of the IDH enzymes helps cell differentiation to proceed normally by providing HDMs with the corresponding cofactors, the activity of the mutant form of the IDH enzyme drastically hinders HDM activity. The wild-type IDH enzyme catalyses the interconversion of isocitrate and alpha-ketoglutarate ( $\alpha$ KG) in cytosol or mitochondria, whereas the neomorphic activity of the mutant form of IDH produces the oncometabolite 2HG from  $\alpha$ KG. As the differentiation process is regulated by  $\alpha$ KG-dependent HDMs, such as JHDM, the oncometabolite 2HG effectively reduces HDM activity to lock *IDH* mutated-cells in a state poised for the acquisition of pluripotency.

We postulated that our stochastic model should be able to simulate how oncometabolite-induced reduction of HDM, along with the activation of reprogramming stimuli (i.e., the stemness-related transcription factors *OCT4* and *SOX2*), actually leads to a significant increase in the speed and efficiency of the nuclear reprogramming phenomenon.

# Supplemental experimental procedures

---

## A. The stochastic model

We employed two different techniques to analyse the diverse aspects of our stochastic model of oncometabolic nuclear reprogramming, namely, direct numerical simulation using Gillespie's stochastic simulation algorithm (SSA) (see **section A. 1**) and the semiclassical approximation (see **section A. 2**). Gillespie's SSA was utilised to study the kinetics of the reprogramming process, whereas the semiclassical approximation enabled us to explore issues such as the enlargement of the basin of attraction of the stemness state upon oncometabolic reduction of HDM activity.

**A. 1. Gillespie's stochastic simulation algorithm.** Gillespie's SSA is a numerical simulation technique that enables us to generate exact sample paths whose probability density is the solution of the master equation [1]. The SSA is a standard technique in the simulation of Markov stochastic processes, and therefore we will not provide a full description here. In short, the algorithm is based on a reformulation of the process described by the master equation, whereby its evolution is driven by the iteration of the following two steps: (i) generation of the exponentially distributed waiting time until the next event and (ii) random selection of the reaction channel to fire up once the waiting time has elapsed (Gillespie, 1976). We used the SSA to produce the numerical results corresponding to the process determined by the reaction rates given in **Tables S1, S2, S3 and S4**, as shown in **Figs. 1B-D, 2A and S3**.

**A. 2. Semiclassical approximation: optimal reprogramming paths.** An alternative way to analyse the dynamics of continuous-time Markov processes on a discrete space of states is to derive an equation for the generating function,  $G(p_1, \dots, p_n, t)$  of the corresponding probabilistic density:

$$G(p_1, \dots, p_n, t) = \sum_x p_1^{x_1} p_2^{x_2} \dots p_n^{x_n} P(x_1, \dots, x_n, t) \quad [18]$$

where  $P(x_1, \dots, x_n, t) = P(X, t)$  is the solution of the master equation (1).  $P(x_1, \dots, x_n, t)$  satisfies a partial differential equation that can be derived from the master equation [1]. This partial differential equation (PDE) is the basic element of the so-called momentum representation of the master equation (1). The corresponding PDE can be solved explicitly only in a few simple cases. However, although closed analytical solutions are rarely available, the PDE for the generating function admits a Wentzel-Kramers-Brillouin (WKB) perturbative solution (Assaf et al., 2010; Dykman et al., 1995; Kubo et al., 1973).

The (linear) PDE that governs the evolution of the generating function, which is derived from the master equation (1) by multiplying both sides by  $p_1^{x_1} p_2^{x_2} \dots p_n^{x_n}$  and summing over all the values of  $X = (x_1, \dots, x_n)$ , can be written as follows:

$$\frac{\partial G}{\partial t} = H(p_1, \dots, p_n, \partial_{p_1}, \dots, \partial_{p_n}) G(p_1, \dots, p_n, t) \quad [19]$$

where the operator  $H$  is determined by the reaction rates of the master equation [1]. Furthermore, the solution to this equation must satisfy the normalisation condition  $G(p_1 = 1, \dots, p_n = 1, t) = 1$  for all  $t$ . This PDE can be solved analytically only in a few simple cases. However, although closed analytical solutions are rarely available, the PDE for the generating function admits a WKB perturbative solution (Assaf et al., 2010; Dykman et al., 1995; Kubo et al., 1973).

From the mathematical point of view, [19] is a Schrödinger-like equation, and, therefore, there is a plethora of methods at our disposal to analyse it. In particular, when the fluctuations are assumed to be small, it is common to resort to WKB perturbation expansion of the solution of [19]. This approach is based on the WKB-like Ansatz that  $G(p_1, \dots, p_n, t) = e^{-S(p_1, \dots, p_n, t)}$ . By substituting this Ansatz into [19], we obtain the following Hamilton-Jacobi equation for the action  $S(p_1, \dots, p_n, t)$ :

$$\frac{\partial S}{\partial t} = -H\left(p_1, \dots, p_n, \frac{\partial S}{\partial p_1}, \dots, \frac{\partial S}{\partial p_n}\right) \quad [20]$$

Instead of directly seeking the explicit solution of [20], we exploit the Hamilton approach by using the Feynman path-integral representation, which yields a solution to [19] of the following type (Feynman and Hibbs, 1965; Dickman and Vidigal, 2003; Täuber et al., 2005):

$$G(p_1, \dots, p_n, t) = \int_0^t e^{-S(p_1, \dots, p_n, q_1, \dots, q_n, s)} Dq(s) Dp(s) \quad [21]$$

with  $S(p_1, \dots, p_n, q_1, \dots, q_n)$  given by:

$$S(p_1, \dots, p_n, t) = \int_0^t \left( -H(p_1, \dots, p_n, q_1, \dots, q_n) + \sum_{i=1}^n \dot{p}_i(s) q_i(s) \right) + S(p_1, \dots, p_n, t=0) \quad [22]$$

where the position operators in the momentum representation are defined as  $q_i \equiv \partial p_i$ , with the commutation relation  $[p_i, q_j] = \delta_{i,j}$ .

The so-called semi-classical approximation consists of approximating the path integral in [21] by Sakurai (1994):

$$G(p_1, \dots, p_n, t) = e^{-So[p_1, \dots, p_n, t]} \quad [23]$$

where  $So(p_1, \dots, p_n, t)$  is the functional action calculated by integrating [22] over the solution of the corresponding Hamilton equations, i.e., the orbits that maximise the action  $S$  :

$$\frac{dp_i}{dt} = -\frac{\partial H}{\partial q_i} \quad [24]$$

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i} \quad [25]$$

where  $[q_i, p_j]$  are the generalised coordinates corresponding to chemical species  $i = 1, \dots, n$ . These equations are formally solved with boundary conditions  $q_i(0) = n_i(0)$  such that  $p_i(t) = p$  (Elgart and Kamenev, 2004).

This approach has been recently applied to the analysis of complex protein-interaction systems involving separation of time scales. In particular, the semi-classical approximation has been used to formulate a stochastic version of the quasi-steady approximation (Alarcón, 2014), a framework that we use to simplify the analysis of our system.

**A.2.1. Quasi-steady state approximation.** The system of Hamilton equations [24] and [25] with Hamilton given by equations [A2] to [A9] is far too complex to analyse. To gain insight into the behaviour of the system, we simplify it by performing a quasi-steady state approximation (QSSA). This technique has been extensively used in biochemical modelling (Keener and Sneyd, 1998) and relies on the existence of different time scales whereby one can distinguish between slow and fast variables. Once one has sorted the system variables according to this criterion, one proceeds to make an adiabatic approximation where the fast variables are assumed to be in equilibrium with the slow ones. This procedure reduces the dimension of the system (i.e., number of variables) as well as the number of independent parameters to be determined.

We have recently developed a stochastic QSSA for the semi-classical treatment of complex networks of protein/gene interactions based on the Briggs-Haldane analysis of the Michaelis-Menten system for enzyme catalysis (Alarcón, 2014). According to this analysis, a systematic separation of the system variables into slow and fast can be achieved provided that the number of enzyme molecules is much smaller than the number of substrate molecules. To apply the method put forward in Alarcón (2014) to our current stochastic model, we make the following assumption regarding the characteristic scales for the different values of our stochastic formulation: we assume that the number of all the proteins in our model, i.e.,  $O$ ,  $S$ ,  $L_1$  and  $L_2$  and their counterparts in the semi-classical approximation  $q_1, q_3, q_5$  and  $q_7$ , have

the same characteristic scales, which we denote by  $y_0$ . Therefore, once transient regimens have been overcome,

$$x_i = \frac{q_i}{y_0} = O(1), i=1,3,5,7 \quad [26]$$

We further assume that the quantities  $T_o, T_s, T_l$  and  $T_2$  (see equations [13] to [16]) are all of the same order of magnitude and all have the same characteristic scale,  $T_o$ .  $T_o$  is the natural scale for the variables corresponding to the free and bound binding sites in the gene promoter regions as well as for the epigenetic regulation variables ( $M_i, U_i$  and  $A_i$  or the counterparts  $q_i, i=22, \dots, 27$ ), i.e.,:

$$x_i = \frac{q_i}{T_o} = O(1), i \neq 1,3,5,7 \quad [27]$$

The separation of time scales necessary to apply the QSSA emerges if we assume the following:

$$\varepsilon = T_o \ll 1 \quad [28]$$

To proceed further, we consider the Hamilton equations [24] and [25] and re-scale the time variable,  $\tau = k_{11}T_o y_0^2 t$ , and the Hamiltonian  $H_k(p, q) = k_{11}T_o y_0^2 H_\kappa(p, x)$ , where  $x_i$  are the re-scaled variables defined in [26] and [27]. When re-scaling the Hamiltonian, the parameters  $k_i$  featured in [A4] and [A9] must also be re-scaled and are thus substituted by the corresponding parameters  $\kappa_i$  given in **Table S6**. For simplicity, we split the re-scaled Hamiltonian into two contributions,  $H_\kappa(p, x) = H_{GR}(p, x) + H_{ER}(p, x)$ , which we analyse separately.  $H_{GR}(p, x) = H_{SD}(p, x) + H_{PO}(p, x)H_{PS}(p, x) + \sum_i H_{Pi}(p, x)$  is the Hamiltonian associated with gene regulation, and  $H_{ER}(p, x) = \sum_i H_{Ei}(p, x)$  is the Hamiltonian corresponding to the epigenetic regulation of the *LSGs*.

**Table S6**

Re-scaled parameter

---

$$\rho_1 = \frac{R_1}{k_{11}T_o y_0^2}$$

$$\kappa_i = \frac{k_i}{k_{11}T_o y_0} \quad \text{If } i=2,4,6,8$$



$$\begin{aligned}
\kappa_i &= \frac{k_i}{k_{11}} && \text{If } i=13,15,17,19,21,23,25,27,29,31,33,35 \\
\kappa_i &= \frac{k_i T_0}{k_{11} y_0} && \text{If } i=35+8(j-1), 36+8(j-1), 37+8(j-1), 38+8(j-1), j=1,2 \\
\kappa_i &= \frac{k_i}{k_{11} y_0} && \text{If } i=39+8(j-1), 40+8(j-1), 41+8(j-1), 42+8(j-1), j=1,2 \\
\kappa_i &= \frac{k_i}{k_{11} y_0^2} && \text{otherwise}
\end{aligned}$$


---

**Re-scaled parameters that result upon re-scaling of the Hamiltonian (A2):**  
 $H_k(p, q) = k_{11} T_0 y_0^2 H_\kappa(p, x)$

The QSSA for the gene regulation Hamiltonian,  $H_{GR}(p, x)$ , derived in detail in **Appendix B**, is given by the following system of ordinary differential equations:

$$\frac{dx_l}{d\tau} = \rho_l + p_2 \frac{\kappa_1}{\kappa_{12}} \frac{\vartheta_O x_1 x_3}{1 + \frac{1}{\kappa_{12}} x_1 x_3 + \frac{\kappa_{13}}{\kappa_{14}} x_5^2 + \frac{\kappa_{15}}{\kappa_{16}} x_7^2} - \kappa_2 x_1 \quad [29]$$

$$\frac{dx_3}{d\tau} = \rho_l + p_4 \kappa_3 \frac{\kappa_{17}}{\kappa_{18}} \frac{\vartheta_S x_1 x_3}{1 + \frac{\kappa_{17}}{\kappa_{18}} x_1 x_3 + \frac{\kappa_{19}}{\kappa_{20}} x_5^2 + \frac{\kappa_{21}}{\kappa_{22}} x_7^2} - \kappa_4 x_3 \quad [30]$$

$$\frac{dx_5}{d\tau} = p_6 \kappa_5 \frac{\kappa_{25}}{\kappa_{26}} \frac{\vartheta_1 (x_{24}) x_5^2}{1 + \frac{\kappa_{23}}{\kappa_{24}} x_1 x_3 + \frac{\kappa_{25}}{\kappa_{26}} x_5^2 + \frac{\kappa_{27}}{\kappa_{28}} x_7^2} - \kappa_6 x_5 \quad [31]$$

$$\frac{dx_7}{d\tau} = p_8 \kappa_7 \frac{\kappa_{31}}{\kappa_{32}} \frac{\kappa_2 (x_{27}) x_7^2}{1 + \frac{\kappa_{29}}{\kappa_{30}} x_1 x_3 + \frac{\kappa_{33}}{\kappa_{34}} x_5^2 + \frac{\kappa_{31}}{\kappa_{32}} x_7^2} - \kappa_8 x_7 \quad [32]$$

for the re-scaled number of Oct4, Sox2, LSG1 and LSG2 proteins, respectively, where we have defined  $\vartheta_O = T_O / T_0$ ,  $\vartheta_S = T_S / T_0$ ,  $\vartheta_1 = T_1 / T_0$ , and  $\vartheta_2 = T_2 / T_0$ . It should be noted that  $x_{24}$  and  $x_{27}$  are the re-scaled canonical coordinates corresponding to  $A_1$  and  $A_2$ , i.e., the acetylation levels of LSG1 and LSG2, respectively. The corresponding QSSA for the momenta yields the following results:

$$p_1 = p_3 = p_5 = p_7 = 1 \quad [33]$$

$$p_2 = p_{10} = p_{11} = p_{12} = c_O \quad [34]$$

$$p_4 = p_{13} = p_{14} = p_{15} = c_S \quad [35]$$

$$p_6 = p_{16} = p_{17} = p_{18} = c_1 \quad [36]$$

$$p_8 = p_{19} = p_{20} = p_{21} = c_2 \quad [37]$$

where  $c_0$ ,  $c_s$ ,  $c_1$ ,  $c_2$  are given in **Appendix C**, where we show that they are determined by the distribution of the number of binding sites in the promoter region of the respective genes in a population of cells.

To close the QSSA system equations [29] to [32], it is necessary to determine the dynamic of  $x_{24}$  and  $x_{27}$ , for which we turn now to the analysis of the epigenetic Hamiltonian,  $H_{ER}(p, x) = H_{E_1}(p, x) + H_{E_2}(p, x)$  (see **Appendix A**). A direct numerical solution of the corresponding Hamilton equations [24] and [25] shows that for  $x_{22}$ ,  $x_{26}$  and  $x_{27}$  to remain positive, the corresponding momenta must satisfy the following:

$$p_{22} = p_{23} = p_{24} = \text{const.} = c_{E_1} \quad [38]$$

$$p_{25} = p_{26} = p_{27} = \text{const.} = c_{E_2} \quad [39]$$

where the constants  $c_{E_1}$  and  $c_{E_2}$  are resolved in **Appendix C**, where we show that they are determined by the distribution of the number of modification sites in the nucleosomes of the respective genes in a population of cells. Considering this fact, we can write the corresponding evolution equations for  $x_{22}$ ,  $x_{23}$  and  $x_{24}$  and  $x_{25}$ ,  $x_{26}$  and  $x_{27}$ :

$$\frac{dx_{22}}{d\tau} = -\kappa_{35}h_2c_{E_1}x_{22}x_{24} - \kappa_{39}x_{22} + \kappa_{36}h_1c_{E_1}x_{22}x_{23} + \kappa_{41}h_1x_{23} \quad [40]$$

$$\frac{dx_{24}}{d\tau} = -\kappa_{37}h_{24}c_{E_1}x_{22}x_{24} - \kappa_{40}x_{24} + \kappa_{38}h_3c_{E_1}x_{24}x_{23} + \kappa_{42}h_3x_{23} \quad [41]$$

$$x_{23} = 1 - x_{22} - x_{24} \quad [42]$$

and

$$\frac{dx_{25}}{d\tau} = -\kappa_{43}h_4c_{E_2}x_{25}x_{26} - \kappa_{47}x_{25} + \kappa_{44}h_1c_{E_2}x_{25}x_{27} + \kappa_{49}h_1x_{26} \quad [43]$$

$$\frac{dx_{27}}{d\tau} = -\kappa_{45}h_4c_{E_2}x_{25}x_{27} - \kappa_{48}x_{27} + \kappa_{46}h_3c_{E_2}x_{27}x_{26} + \kappa_{50}h_3x_{26} \quad [44]$$

$$x_{26} = 1 - x_{25} - x_{27} \quad [45]$$

respectively.

### A. 3. Parameter values

The last element that remains to be addressed regarding our model formulation concerns the values of the biophysical parameters (e.g., binding and unbinding rates, protein synthesis and degradation rates) that determine the behaviour of our system. To fix the value of the parameters in our model, we require the following three properties be satisfied:

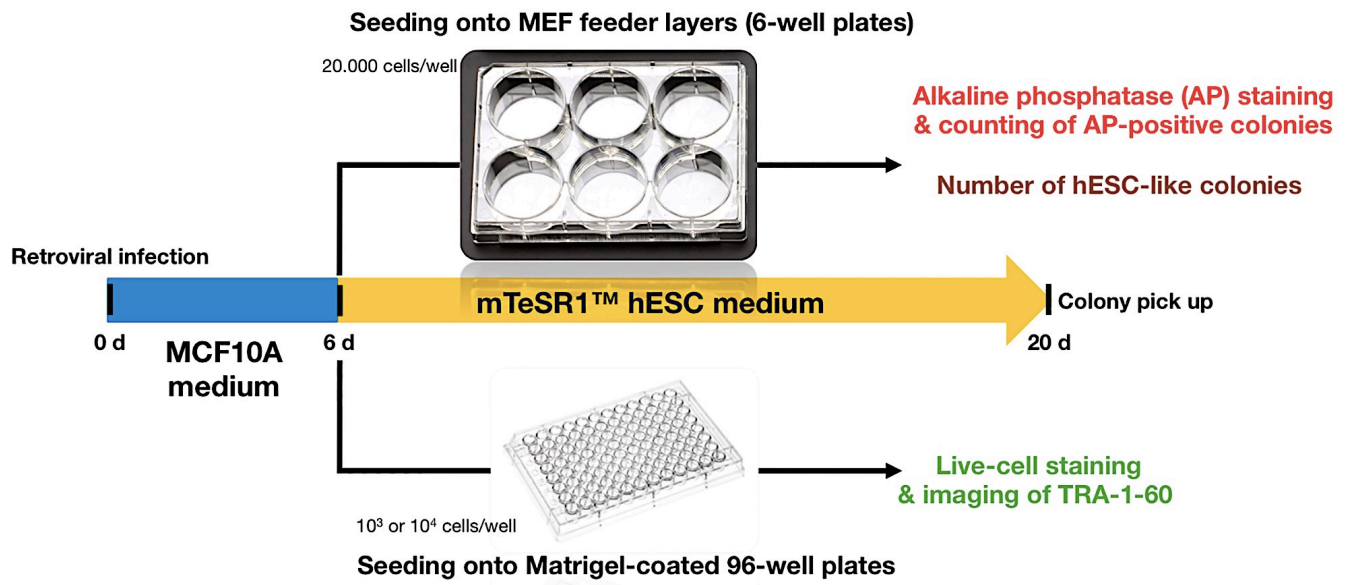
- (1) The epigenetic regulation models, equations [40] and [41] and [43] to [45], should be in their bistable regimes, where two stable steady states corresponding to methylated and acetylated states exist.
- (2) A baseline scenario must exist where the acetylated states of both *LSGs* are the more stable ones, which, in turn, should give rise to cells that exist in a differentiated state.
- (3) The gene regulation model, equations [29] to [32], must be such that, in the baseline scenario, three stable steady-states exist: a pluripotent state,  $x_p^*$ , such that the stemness-related genes (*OCT4* and *SOX2*) have positive levels of expression and the *LSGs* are not being expressed, and two differentiation states,  $x_1^*$  and  $x_2^*$ , where *LSG1* and *LSG2*, respectively, have positive levels of expression, whereas all the other genes have vanishing levels of expression.

Our so-called baseline scenario is characterised by (i) normal levels of the oncometabolite 2HG (i.e., normal activity of the histone de-methylating enzymes) and (ii) no induction of the stemness-related genes (i.e.,  $\rho_1 = \rho_2 = 0$ ). These requirements are not sufficient to uniquely determine the values of all the parameters in our model, but they help us establish the region of parameter values where our stochastic model makes biological sense.

To further simplify the analysis and without loss of generality in our results, we assume the following: (i) all the dimer-promoter binding constants are the same; (ii) all the dimer-promoter unbinding constants are the same; (iii) all the gene products are synthesised at the same time; and (iv) all the proteins are degraded at the same rate. We make a similar simplifying assumption regarding the epigenetic regulatory system, namely, that all the rates corresponding to recruited modification have the same value. Similarly, all the rates of unrecruited modification are assumed to have the same value.

The conditions that the parameter values of our stochastic model must satisfy for the three properties listed above are derived in **Appendix D** and **Appendix E**. Parameter values satisfying such conditions and compatible with the baseline scenario are given in **Table S7**.

## B. Experiments on living cells



**Figure S1.** A timeline for the overall iPS cell derivation protocol employed in the proof-of-concept validation experiments on living cells is outlined.

# Supplemental appendices

---

## APPENDIX A. The Hamiltonian function

As the master equation [1] is linear, we can arrange it by splitting its right-hand side operator as follows:

$$\begin{aligned}
 \frac{\partial P(X,t)}{\partial t} = & \sum_{i=1}^8 W_i(X-r_i)P(X-r_i,t) - W_i(X)P(X,t) \\
 & + \sum_{i=11}^{16} W_i(X-r_i)P(X-r_i,t) - W_i(X)P(X,t) \\
 & + \sum_{i=17}^{22} W_i(X-r_i)P(X-r_i,t) - W_i(X)P(X,t) \\
 & + \sum_{i=23}^{28} W_i(X-r_i)P(X-r_i,t) - W_i(X)P(X,t) \\
 & + \sum_{i=29}^{34} W_i(X-r_i)P(X-r_i,t) - W_i(X)P(X,t) \\
 & + \sum_{i=35}^{42} W_i(X-r_i)P(X-r_i,t) - W_i(X)P(X,t) \\
 & + \sum_{i=43}^{50} W_i(X-r_i)P(X-r_i,t) - W_i(X)P(X,t)
 \end{aligned} \tag{A1}$$

where the first summation corresponds to the processes described in **Table S4**, i.e., protein synthesis and degradation. The second to fifth summations correspond to the competitive binding to the gene promoters as given in **Tables S1, S2 and S3**. The last two summations in [A1] correspond to the dynamics of epigenetic regulation of each of the LSGs and are determined by the rates given in **Table S5**.

To obtain the corresponding Hamiltonian, we begin by multiplying both sides of [A1] by  $p_{21}, \dots, p_{27}$  and sum over all the possible values of the state vector  $X = (O, S, \dots, U_2, A_2)$ . This procedure yields a partial differential equation for the probability generating function,  $G(p_{21}, \dots, p_{27})$ , of the type of equation [19], where the time evolution of the generating function is driven by the operator  $G(p_{21}, \dots, p_{27}, \partial_{p1}, \dots, \partial_{p27})$ . The representation of this operator where  $\partial_{pi}$  are represented by  $q_i = \partial_{pi}$  yields the Hamiltonian corresponding to [A1]. Following the arrangement of [A1], the Hamiltonian  $H_\kappa(p, q)$ , can be written as follows:

$$H_\kappa(p, x) = H_{SD}(p, x) + H_{PO}(p, x) + H_{PS}(p, x) + \sum_i H_{Pi}(p, x) + H_{Ei}(p, x) \tag{A2}$$

where  $H_{P_O}(p, q)$ ,  $H_{P_S}(p, q)$ ,  $H_{P_i}(p, q)$ ,  $H_{SD}(p, q)$  and  $H_{E_i}(p, q)$  are the Hamiltonians corresponding to the processes described in **Tables S1, S2, S3, S4** and **S5**, respectively. These Hamiltonians can be easily computed from **[A1]** and are given by:

$$\begin{aligned} H_{SD}(p, q) = & k_1 p_2 (p_1 - 1) q_2 + R(p_1 - 1) + \\ & k_2 (1 - p_1) q_1 + k_3 p_4 (p_3 - 1) q_4 + R(p_3 - 1) + k_4 (1 - p_3) q_3 \\ & + k_5 p_6 (p_5 - 1) q_6 + k_6 (1 - p_5) q_6 + k_7 p_8 (p_7 - 1) q_8 + k_8 (1 - p_7) q_7 \end{aligned} \quad [\text{A3}]$$

where the pairs  $(p_i, q_i)$  with  $i = 1, \dots, 8$  are the generalised coordinates corresponding to  $O$ ,  $B_{O,O-S}$ ,  $S$ ,  $B_{S,O-S}$ ,  $L_1$ ,  $B_{1,1-1}$ ,  $L_2$  and  $B_{2,2-2}$ , respectively,

$$\begin{aligned} H_{P_O}(p, q) = & k_{11} (p_2 - p_9 p_{10}) q_9 q_{10} + \\ & k_{12} (p_9 p_{10} - p_2) q_2 + k_{13} (p_{11} - p_5^2 p_{10}) q_5^2 q_{10} + k_{14} (p_5^2 p_{10} - p_{11}) q_{11} \\ & + k_{15} (p_{12} - p_7^2 p_{10}) q_7^2 q_{10} + k_{16} (p_7^2 p_{10} - p_{12}) q_{12} \end{aligned} \quad [\text{A4}]$$

where  $(p_i, q_i)$ ,  $i = 10, 11, 12$  are the generalised coordinates corresponding to  $F_O$ ,  $B_{O,1-1}$  and  $B_{O,2-2}$ , respectively,

$$\begin{aligned} H_{P_S}(p, q) = & k_{14} (p_4 - p_9 p_{13}) q_9 q_{13} + \\ & k_{18} (p_9 p_{10} - p_4) q_4 + k_{19} (p_{14} - p_5^2 p_{13}) q_5^2 q_{13} + k_{20} (p_5^2 p_{13} - p_{14}) q_{14} \\ & + k_{21} (p_{15} - p_7^2 p_{13}) q_7^2 q_{13} + k_{22} (p_7^2 p_{10} - p_{15}) q_{15} \end{aligned} \quad [\text{A5}]$$

where  $(p_i, q_i)$ ,  $i = 13, 14, 15$  are the generalised coordinates corresponding to  $F_S$ ,  $B_{S,1-1}$  and  $B_{S,2-2}$ , respectively,

$$\begin{aligned} H_{P_1}(p, q) = & k_{23} (p_{18} - p_9 p_{16}) q_9 q_{16} + \\ & k_{24} (p_9 p_{16} - p_{18}) q_{18} + k_{25} (p_6 - p_5^2 p_{16}) q_5^2 q_{16} + k_{26} (p_5^2 p_{13} - p_6) q_6 \\ & + k_{27} (p_{17} - p_7^2 p_{16}) q_7^2 q_{16} + k_{28} (p_7^2 p_{16} - p_{17}) q_{17} \end{aligned} \quad [\text{A6}]$$

$$\begin{aligned} H_{P_2}(p, q) = & k_{29} (p_{21} - p_9 p_{19}) q_9 q_{19} + \\ & k_{30} (p_9 p_{19} - p_{21}) q_{21} + k_{31} (p_8 - p_7^2 p_{19}) q_7^2 q_{19} + k_{32} (p_7^2 p_{19} - p_8) q_8 \\ & + k_{33} (p_{20} - p_5^2 p_{19}) q_5^2 q_{19} + k_{34} (p_5^2 p_{19} - p_{20}) q_{20} \end{aligned} \quad [\text{A7}]$$



where the pairs  $(p_i, q_i)$  with  $i = 16, \dots, 21$  are the generalised coordinates corresponding to  $F_1, B_{1,2-2}, B_{1,0-S}, F_2, B_{2,1-1}$  and  $B_{2,0-S}$  respectively. Finally, the Hamiltonian functions corresponding to the epigenetic regulation of the *LSGs* (see **Table S5**) are given by:

$$\begin{aligned} H_{E_1}(p, q) = & (p_{23} - p_{22})(k_{35}h_2p_{24}q_{22}q_{24} + k_{39}q_{22}) + (p_{22} - p_{23})(k_{36}h_1p_{22}q_{22}q_{23} + k_{41}h_1q_{23}) \\ & + (p_{23} - p_{24})(k_{37}h_4p_{22}q_{22}q_{24} + k_{40}q_{24}) + \\ & (p_{23} - p_{24})(k_{38}h_3p_{24}q_{22}q_{24} + k_{42}h_3q_{23}) \end{aligned} \quad [\text{A8}]$$

$$\begin{aligned} H_{E_2}(p, q) = & (p_{26} - p_{25})(k_{35}h_2p_{27}q_{25}q_{27} + k_{39}q_{25}) + (p_{25} - p_{26})(k_{36}h_1p_{25}q_{25}q_{26} + k_{41}h_1q_{26}) \\ & + (p_{26} - p_{27})(k_{37}h_4p_{25}q_{25}q_{27} + k_{40}q_{27}) + (p_{27} - p_{26})(k_{38}h_3p_{27}q_{25}q_{27} + k_{42}h_3q_{26}) \end{aligned} \quad [\text{A9}]$$

The pairs  $(p_i, q_i)$  with  $i = 22, \dots, 27$  appearing in [A8] and [A9] are the generalised coordinates corresponding to  $M_1, U_1, A_1, M_2, U_2$  and  $A_2$ , respectively.

## APPENDIX B. QSSA of the Hamilton equations for the gene regulation Hamiltonian

We begin QSSA analysis by studying the gene regulation Hamiltonian,  $H_{GR}(p, x)$ . By introducing re-scaled variables and the re-scaled gene regulation Hamiltonian in [24] and [25], we obtain:

$$s_i \frac{dp_i}{dt} = -y_0 \frac{\partial H_{GR}}{\partial x_i} \quad [\text{B1}]$$

$$s_i \frac{dx_i}{d\tau} = y_0 \frac{\partial H_{GR}}{\partial p_i} \quad [\text{B2}]$$

where  $i = 1, \dots, 21$  (i.e., all the variables except the ones characterising the epigenetic states of the *LSGs*) and

$$s_i = y_0 \text{ if } i = 1, 3, 5, 7 \text{ (} T_0 \text{ otherwise),} \quad [\text{B3}]$$

which separates the variables of our system into slow variables:

$$\frac{dp_i}{d\tau} = \frac{-\partial H_{GR}}{\partial x_i} \quad [\text{B4}]$$

$$\frac{dx_i}{d\tau} = \frac{\partial H_{GR}}{\partial p_i} \quad [\text{B5}]$$

$i = 1, 3, 5, 7$ , and fast variables:

$$\varepsilon \frac{dp_i}{d\tau} = \frac{-\partial H_{GR}}{\partial x_i} \quad [\text{B6}]$$

$$\varepsilon \frac{dx_i}{d\tau} = \frac{\partial H_{GR}}{\partial p_i} \quad [\text{B7}]$$

otherwise. Finally, the QSSA system is given by:

$$\frac{dp_i}{d\tau} = \frac{-\partial H_{GR}}{\partial x_i} \quad [\text{B8}]$$

$$\frac{dx_i}{d\tau} = \frac{\partial H_{GR}}{\partial p_i} \quad [\text{B9}]$$

$$0 = \frac{-\partial H_{GR}}{\partial x_i} \quad [\text{B10}]$$

$$0 = \frac{\partial H_{GR}}{\partial p_i} \quad [\text{B11}]$$

for  $i \neq 1, 3, 5, 7$ .

In addition, we also have several constraints that help us solve the QSSA system [B8] and [B9]:

(i) Our system is conservative, and therefore  $H_{GR}(p(t), x(t))$  must be constant along the solutions of [B8] and [B9], i.e.,  $H_{GR}(p(t), x(t)) = E_{GR}$ .  $E_{GR}$  is an arbitrary quantity, so for simplicity, we choose  $E_{GR} = 0$ .

(ii) We also need to consider the constraints regarding the number of binding sites in each promoter region. In re-scaled variables, [13] to [16] read:

$$x_2 + x_{10} + x_{11} + x_{12} = \vartheta_O \quad [\text{B12}]$$

$$x_4 + x_{13} + x_{14} + x_{15} = \vartheta_S \quad [\text{B13}]$$

$$x_6 + x_{16} + x_{17} + x_{18} = \vartheta_1 \quad [\text{B14}]$$

$$x_8 + x_{19} + x_{20} + x_{21} = \vartheta_2 \quad [\text{B15}]$$

where  $\vartheta_O = T_O / T_0$ ,  $\vartheta_S = T_S / T_0$ ,  $\vartheta_1 = T_I / T_0$ , and  $\vartheta_2 = T_2 / T_0$ . It is easy to verify from [B10] and [B11] with  $i = 2, 4, 6, 8$  that for the conservation laws [B12] to [B15] to hold,  $p_1 + p_3 + p_5 + p_7 = 1$  must be satisfied.

It is also easy to verify that **[B8]** to **[B11]** imply that  $H_{p_O}(p, x) = H_{p_S}(p, x) = H_{p_1}(p, x) = H_{p_2}(p, x) = 0$ . Therefore, under QSSA conditions, the gene regulation Hamiltonian reduces to:

$$H_{GR}(p, q) \approx H_{SD}(p, q) = k_1 p_2 (p_1 - 1) q_2 + R (p_1 - 1) + k_2 (1 - p_1) q_1 + k_3 p_4 (p_3 - 1) q_4 + R (p_3 - 1) + k_4 (1 - p_3) q_3 \quad \text{[B16]}$$

$$+ k_5 p_6 (p_5 - 1) q_6 + k_6 (1 - p_5) q_6 + k_7 p_8 (p_7 - 1) q_8 + k_8 (1 - p_7) q_7 \quad \text{[B17]}$$

which implies that the condition  $E_{GR} = 0$  is trivially satisfied when  $p_1 + p_3 + p_5 + p_7 = 1$ . Furthermore, **[B11]** for  $i = 2, 11, 12$ , read:

$$\kappa_1 p_2 (1 - p_1) + \kappa_{12} (p_1 - p_1 p_3 p_{10}) = 0 \quad \text{[B18]}$$

$$\kappa_{14} (p_{11} - p_5^2 p_{10}) = 0 \quad \text{[B19]}$$

$$\kappa_{16} (p_{12} - p_7^2 p_{10}) = 0 \quad \text{[B20]}$$

which, because  $p_1 + p_3 + p_5 + p_7 = 1$ , reduces to  $p_2 + p_{10} + p_{11} + p_{12} = c_0$  where  $c_0$  is a constant to be determined later in our calculation. Similarly, the fact that  $p_1 + p_3 + p_5 + p_7 = 1$  implies that:

$$p_4 = p_{13} = p_{14} = p_{15} = c_s \quad \text{[B21]}$$

$$p_6 = p_{16} = p_{17} = p_{18} = c_l \quad \text{[B22]}$$

$$p_8 = p_{19} = p_{20} = p_{21} = c_2 \quad \text{[B23]}$$

where  $c_s$ ,  $c_l$  and  $c_2$  are constants to be determined later.

To proceed further, consider **[B8]**, which reduces to:

$$\frac{dx_l}{d\tau} = \rho_l + \kappa_1 p_2 x_2 - \kappa_2 x_1 \quad \text{[B24]}$$

$$\frac{dx_3}{d\tau} = \rho_l + \kappa_3 p_4 x_4 - \kappa_4 x_3 \quad \text{[B25]}$$

$$\frac{dx_5}{d\tau} = \kappa_5 p_6 x_6 - \kappa_6 x_5 \quad \text{[B26]}$$

$$\frac{dx_7}{d\tau} = \kappa_7 p_8 x_8 - \kappa_8 x_7 \quad \text{[B27]}$$

Let us focus on **[B24]**, which determines the time evolution of  $x_l$ . Moreover, according to **[B10]** with  $i = 2, 11, 12$ , we have:

$$x_2 = \frac{1}{\kappa_{12}} x_1 x_3 x_{10} \quad [\text{B28}]$$

$$x_{11} = \frac{\kappa_{13}}{\kappa_{14}} x_5^2 x_{10} \quad [\text{B29}]$$

$$x_{12} = \frac{\kappa_{15}}{\kappa_{16}} x_7^2 x_{10} \quad [\text{B30}]$$

[B12] and [B28] enable us to write  $x_{10}$  as a function of  $x_1, x_3, x_5$  and  $x_7$  :

$$x_{10} = \frac{\vartheta_0}{1 + \frac{1}{\kappa_{12}} x_1 x_3 + \frac{\kappa_{13}}{\kappa_{14}} x_1 x_3 + \frac{\kappa_{15}}{\kappa_{16}} x_7^2} \quad [\text{B31}]$$

which, in turn, enables us to express  $x_2$  as a function of  $x_1, x_3, x_5$  and  $x_7$  :

$$x_2 = \frac{1}{\kappa_{12}} \frac{\vartheta_0 x_1 x_3}{1 + \frac{1}{\kappa_{12}} x_1 x_3 + \frac{\kappa_{13}}{\kappa_{14}} x_5^2 + \frac{\kappa_{15}}{\kappa_{16}} x_7^2} \quad [\text{B32}]$$

Finally, [B24] and [B32] lead to [29]. The derivation of equations [30] to [32] is completely analogous.

### APPENDIX C. Determination of the values and physical interpretation of $p_2, p_4, p_6$ and $p_8$

This appendix is devoted to determining the values of the constants  $p_2, p_4, p_6$  and  $p_8$ , which we need to close the QSSA of our system given by equations [29] to [32]. The procedure to calculate the value of these quantities also enables us to describe a physical interpretation of their meaning.

Using the fact that  $E_{GR} = 0$  and interpreting by parts, we can re-write the action functional  $S$  [22] as follows:

$$S_{GR}(p_1, \dots, p_{21}, \tau) = y_0 \sum_{i=1}^{21} \int_0^\tau x_i(s) \left( \frac{s_i}{y_0} \frac{dp_i}{ds} \right) ds + S(p_1, \dots, p_{21}, \tau = 0) \quad [\text{C1}]$$

where  $s_i$  is defined in [B3]. As the QSSA implies that  $p_i = \text{constant}$  for all  $i$ , the QSSA approximation of [C1],  $S_{GRQ}(p, t)$ , reduces to:

$$S_{GRQ}(p_1, \dots, p_{21}, \tau) = S_{GRQ}(p_1, \dots, p_{21}, \tau = 0) \quad [\text{C2}]$$

We further assume that the initial value of each of the state variables is an independent random variable, which implies that  $S_{GRQ}(p_1, \dots, p_{21}, \tau) = \sum_{i=1}^{21} S_i(p_i)$ . As  $p_1 = p_3 = p_5 = p_7 = 1$ , the corresponding functions  $S_1 = S_3 = S_5 = S_7 = 0$ . The resulting  $S_{GRQ}(p, \tau)$  can therefore be written as:

$$S_{GRQ}(p_{21}, \dots, p_{21}, \tau) = S_O(c_O) + S_S(c_S) + S_{L_1}(c_1) + S_{L_2}(c_2) \quad [\text{C3}]$$

where  $S_O(c_O)$  is the negative of the logarithm of the generating function of the probability density of the total number of binding sites within the promoter of *OCT4*.  $S_S(c_S)$ ,  $S_{L_1}(c_1)$  and  $S_{L_2}(c_2)$  are the corresponding quantities for the total number of binding sites in the promoters of *SOX2*, *LSG1* and *LSG2*, respectively. To obtain this result, we used [B21] and the well-known fact regarding the properties of the generating function of a random variable  $G_N(p)$ , which is the sum of  $n$  independent random variables given by:

$$G_N(p) = \prod_{i=1}^n G_i(p) \quad [\text{C4}]$$

where  $G_i(p)$  is the generating function of the probability density of each independent random variable (Grimmett and Stirzaker, 1992).

As [C3] implies that the total number of binding sites in each of the four promoters is an independent random variable, we can proceed to calculate, say,  $c_O$  by setting  $c_S = c_1 = c_2 = 1$ , i.e.,  $S_S(c_S = 1) = S_{L_1}(c_1 = 1) + S_{L_2}(c_2 = 1) = 0$ , which is equivalent to taking the marginal probability of the number of binding sites in the promoter of *OCT4*,  $b_O$ . The calculation of the other three quantities,  $c_S$ ,  $c_1$ , and  $c_2$ , is completely analogous.

As an immediate consequence of the definition of the probability generation function, one can use Cauchy's formula to obtain the probability of the number of binding sites in the promoter of *OCT4*,  $b_O$ , i.e., (2)

$$P(b_O, \tau) = \frac{1}{2\pi i} \oint \left( \frac{G_O(p, \tau)}{p^{b_O+1}} \right) dp = \frac{1}{2\pi i} \oint \left( \frac{e^{-s_O(p, \tau)}}{p^{b_O+1}} \right) dp = \frac{1}{2\pi i} \oint \left( \frac{e^{-s_O(p, \tau) - b_O \log(p)}}{p} \right) dp \quad [\text{C5}]$$

Let us now define the function  $f(p)$  through the relation  $E_{of}(p) = E_o \left( s_o(p) + \frac{b_o}{E_o} \log p \right)$ , where  $E_o$  is the average of the total number of binding sites in the promoter of *OCT4*.

$$P(b_o, \tau) = \frac{1}{2\pi i} \oint \left( \frac{e^{-E_o f(p)}}{p} \right) dp \quad [\text{C6}]$$

If  $E_o \gg 1$ , then we can apply the method of Laplace (or the method of the stationary phase) to approximately solve the integral (Ablowitz and Fokas, 2003; Murray, 1984). According to this method, if  $E_o \gg 1$ , the only contribution to the integral corresponds to the value of  $p = c_o$  for which  $f(p)$  exhibits a maximum. The value of  $c_o$  is therefore found by maximising  $f(p)$ , i.e.:

$$c_o \left| \frac{ds_o}{dp} \right|_{p=c_o} = \frac{-b_o}{E_o} \quad [\text{C7}]$$

Similarly,

$$c_s \left| \frac{ds_s}{dp} \right|_{p=c_s} = \frac{-b_s}{E_s} \quad [\text{C8}]$$

$$c_1 \left| \frac{ds_1}{dp} \right|_{p=c_1} = \frac{-b_1}{E_1} \quad [\text{C9}]$$

$$c_2 \left| \frac{ds_2}{dp} \right|_{p=c_2} = \frac{-b_2}{E_2} \quad [\text{C10}]$$

If we consider that the number of binding sites in each of the four promoters is distributed according to a Poisson distribution, we have the following:

$$S_o(p) = -E_o(p-1), s_o(p) = -(p-1) \quad [\text{C11}]$$

$$S_s(p) = -E_s(p-1), s_s(p) = -(p-1) \quad [\text{C12}]$$

$$S_1(p) = -E_1(p-1), s_1(p) = -(p-1) \quad [\text{C13}]$$

$$S_2(p) = -E_2(p-1), s_2(p) = -(p-1) \quad [\text{C14}]$$



which, according to equations [C7] to [C10], yields  $c_o = \frac{b_o}{E_o}, c_s = \frac{b_s}{E_s}, c_1 = \frac{b_1}{E_1}, c_2 = \frac{b_2}{E_2}$ .

*Physical interpretation.* The derivation of the previous section implies that the values of  $c_o, c_s, c_1$  and  $c_2$  are determined by the probability distribution of binding sites in the promoter of the corresponding gene.

The interpretation of this result is that these quantities are non-uniformly distributed over a population of cells. Thus, we look at a population of cells, the proportion of cells with a given number of binding sites in the promoter of, say, *OCT4* is given by [C6]. In other words, each cell randomly picked from such population will have a different value of binding sites in the *OCT4* promoter distributed according to [C6]. In the particular case that the number of binding sites is distributed according to a Poisson distribution, we find that, for example,  $c_o$  depends on the ratio between the actual binding sites in a cell within the population,  $b_o$ , and its average over the population, so that individual cells are characterised by a parameter that is determined by whether its number of binding sites in the *OCT4* promoter is above, exactly equal to, or below average.

#### Appendix C1. Determination of $C_{E_1}$ and $C_{E_2}$

To determine the values of the constants  $C_{E_1}$  and  $C_{E_2}$ , we follow the same general procedure as in the previous case. In the main text, we have established that for  $x_{22}, x_{23}$  and  $x_{24}$  and  $x_{25}, x_{26}$  and  $x_{27}$  to be positive, the corresponding momenta must satisfy:

$$p_{22} = p_{23} = p_{24} = \text{const.} = c_{E_1} \quad [\text{C15}]$$

$$p_{25} = p_{26} = p_{27} = \text{const.} = c_{E_2} \quad [\text{C16}]$$

These equations have the consequence that  $E_{ER} = H_{ER}(p, x) = 0$  (see equations [A6] and [A7]), which implies, in turn, that the corresponding action functional,  $S_{ER}(p, t)$ , can be expressed as:

$$S_{ER}(p_{22}, \dots, p_{27}, \tau) = y_0 \sum_{i=22}^{27} \int_0^\tau x_i(s) \left( \frac{s_i}{y_0} \frac{dp_i}{d\tau} \right) ds + S_{E_1}(p_{22}, \dots, p_{24}, \tau = 0) + S_{E_2}(p_{25}, \dots, p_{27}, \tau = 0) \quad [\text{C17}]$$

where  $s_i$  is given by [B3]. Furthermore because all the momenta in [C17] are constant, the epigenetic regulation action,  $S_{ER}(p_{22}, \dots, p_{27}, \tau)$ , reduces to:

$$S_{ER}(p_{22}, \dots, p_{27}, \tau) = S_{E_1}(p_{22}, \dots, p_{24}, \tau = 0) + S_{E_2}(p_{25}, \dots, p_{27}, \tau = 0) \quad [\text{C18}]$$

Following the same procedure as in the previous section, we finally obtain that  $C_{E_1}$  and  $C_{E_2}$  are given by:

$$c_{E_1} \left| \frac{ds_{E_1}}{dp} \right|_{p=c_{E_1}} = \frac{-b_{E_1}}{E_{E_1}} \quad [\text{C19}]$$

$$c_{E_2} \left| \frac{ds_{E_2}}{dp} \right|_{p=c_{E_2}} = \frac{-b_{E_2}}{E_{E_2}} \quad [\text{C20}]$$

where  $S_{E_1}(p) = E_{E_1} s_{E_1}(p)$  is the negative of the logarithm of the generating function of the probability density function of the number of modification sites of the *LSG1* gene (reciprocally, for *LSG2*).

The physical interpretation is also analogous to the one in the previous section, namely, each cell randomly picked from a cell population will have a different number of modification sites in *LSG1* and *LSG2*, distributed according to the corresponding probability distribution function.

Again, if the number of modification sites in *LSG1* and *LSG2* is distributed according to a Poisson distribution, then  $c_{E_1} = \frac{b_{E_1}}{E_{E_1}}$  and  $c_{E_2} = \frac{b_{E_2}}{E_{E_2}}$ .

#### APPENDIX D. Stability analysis of the epigenetic regulation model in the baseline scenario

To simplify our analysis, and without loss of generality, we assume that  $\kappa_{35} = \kappa_{36} = \kappa_{37} = \kappa_{38} = d_1$  and  $\kappa_{39} = \kappa_{40} = \kappa_{41} = \kappa_{42} = d_2$ . Similarly,  $\kappa_{43} = \kappa_{44} = \kappa_{45} = \kappa_{46} = d_1$  and  $\kappa_{47} = \kappa_{48} = \kappa_{49} = \kappa_{50} = d_2$ . The steady states of the system will be expressed as given by  $x^* = (1 - \alpha - \beta, \alpha, \beta)$ , where  $\alpha$  and  $\beta$  are the steady-state values of  $x_{23}$  (normalised number of unmodified sites) and  $x_{24}$  (normalised number of acetylated sites), respectively, which must satisfy:

$$-d_1 h_2 c_{E_1} (1 - \alpha - \beta) \beta - d_2 (1 - \beta - \alpha) + d_1 h_1 c_{E_1} (1 - \alpha - \beta) \alpha + d_2 h_1 \alpha = 0 \quad [\text{D1}]$$

$$-d_1 h_4 c_{E_1} (1 - \alpha - \beta) \beta - d_2 \beta + d_1 h_3 c_{E_1} \beta \alpha + d_2 h_3 \alpha = 0 \quad [\text{D2}]$$

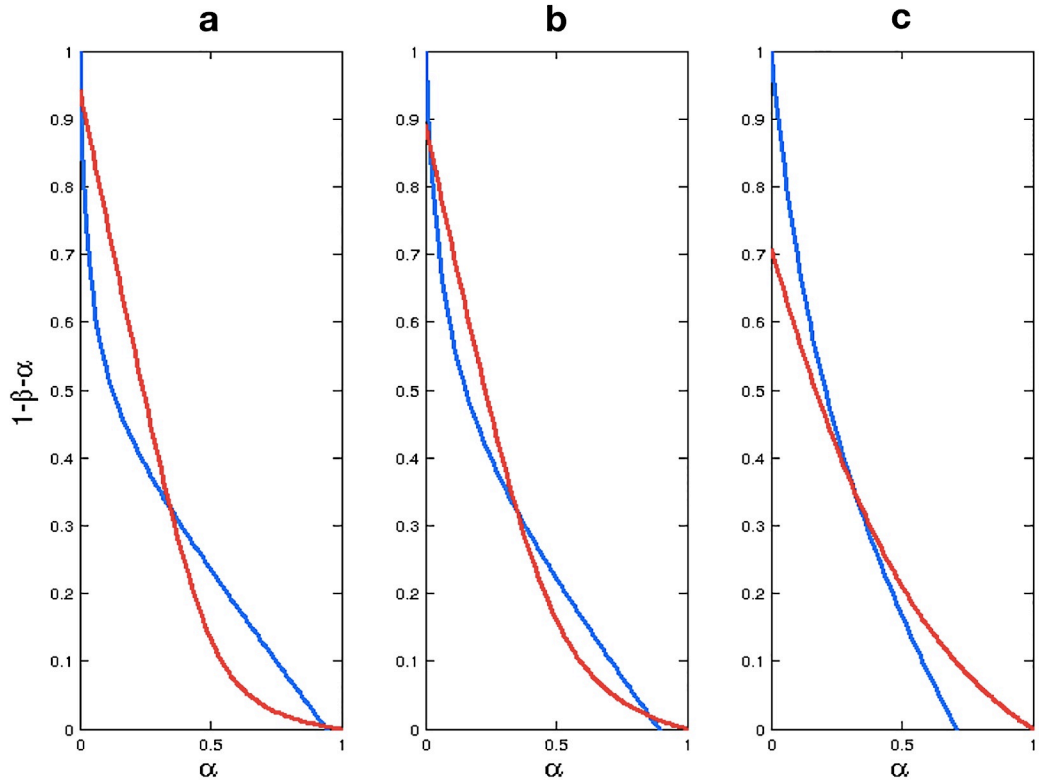
The number of steady states is determined by the intersections between the null-clines:

$$-d_1 h_2 c_{E_1} (1 - \alpha - \beta) \beta - d_2 (1 - \beta - \alpha) + d_1 h_1 c_{E_1} (1 - \alpha - \beta) \alpha + d_2 h_1 \alpha = 0 \quad [\text{D3}]$$

$$\alpha = \frac{d_1 h_4 c_{E_1} (1 - \beta) \beta + d_\beta}{d_1 c_{E_1} (h_3 + h_4) \beta + d_2 h_3}$$

[D4]

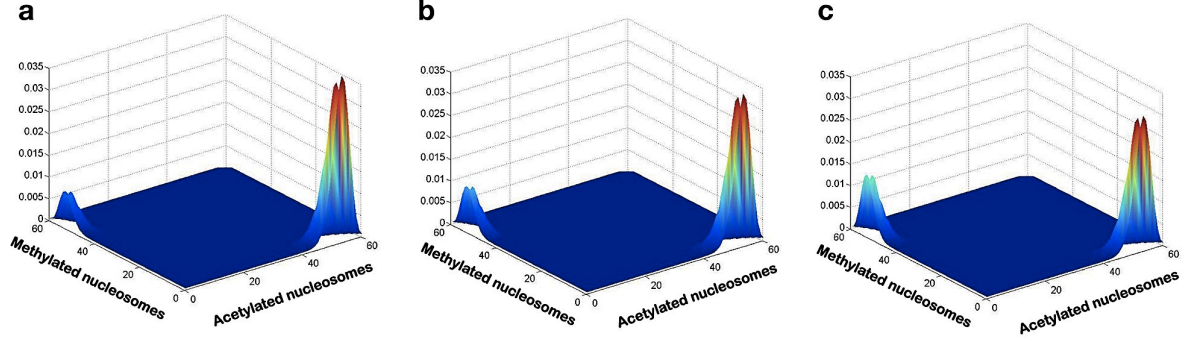
The result is illustrated in **Fig. S2**, in which we show the null-clines for different values of the parameters  $\delta_1 = d_1 c_{E_1}$ , with all other parameter values fixed. We can see that this parameter is a bifurcation control parameter, as modifying its value drives the system from bistability (**Fig. S2 (a) and (b)**) to monostability (**Fig. S2 (c)**) via a saddle-node bifurcation (Strogatz, 1994). In the bistable regime shown in **Fig. S2 (a) and (b)**, the two stable states correspond to the two states of epigenetic regulation: methylated (high steady-state value of  $x_{22}(x_{25})$  and low stationary value of  $x_{26}(x_{28})$ ) or acetylated (low steady-state value of  $x_{22}(x_{25})$  and high stationary value of  $x_{26}(x_{28})$ ).



**Figure S2. Plots showing the null-clines corresponding to the epigenetic regulation system (40) and (41) for different values of  $\delta_1$ .** Panel (a) corresponds to  $\delta_1 = 0.054$ , panel (b) to  $\delta_1 = 0.026$ , and panel (c) to  $\delta_1 = 0.0054$ . Other parameter values are taken from **Table E1**.

Our baseline scenario also requires that, in the absence of *OCT4* and *SOX2* induction and 2HG-induced reduction of HDM activity, our system must produce differentiated cells, which requires the epigenetic regulation model of both *LSG1* and *LSG2* to be acetylated so that transcription factors can access their promoters, whereupon differentiation ensues. This requirement is achieved by biasing (reducing) the activity parameter of the corresponding histone deacetylases. **Fig. S3** shows a stochastic simulation for the joint probability density

for the number of methylated and acetylated sites of our stochastic epigenetic regulation model. **Fig. S3** shows the corresponding result for the baseline scenario, which can be interpreted as implying that the time spent by the epigenetic regulatory system in the acetylated state is overwhelmingly longer than the time spent in the methylated state.



**Figure S3. Stochastic simulation results for the probability density of our stochastic model of epigenetic regulation.** Panel (a) shows a baseline scenario that corresponds to the parameter values given in **Table E1**. Panels (b) and (c) show two cases of 2HG-induced reduction of HDM activity with respect to the baseline scenario (b:  $h_2 = 0.975$  and c:  $h_2 = 0.97$ ).

The parameters corresponding to our baseline scenario are chosen so that (i) the epigenetic regulation system [40]-[41] and [43]-[44], corresponding to *LSG1* and *LSG2*, respectively, are in the bistable regime (see **Fig. S2**), and (ii) the acetylated state is the most stable of the two stable states (see **Fig. S3**). The parameter values compatible with these general properties are given in **Table S7**.

## APPENDIX E. Stability analysis of the gene regulatory model in the base-line scenario

Considering the following simplifying assumptions: (i) all the dimer-promoter binding constants are the same, (ii) all the dimer-promoter unbinding constants are the same, (iii) all the gene products are synthesised at the same rate, and (iv) all the proteins are degraded at the same rate, the system of equations [29] to [32] can be re-written as follows:

$$\frac{dx_1}{d\tau} = p_2 b_1 \frac{x_1 x_3}{1 + a(x_1 x_3 + x_5^2 + x_7^2)} - c x_1 \quad [\text{E1}]$$

$$\frac{dx_3}{d\tau} = p_4 b_1 \frac{x_1 x_3}{1 + a(x_1 x_3 + x_5^2 + x_7^2)} - c x_3 \quad [\text{E2}]$$

$$\frac{dx_5}{d\tau} = p_6 b_1 \frac{x_5^2}{1 + a(x_1 x_3 + x_5^2 + x_7^2)} - c x_5 \quad [\text{E3}]$$

$$\frac{dx_7}{d\tau} = p_8 b_1 \frac{x_7^2}{1 + a(x_1 x_3 + x_5^2 + x_7^2)} - c x_7 \quad [\text{E4}]$$

where

$$a = \frac{1}{\kappa_{12}} = \frac{\kappa_{13}}{\kappa_{14}} = \frac{\kappa_{15}}{\kappa_{16}} = \frac{\kappa_{17}}{\kappa_{18}} = \frac{\kappa_{19}}{\kappa_{20}} = \frac{\kappa_{21}}{\kappa_{22}} = \frac{\kappa_{23}}{\kappa_{24}} = \frac{\kappa_{25}}{\kappa_{26}} = \frac{\kappa_{27}}{\kappa_{28}} = \frac{\kappa_{29}}{\kappa_{30}} = \frac{\kappa_{31}}{\kappa_{32}} = \frac{\kappa_{33}}{\kappa_{34}} \quad [\text{E5}]$$

$$b_1 = \vartheta_O \frac{\kappa_1}{\kappa_1} = \vartheta_O \kappa_3 \frac{\kappa_{17}}{\kappa_{18}} \quad [\text{E6}]$$

$$b_2 = \vartheta_1 \kappa_5 \frac{\kappa_{25}}{\kappa_{26}} = \vartheta_2 \kappa_7 \frac{\kappa_{31}}{\kappa_{32}} \quad [\text{E7}]$$

$$c = \kappa_2 = \kappa_4 = \kappa_6 = \kappa_8 \quad [\text{E8}]$$

In this appendix, we analyse the conditions for the existence and stability of three fixed points of the form:

$$x_p = (\alpha, \alpha, 0, 0) \quad [\text{E9}]$$

$$x_1 = (0, 0, \beta, 0) \quad [\text{E10}]$$

$$x_2 = (0, 0, 0, \beta) \quad [\text{E11}]$$

where the first, second, third, and fourth components correspond to the steady-state value of *OCT4*, *SOX2*, *LSG1* and *LSG2*, respectively;  $x_p$  corresponds to the pluripotency steady state, and  $x_1$  and  $x_2$  correspond to each of the differentiated state states.

We also examine the existence and stability of a fixed point given by:

$$x_U = (\alpha, \alpha, \beta, 0) \quad [\text{E12}]$$

*Existence and linear stability of  $x_p$ ,  $x_1$  and  $x_2$ .* We begin our analysis by examining the properties of  $x_p$ . By using the steady-state version of equations [E1] to [E4] (i.e.,  $\frac{dx_1}{d\tau} = 0$ ), we have that, from [E1] and [E2],  $\alpha$  must satisfy:

$$c a \alpha^2 - p_2 b_1 \alpha + c = 0 \quad [\text{E13}]$$

$$c a \alpha^2 - p_4 b_1 \alpha + c = 0 \quad [\text{E14}]$$

which implies that  $p_2 = p_4$  ‡.

‡ If  $p_2 \neq p_4$ , we have a steady-state  $x_p = (\alpha_1, \alpha_2, 0, 0)$ , which has the same properties as  $x_p$

The equations [E3] and [E4] are trivially satisfied. If this restriction on  $p_2$  and  $p_4$  is satisfied, then  $\alpha$  is given by:

$$\alpha = \frac{p_2 b_1}{2ac} \left( 1 \pm \sqrt{1 - 4 \frac{ac^2}{p_2 b_1}} \right) \quad [\text{E15}]$$

According to this equation, for  $\alpha$  to be real, the following condition must be satisfied:

$$4 \frac{ac^2}{p_2 b_1} \leq 1 \quad [\text{E16}]$$

Provided that [16], the linear stability of  $x_p$  is determined by the sign of the eigenvalues of the corresponding Jacobian matrix,  $J_p$ , obtained by linearising the right hand side of equations [E1] to [E4] around  $x_p$ , which is given by:

$$J_p = \begin{pmatrix} A_p & 0 & 0 \\ 0 & -c & 0 \\ 0 & 0 & -c \end{pmatrix} \quad [\text{E17}]$$

where  $A_p$  is a 2 x 2 matrix corresponding to the linearisation of equations [E1] and [E2].

$$A_p = \begin{pmatrix} \frac{p_2 b_1 \alpha}{(1 + a\alpha^2)^2} - c & \frac{p_2 b_1 \alpha}{(1 + a\alpha^2)^2} \\ \frac{p_2 b_1 \alpha}{(1 + a\alpha^2)^2} & \frac{p_2 b_1 \alpha}{(1 + a\alpha^2)^2} - c \end{pmatrix} \quad [\text{E18}]$$

Given the block structure of  $J_p$ , its eigenvalues are the eigenvalues of each block, i.e.:

$$\lambda_1 = \lambda_2 = \lambda_3 = -c \quad [\text{E19}]$$

$$\lambda_4 = 2 \frac{p_2 b_1 \alpha}{(1 + a\alpha^2)^2} - c \quad [\text{E20}]$$

As for  $x_p$  to be stable, all the eigenvalues of the Jacobian  $J_p$  must be negative, we have a second condition that the system must satisfy:



$$c > 2 \frac{p_2 b_2 \alpha}{(1 - a \alpha^2)^2} \quad [\text{E21}]$$

An analogous analysis regarding the conditions for the existence of  $x_1$  and  $x_2$  leads to the following conditions:

$$4 \frac{ac^2}{p_6 b_2} \leq 1 \quad [\text{E22}]$$

$$4 \frac{ac^2}{p_8 b_2} \leq 1 \quad [\text{E23}]$$

respectively. Similarly, for them to be stable, the following conditions must hold:

$$c > 2 \frac{p_6 b_2 \beta}{(1 - a \beta^2)^2} \quad [\text{E24}]$$

$$c > 2 \frac{p_8 b_2 \beta}{(1 - a \beta^2)^2} \quad [\text{E25}]$$

The less restrictive pair of existence-stability conditions gives the region of coexistence of the three states as stable steady states. Parameter values compatible with such situations are given in **Table S7**.

**Table S7**

Parameter	Value	Description	References
$T_0$	60	Characteristic scale for number of binding sites	[Dodd et al., 2007]
$\gamma_0$	600	Characteristic scale for number of protein	
$h_1$	1	Activity of histone methyl transferases	
$h_2$	1	Activity of histone demethylases	
$h_3$	1	Activity of histone acetyl transferases	
$h_4$	0.9	Activity of histone deacetylases	
$d_1$	0.027	Rate of recruited nucleosome modifications	
$d_2$	0.003	Rate of unrecruited nucleosome modifications	
$d_3$	0.425	See <b>Appendix C</b>	
$c_{E_1}$	0.425	See <b>Appendix C</b>	
$c_{E_1}$	$9 \cdot 10^5$	Re-scaled dimer-promoter affinity (See <b>Appendix E</b> )	
$b_1$	5	Re-scaled rate of OCT4 and SOX2 protein production (See <b>Appendix E</b> )	

$b_2$	10	Re-scaled rate of LSG1 and LSG2 protein production (See <b>Appendix E</b> )
$c$	$1.1 \cdot 10^{-6}$	Re-scaled rate of protein degradation (See <b>Appendix E</b> )
$c_o$	1	See <b>Appendix C</b>
$c_s$	1	See <b>Appendix C</b>
$c_2$	1	See <b>Appendix C</b>
$c_2$	1	See <b>Appendix C</b>

---

**Parameter values corresponding to the baseline scenario**

**APPENDIX F. Competition and invasion for clones originated from CSCs *de novo* generated by nuclear reprogramming**

We have analysed a mathematical model of competition between two cell populations, i.e., a native, “resident” population sustained by normal stem cells (SCs) and an “invader” population of cancer stem cells (CSCs) *de novo* generated by oncometabolic nuclear reprogramming. We focused on the so-called *neutral scenario*. We assumed that reprogrammed CSCs display the same features as normal SCs, i.e., they possess the same self-renewal and differentiation rates and, crucially, are under the same homeostatic controls than their healthy counterparts. By choosing this setting that corresponds to a base-line case of *least favourable invasion*, we are faced with a situation in which the CSCs are initially under zero-net growth conditions, i.e., identical to that of the native (healthy) cell population. This implies that the early evolution of the invader can be studied in terms of a simple diffusion equation, since its dynamic is essentially dictated by fluctuations due to smallness of the cell population.

**Introduction.** We herein addressed how significant is the increase of reprogramming rate observed in our proof-of-concept studies with live cells in terms of cancer evolution. In our hands, the number of reprogrammed CSC-like colonies increased by >10 times in the presence of the oncometabolite 2HG. In this appendix we propose a simple mathematical model of cancer evolution in which we demonstrate that the probability of spontaneous clearance of an invader generated by a colony of *de novo* reprogrammed CSCs decreases exponentially with the size of the colony.

**Model formulation**

*Model hypotheses*

Our modelling approach is based on a wealth of papers that have approached the problem of regulation of differentiation cascades (Marciniak-Czochra et al., 2009; Pepper et al., 2007; Rodriguez-Brenes et al., 2010; Sánchez-Taltavull and Alarcón, 2014).

(i) We assume that a hierarchical differential cascade maintains each population:

$$\begin{aligned} \text{Resident: } & \text{SC} \rightarrow \text{TAC1} \rightarrow \text{TAC2} \cdots \rightarrow \text{MC} \\ \text{Invader: } & \text{CSC} \rightarrow \text{TAC1} \rightarrow \text{TAC2} \cdots \rightarrow \text{MC} \end{aligned}$$

SC: stem cell, CSC: cancer stem cell, TAC: transient amplifying cell, MC: mature (fully differentiated) cell.

We assume that the length of the differentiation cascade (i.e. the number of stages between SC or CSC and MC),  $L$ , is the same for both populations. SCs and CSCs are the only cell types with the ability for self-renewal whereas TACs are assumed to differentiate upon proliferation.

(ii) Both populations are assumed to compete for a common pool of resources. As far as the model is concerned, this is implemented through a carrying capacity that limits the size of the total population (i.e., resident *plus* invader).

(iii) We consider the most disfavourable situation from the point of view of the invader, i.e. we will assume a neutral scenario in which all the parameter values, e.g. birth rates, death rates, differentiation rates, etc. are set to be the same for both populations. We examine the ability for an invader to take over the resident population as a function of the number of CSCs (equivalently, of the rate of reprogramming of somatic MCs).

(iv) Population dynamics (Pepper et al., 2007, Sánchez-Taltavull and Alarcón, 2014):

(a) SCs and CSCs undergo:

- Asymmetric division at a rate:

$$(1 - \epsilon)p_0 e^{-M/K},$$

where  $M = M_r + M_i$  with  $M_r$  and  $M_i$  are the number of mature cells of the resident and invader, respectively.  $K$  is the carrying capacity and the probability of SC (CSC) symmetric division. Note that if  $M \gg K$  (i.e. the mature population is much larger than the carrying capacity) the division rate tends to zero, whereas when  $M \ll K$  the proliferation rate approaches its maximum value.

- Symmetric self-renewal at a rate:

$$(1 - d_0)\epsilon p_0 e^{-M/K}$$

where  $d_0$  is the probability of symmetric SC (CSC) differentiation.

- Symmetric differentiation at a rate:

$$d_0 \epsilon p_0 e^{-M/K}$$

- Apoptosis at a rate  $l_0$ . We will assume that both SCs and CSCs undergo cell death at a very slow rate.

(b) For simplicity and without loss of generality we assume that all TACs differentiate and die at the same rate, regardless of their position in the cascade:

(c) Symmetric differentiation at a constant rate  $d_1$

(d) Apoptosis at a rate  $\lambda_1$

(e) Mature cells undergo death at a constant rate  $\lambda_3$

**Table S8**

Event (Resident population)	Transition rate	$r_k$
Asymmetric SC division	$W_{r_1} = (1 - \epsilon)p_0 e^{-M/K} X_{r_1}$	$r_{r_1} = (0, +1, 0, \dots, 0)$
Symmetric SC self-renewal	$W_{r_2} = \epsilon(1 - d_0)p_0 e^{-M/K} X_{r_1}$	$r_{r_2} = (+1, 0, 0, \dots, 0)$
Symmetric SC differentiation	$W_{r_3} = \epsilon d_0 p_0 e^{-M/K} X_{r_1}$	$r_{r_3} = (-1, +2, 0, \dots, 0)$
SC apoptosis	$W_{r_4} = \lambda_0 X_{r_1}$	$r_{r_4} = (-1, 0, 0, \dots, 0)$
Symmetric $j$ th-TAC differentiation	$W_{r_{4+2(j-2)+1}} = d_1 X_{r_j}$	$r_{r_{4+2(j-2)+1}} = (0, 0, \dots, -1, +2, \dots, 0)$
TAC apoptosis	$W_{r_{4+2(j-2)+2}} = \lambda_1 X_{r_j}$	$r_{r_{4+2(j-2)+2}} = (0, 0, \dots, -1, 0, \dots, 0)$
MC apoptosis	$W_{r_{4+2(L-3)+3}} = \lambda_3 X_{r_L}$	$r_{r_{4+2(L-3)+3}} = (0, 0, 0, \dots, -1)$
Event (Invader population)	Transition rate	$r_k$
Asymmetric CSC division	$W_{i_1} = (1 - \epsilon)p_0 e^{-M/K} X_{i_1}$	$r_{i_1} = (0, +1, 0, \dots, 0)$
Symmetric CSC self-renewal	$W_{i_2} = \epsilon(1 - d_0)p_0 e^{-M/K} X_{i_1}$	$r_{i_2} = (+1, 0, 0, \dots, 0)$
Symmetric CSC differentiation	$W_{i_3} = \epsilon d_0 p_0 e^{-M/K} X_{i_1}$	$r_{i_3} = (-1, +2, 0, \dots, 0)$
SC apoptosis	$W_{i_4} = \lambda_0 X_{i_1}$	$r_{i_4} = (-1, 0, 0, \dots, 0)$
Symmetric $j$ th-TAC differentiation	$W_{i_{4+2(j-2)+1}} = d_1 X_{i_j}$	$r_{i_{4+2(j-2)+1}} = (0, 0, \dots, -1, +2, \dots, 0)$
TAC apoptosis	$W_{i_{4+2(j-2)+2}} = \lambda_1 X_{i_j}$	$r_{i_{4+2(j-2)+2}} = (0, 0, \dots, -1, 0, \dots, 0)$
MC apoptosis	$W_{i_{4+2(L-3)+3}} = \lambda_3 X_{i_L}$	$r_{i_{4+2(L-3)+3}} = (0, 0, 0, \dots, -1)$

Transition rates associated to the stochastic dynamics of the competition between the resident and the invader hierarchical populations.  $X_{r_l}$  and  $X_{i_l}$  are the number of SC and CSCs, respectively,  $X_{r_j}$  and  $X_{i_j}$ ,  $j=2, \dots, L-1$ , are the number of TACs of the resident and invader, respectively and  $X_{r_L}$  and  $X_{i_L}$  are the number of the mature cells in either population. The total mature population,  $M$ , is given by  $M=X_{r_L}+X_{i_L}$ .

### Stochastic dynamics

The stochastic population dynamics of the system described in the above mentioned section is described in terms of the Master Equation [1] (Gardiner, 2009):

$$\frac{\partial P(X, t)}{\partial t} = \sum_{i=1}^R (W_i(X - r_i)P(X - r_i, t) - W_i(X)P(X, t)) \quad \text{[F1]}$$

The associated rates,  $W_i$ , and stoichiometric vectors,  $r_i$ , are given in **Table S8**.  $X$  is the  $2L$ -dimensional state vector whose entries are the numbers of each cellular type in the population.

### Setup and initial conditions

We consider a setup or initial condition in which the resident population is let to evolve until it reaches a steady state. This steady state is described (on average) by the mean-field limit of the stochastic dynamics, as shown below. Once the resident steady state has been reached, we evaluate the behaviour of the invader population as a function of the number of *de novo* generated CSCs.

### Model analysis

We start by discussing the behaviour of mean-field limit, in particular regarding its steady states. We then proceed to study the diffusion limit of the Master Equation [1], which provides closed form solutions for the clearance probability of the invader.

### Mean-field limit

- *Mean-field limit in the absence of invader.* We start by the describing the mean-field limit (Gillespie, 1976) associated to the resident population in the absence of the invader. The Master Equation **[F1]** is then given by the following system of ordinary differential equations:

$$\frac{dx_{r_1}}{dt} = (\epsilon(1 - d_0)p_0e^{-x_{r_L}} - \epsilon d_0 p_0 e^{-x_{r_L}} - \lambda_0) x_{r_1} \quad \text{[F2]}$$

$$\frac{dx_{r_2}}{dt} = 2\epsilon d_0 p_0 e^{-x_{r_L}} x_{r_1} - (d_1 + \lambda_1) x_{r_2} \quad \text{[F3]}$$

$$\frac{dx_{r_j}}{dt} = 2d_1 x_{r_{j-1}} - (d_1 - \lambda_1) x_{r_j}, \quad j = 3, \dots, L - 1 \quad \text{[F4]}$$

$$\frac{dx_{r_L}}{dt} = 2d_1 x_{r_{L-1}} - \lambda_2 x_{r_L} \quad \text{[F5]}$$

where  $x_{ri} = X_{ri}/K$ . Eqs. **[F2]** to **[F5]** have two steady states: the trivial equilibrium  $(x_{r_1}, x_{r_2}, \dots, x_{r_j}, \dots, x_{r_L}) = (0, 0, \dots, 0, \dots, 0)$  and the stable positive equilibrium given by:

$$x_{r_1} = \frac{d_1 + \lambda_1}{2\epsilon d_0 p_0 e^{-x_{r_L}}} x_{r_2}, \quad \text{[F6]}$$

$$x_{r_j} = \frac{d_1 + \lambda_1}{2d_1} x_{r_{(j+1)}}, \quad j = 3, \dots, L - 1, \quad \text{[F7]}$$

$$x_{r_L} = -\log \left( \frac{\lambda_0}{\epsilon(1 - 2d_0)p_0} \right). \quad \text{[F8]}$$

The positive equilibrium exists provided that:

$$\frac{\lambda_0}{\epsilon(1 - 2d_0)p_0} < 1 \quad \text{and} \quad d_0 < \frac{1}{2}$$

Eq. **[F8]** is obtained by imposing that the net growth rate of the SC population is equal to zero:

$$\epsilon(1 - 2d_0)p_0e^{-x_{r_L}} - \lambda_0 = 0.$$

- *Mean-field limit in the presence of the invader.* The presence of the invader substantially changes the structure of the steady states of the model. First, the mean-field equations for the whole system (resident *plus* invader) are:

$$\frac{dx_{r_1}}{dt} = (\epsilon(1 - d_0)p_0e^{-(x_{r_L}+x_{i_L})} - \epsilon d_0 p_0 e^{-(x_{r_L}+x_{i_L})} - \lambda_0) x_{r_1} \quad [\text{F9}]$$

$$\frac{dx_{r_2}}{dt} = 2\epsilon d_0 p_0 e^{-x_{r_L}} x_{r_1} - (d_1 + \lambda_1) x_{r_2} \quad [\text{F10}]$$

$$\frac{dx_{r_j}}{dt} = 2d_1 x_{r_{j-1}} - (d_1 - \lambda_1) x_{r_j}, \quad j = 3, \dots, L-1 \quad [\text{F11}]$$

$$\frac{dx_{r_L}}{dt} = 2d_1 x_{r_{L-1}} - \lambda_2 x_{r_L} \quad [\text{F12}]$$

$$\frac{dx_{i_1}}{dt} = (\epsilon(1 - d_0)p_0e^{-(x_{r_L}+x_{i_L})} - \epsilon d_0 p_0 e^{-(x_{r_L}+x_{i_L})} - \lambda_0) x_{i_1} \quad [\text{F13}]$$

$$\frac{dx_{i_2}}{dt} = 2\epsilon d_0 p_0 e^{-(x_{r_L}+x_{i_L})} x_{i_1} - (d_1 + \lambda_1) x_{i_2} \quad [\text{F14}]$$

$$\frac{dx_{i_j}}{dt} = 2d_1 x_{i_{j-1}} - (d_1 - \lambda_1) x_{i_j}, \quad j = 3, \dots, L-1 \quad [\text{F15}]$$

$$\frac{dx_{i_L}}{dt} = 2d_1 x_{i_{L-1}} - \lambda_2 x_{i_L} \quad [\text{F16}]$$

where  $x_{r_l} = X_{r_l}/K$  and  $x_{i_l} = X_{i_l}/K$ .

As in the previous case, Eqs. [F9] to [F16], have two steady states: the trivial equilibrium  $(x_{r_1}, x_{r_2}, \dots, x_{r_j}, \dots, x_{r_L}, x_{i_1}, x_{i_2}, \dots, x_{i_j}, \dots, x_{i_L}) = (0, 0, \dots, 0, \dots, 0, 0, 0, \dots, 0, \dots, 0)$ . The associated positive equilibrium is given by:

$$x_{r_1} = \frac{d_1 + \lambda_1}{2\epsilon d_0 p_0 e^{-(x_{r_L}+x_{i_L})}} x_{r_2}, \quad [\text{F17}]$$

$$x_{r_j} = \frac{d_1 + \lambda_1}{2d_1} x_{r_{(j+1)}}, \quad j = 3, \dots, L-1, \quad [\text{F18}]$$

$$x_{i_1} = \frac{d_1 + \lambda_1}{2\epsilon d_0 p_0 e^{-(x_{r_L}+x_{i_L})}} x_{i_2}, \quad [\text{F19}]$$

$$x_{i_j} = \frac{d_1 + \lambda_1}{2d_1} x_{i_{(j+1)}}, \quad j = 3, \dots, L-1, \quad [\text{F20}]$$

$$x_{r_L} + x_{i_L} = -\log \left( \frac{\lambda_0}{\epsilon(1 - 2d_0)p_0} \right) \equiv \Omega. \quad [\text{F21}]$$

Note that in the presence of the invader, the system does not exhibit a unique positive equilibrium but rather a continuum of equilibria determined by Eq. [F21]. Eq. [F21] is a restatement of the requirement that the growth rate of the SCs and CSCs is equal to zero. The positive equilibria exist provided that:

$$\frac{\lambda_0}{\epsilon(1 - 2d_0)p_0} < 1 \text{ and } d_0 < \frac{1}{2}$$

## Neutral dynamics: Fokker-Planck equation with demographical noise

We now formulate a diffusion approximation of the stochastic process under the following conditions: Once the resident population has settled down onto its steady state, we introduce a number of *de novo* reprogrammed CSCs and study the behaviour of the invader population generated by the CSCs. The diffusion approximation allows us to find closed form analytical expressions for how the probability of clearance of the invader varies as the number of CSCs changes.

Before proceeding with the formal derivation of the diffusion approximation, an important caveat, which greatly simplifies the analysis, is in order. The scenario we are contemplating, in which a number of CSCs appear in the system (in steady state) as a result of oncometabolic reprogramming of somatic mature cells, implies that both the SC (i.e. the native stem cells of the healthy, resident tissue) and the CSC populations are under conditions of zero net growth rate. Initially,  $X_{iL}$ , the number of fully differentiated cells associated to the invader, is  $X_{iL}=0$  and  $X_{iL}$  is similar to  $N$  where  $N \equiv \Omega K$ , which satisfies the equilibrium condition whereby the growth rate of the CSC population is zero, i.e., the dynamics of the CSC population is critical. Since this condition is initially fulfilled, the whole system continues to evolve under conditions so that Eqs. [17] to [21] are satisfied at any later time.

Under these conditions, the dynamics of the CSC population is critical with an average total mature population given by Eq. [21]. Therefore, its dynamics is entirely dominated by fluctuations. Furthermore, if we assume that  $K$  is big enough so that  $X_{rL} + X_{iL}$  is similar to  $N$ (=cnt.), we can assume that the dynamics of the CSCs is decoupled and can be studied independently.

*Diffusion approximation: system size expansion.* The systems size expansion is a well-established technique, originally proposed by Kampen et al. (2008), to extract the mean-field limit, i.e., the deterministic behaviour associated to infinite systems where no fluctuations are present, and its first stochastic correction given by the Fokker-Planck equation for the Gaussian fluctuations around the mean-field behaviour (Gardiner, 2006; Kampen et al., 2008). This approximation has been widely used and successfully applied to many different types of problems.

The original approach of Kampen et al. (2008) has been recently critiqued by Di Patti et al (2011), who noted that it breaks down for systems with absorbing states (e.g. extinctions) in the proximity of the absorbing state. They further propose a modification of the basic size expansion that exhibits improved accuracy for systems with absorbing states. The system size expansion is based on the following assumption regarding the stochastic process:

$$X(t) = Nx(t) + N^\alpha \xi(t),$$

where  $N$  is a measure of system size,  $X(t)$  is the mean-field limit,  $\xi(t)$  is a stochastic correction to the mean-field. In the original proposal by Kampen et al. (2008)  $\alpha = 1/2$ , which implies that the size of the effects of noise relative to the size of the systematic (mean-field part) decreases as  $N^{-1/2}$  as system size grows. This Ansatz produces a consistent expansion which, at the lowest order, yields a linear-noise (i.e. independent of  $\xi(t)$ ) Fokker-Planck equation (FPE) for the probability density function (PDF) of  $\xi$ .

Di Patti et al. (2011) have shown that, in the case in which the stochastic process has an absorbing state, the Van Kampen's Ansatz is no longer consistent. On the contrary, for the system size expansion to be consistent they show that  $\alpha = 0$ , i.e. the size of the effects of noise is independent of system size. In this case, the Fokker-Planck equation one obtains for the PDF of  $\xi(t)$  is no longer a linear-noise equation. On the contrary, one obtains an FPE with a particular type of non-linear (i.e. dependent on the random variable  $\xi(t)$ ) noise which we refer to as demographic noise. We proceed now to formulate the associated FPE for our system. For the technical details regarding the expansion, we refer the reader to Di Patti et al. (2011).

Assuming that  $X_{rL} + X_{iL} \cong N$ , we can define a stochastic dynamics for the evolution of the number of CSCs defined by the transition rates:

$$\begin{aligned} W_1(X_{i1}) &= e^{-\Omega} \epsilon p_0 (1 - d_0) X_{i1}, \quad r_1 = +1, \\ W_2(X_{i1}) &= e^{-\Omega} \epsilon p_0 d_0 X_{i1}, \quad r_2 = -1, \\ W_3(X_{i1}) &= \lambda_0 X_{i1}, \quad r_3 = -1. \end{aligned}$$

We further define  $T_{\pm}(z_{iL})$  where  $z_{iL} = X_{iL}/N$  as:

$$T_+(z_{i1}) = e^{-\Omega} \epsilon p_0 (1 - d_0) z_{i1}, \quad T_-(z_{i1}) = (e^{-\Omega} \epsilon p_0 d_0 + \lambda_0) z_{i1}. \quad [\text{F22}]$$

Finally, we expand the quantities  $T_{\pm}(z_{iL}) = T_{\pm}(\phi + N^{\alpha-1} \xi)$  as a power series of  $y \equiv N^{\alpha-1} \xi(t)$ :

$$T_{\pm}(\phi + y) = \sum_{k=0} T_{\pm}^{(k)}(\phi) \frac{y^k}{k!} \quad [\text{F23}]$$

With the above definitions, Di Patti et al. (2011) have shown that the Fokker-Planck equation for the PDF of  $\xi(t)$ ,  $P(\xi, t)$ , is given by:

$$\frac{\partial P}{\partial t} = \left( T_+^{(1)}(0) - T_-^{(1)}(0) \right) \frac{\partial}{\partial \xi} (\xi P) + \frac{1}{2} \left( T_+^{(1)}(0) + T_-^{(1)}(0) \right) \frac{\partial^2}{\partial \xi^2} (\xi P). \quad [\text{F24}]$$

In our case:

$$T_+^{(1)}(0) + T_-^{(1)}(0) = e^{-\Omega} \epsilon p_0 + \lambda_0,$$

see Eq. [F21], and

$$T_+^{(1)}(0) + T_-^{(1)}(0) = e^{-\Omega} \epsilon p_0 + \lambda_0,$$

and, therefore, Eq. [F24] reads:



$$\frac{\partial P}{\partial t} = \frac{D}{2} \frac{\partial^2}{\partial \xi^2} (\xi P). \quad [\text{F25}]$$

where  $D \equiv e^{-\lambda_0} p_0 + \lambda_0$ . Eq. [F25] is the diffusion approximation for the process of neutral dynamics of the CSC population.

*Survival probability of the CSC population: backward Kolmogorov equation.* In order to study the survival probability of the CSC population, rather than directly using the FPE Eq. [F25], it is more convenient to use an equivalent description, the so-called Kolmogorov backward equation (KBE) (Demetrius et al., 2009; Gardiner, 2009). Mathematically speaking, the KBE is the adjoint equation of the FPE (also known as the Kolmogorov forward equation). The KBE associated to the Fokker-Planck equation Eq. [25] is given by (Demetrius et al., 2009 and Gardiner, 2009):

$$\frac{\partial \Psi}{\partial t} = \frac{D}{2} q \frac{\partial^2 \Psi}{\partial q^2}. \quad [\text{F26}]$$

with natural boundary conditions  $\Psi(q=0, t)=0$  and  $\Psi(q=\infty, t)=1$ . The associated initial condition is  $\Psi(q, t=0)=1$  for positive  $q$ . With these boundary conditions, the solution of the KBE can be interpreted as the probability of survival of a CSC population of initial size  $q$  at time  $t$ . The associated clearance probability at time  $t$ ,  $P_C(q, t)$ , is therefore given by  $P_C(q, t)=1-\Psi(q, t)$ .

It is possible to obtain a closed form, analytical solution for Eq. [F26] by considering its similarity structure (Ockendon et al., 2003). It is immediate to verify that, if  $\Psi(q, t)$  is a solution of Eq. [F26], so is  $\Psi(\mu q, \mu t)$  where  $\mu$  is an arbitrary constant. Therefore, for any given value of  $t$ , we can set  $\mu=t^{-1}$ , so that  $\Psi(q, t)=F(q/t)$  for some function  $F$  to be determined from Eq. [F26]. By writing  $\eta=q/t$ , Eq. [F26] can be re-written as:

$$\frac{D}{2} \frac{d^2 F}{d\eta^2} + \frac{dF}{d\eta} = 0 \quad [\text{F27}]$$

with boundary conditions  $F(\eta=0)=0$  and  $F(\eta=\infty)=1$ . The solution is therefore given by:

$$\Psi(q, t) = F(q/t) = 1 - P_C(q, t) = 1 - (P_1(t))^q = 1 - e^{-\frac{2}{D} \frac{q}{t}} \quad [\text{F28}]$$

where  $P_1(t)$  is the probability of clearance of a single CSC at time  $t$ . The clearance probability is therefore given by  $P_C(q, t)=(P_1(t))^q$ . This result implies that the probability of spontaneous clearance of a population generated by a colony of reprogrammed CSCs decreases exponentially as the size of the colony,  $q$ , increases. This implies that with a 10 to 20-fold increase in nuclear reprogramming frequency, the chances of clearance of the population generated by the oncometabolically *de novo* reprogrammed CSCs decreases by many orders of magnitude.

In summary, we are able to find an analytical solution for spontaneous clearance probability of the CSC population as a function of the size of the *de novo* reprogrammed CSC pool,  $q$ , at time  $t$ ,  $P_C(q, t)=(P_1(t))^q$ , where  $P_1(t)<1$  is the probability of clearance at time  $t$  of a single

CSC. Therefore, from the point of view of cancer evolution, an apparently inefficient 10 to 20-fold increase in the reprogramming frequency is highly significant in terms of the prolonged survival of the initial population generated by the *de novo* reprogrammed CSCs.

## Supplemental references

---

- Ablowitz, M. J., Fokas, A. S. *Complex variables. Introduction and applications*. Cambridge University Press, Cambridge, UK. 2003.
- Alarcón, T. Stochastic quasi-steady state approximations for asymptotic solutions of the chemical master equation. *J. Chem. Phys.* 2014;140:184109.
- Assaf, M., Meerson, B., Sasorov, P. V. Large fluctuations in stochastic population dynamics: momentum space calculations. *J. Stat. Mech.* 2010; P07018.
- Cinquin, O., Demongeot, J. High-dimensional switches and the modelling of cellular differentiation. *J. Theor. Biol.* 2005; 233:391-411.
- Cinquin, O., Page, K. M. Generalized, switch-like competitive heterodimerization networks. *Bull Math. Biol.* 2007; 69:483-494.
- Demetrius, L., Gundlach, V. M., Ochs, G. Invasion exponents in biological networks. *Physica A.* 2009; 388: 651-672.
- Dickman, R., Vidigal, R. *Braz. J. Phys. Path Integrals and Perturbation Theory for Stochastic Processes.* 2003; 33:73-93.
- Dykman, M. I., Horita, T., Ross, J. Statistical distribution and stochastic resonance in a periodically driven chemical system. *J. Chem. Phys.* 1995; 103:966-972.
- Elgart, V., Kamenev, A. Rare event statistics in reaction-diffusion systems. *Phys. Rev. E. Stat. Nonlin. Soft. Matter. Phys.* 2004; 70:041106.
- Feynman, R. P., Hibbs, A. R. *Quantum Mechanics and Path Integrals*, McGraw Hill: 1965 (ISBN 0-07-020650-3), Dover Publications: 2010.
- Gardiner, C. W. *Stochastic methods*. Springer-Verlag, Berlin, Germany. 2009
- Gillespie, D. T. A General Method for Numerically Simulating the Stochastic Time Evolution of Coupled Chemical Reactions. *J. Comp. Phys.* 1976; 22:403-434.
- Grimmett, G. R., Stirzaker, D. R. *Probability and random processes*. 2<sup>nd</sup> ed Oxford University Press. 1992.
- Kampen, N. G. V. *Stochastic processes in Physics and Chemistry*. Elsevier, The Netherlands. 2007.
- Keener, J., Sneyd, J. *Mathematical physiology*. Springer-Verlag, New York, NY, USA. 1998.
- Kubo, R., Matsuo, K., Kitahara, K. Fluctuation and Relaxation of Macrovariables. *J. Stat. Phys.* 1973; 9:51-96.

- MacArthur, B. D., Lemischka, I. R. Statistical mechanics of pluripotency. *Cell*. 2013; 154: 484-489.
- MacArthur, B. D., Please, C. P., Oreffo, R. O. Stochasticity and the molecular mechanisms of induced pluripotency. *PLoS One*. 2008; 3:e3086.
- Marciniak-Czochra, A., Stiehl, T., Ho, A. D., Jäger, W., Wagner, W. Modeling of asymmetric cell division in hematopoietic stem cells--regulation of self-renewal is essential for efficient repopulation. *Stem Cells Dev*. 2009;18:377-385.
- Murray, J. D. *Asymptotic analysis*. Springer-Verlag, New York, NY, USA. 1984.
- Di Patti, F., Azaele, S., Banavar, J. R., Maritan, A. System size expansion for systems with an absorbing state. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys*. 2011; 83:010102.
- Pepper, J. W., Sprouffske, K., Maley, C. C. Animal cell differentiation patterns suppress somatic evolution. *PLoS Comput. Biol*. 2007; 3:e250.
- Rodriguez-Brenes, I. A., Wodarz, D., Komarova, N. L. Minimizing the risk of cancer: tissue architecture and cellular replication limits. *J. R. Soc. Interface*. 2013; 10:20130410.
- Sánchez-Taltavull, D., Alarcón, T. Robustness of differentiation cascades with symmetric stem cell division. *J. R. Soc. Interface*. 2014; 11:20140264.
- Sakurai, J. J. *Modern quantum mechanics* (Addison-Wesley, New York, USA) 1994.
- Strogatz, S. H. *Non-linear dynamics and chaos*. Perseus Books, Reading, Mass. USA. 1994.
- Takahashi, K., Tanabe, K., Ohnuki, M., Narita, M., Ichisaka, T., Tomoda, K., Yamanaka, S. Induction of pluripotent stem cells from adult human fibroblasts by defined factors. *Cell*. 2007; 131:861-872.
- Takahashi, K., Yamanaka, S. Induction of pluripotent stem cells from mouse embryonic and adult fibroblast cultures by defined factors. *Cell*. 2006; 126:663-676.
- Täuber, U. C., Howard, M., Vollmayr-Lee, B. P. Applications of field-theoretic renormalization group methods to reaction-diffusion problems. *J. Phys. A: Math. Gen*. 2005; 38:R79-R131.
- Traulsen, A., Lenaerts, T., Pacheco, J. M., Dingli, D. On the dynamics of neutral mutations in a mathematical model for a homogeneous stem cell population. *J. R. Soc. Interface*. 2012; 10:20120810.